

# Modeling Ordered Choices: A Primer and Recent Developments

William H. Greene<sup>1</sup>  
David A. Hensher<sup>2</sup>

*June 15, 2008*

Version 4

## Abstract

We survey the literature on models for ordered choices, including ordered logit and probit specifications. The contemporary form of the model is presented and analyzed in detail. The historical development of the model is presented as well. We detail a number of generalizations that have appeared in the recent literature. Finally, we propose a new form of the model that accommodates in a natural, internally consistent form, functional form flexibility and individual heterogeneity. Much of this study is pedagogical. However, the last few sections propose new model formulations, and illustrate them with an application to self reported health satisfaction.

<sup>1</sup>Department of Economics, Stern School of Business, New York University, New York, NY 10012, [wgreene@stern.nyu.edu](mailto:wgreene@stern.nyu.edu)

<sup>2</sup>Institute of Transport and Logistics Studies, Faculty of Economics and Business, University of Sydney, NSW 2006 Australia [Davidh@itls.usyd.edu.au](mailto:Davidh@itls.usyd.edu.au)

# Contents

- 1 Introduction 4
2. An Ordered Choice Model for Social Science Applications 6
  - 2.1 A Latent Regression Model for a Continuous Measure 6
  - 2.2 The Observed Discrete Outcome 8
  - 2.3 Probabilities and the Log Likelihood 9
  - 2.4 Analysis of Data on Ordered Choices 10
3. Antecedents and Contemporary Counterparts 11
  - 3.1 The Origin of Probit Analysis: Bliss (1934), Finney (1947) 11
  - 3.2 Social Science Data and Regression Analysis for Binary Outcomes 15
  - 3.3 Analysis of Binary Choice 16
  - 3.4 Ordered Outcomes: Aitchison and Silvey (1957), Snell (1964) 16
  - 3.5 Minimum Chi Squared Estimation of an Ordered Response Model: Gurland et al. (1960) 20
  - 3.6 Individual Data and Polychotomous Outcomes: Walker and Duncan (1967) 21
  - 3.7 McElvey and Zavoina (1975) 22
  - 3.8 Developments Since Zavoina and McElvey 23
  - 3.9 Other Related Models 25
    - 3.9.1 Known Thresholds 26
    - 3.9.2 Nonparallel Regressions 27
4. Estimation, Inference and Analysis Using the Ordered Choice Model 28
  - 4.1 Application of the Ordered Choice Model to Self Assessed Health Status 28
  - 4.2 Distributional Assumptions 28
  - 4.3 The Estimated Ordered Probit (Logit) Model 30
  - 4.4 Interpretation of the Model – Partial Effects and Scaled Coefficients 32
    - 4.4.1 Nonlinearities in the Variables 36
    - 4.4.2 Average Partial Effects 36
    - 4.4.3 Interpreting the Threshold Parameters 38
    - 4.4.4 The Underlying Regression 38
  - 4.5 Inference 39
    - 4.5.1 Inference about Coefficients 39
    - 4.5.2 Testing for Structural Change or Homogeneity of Strata 43
    - 4.5.3 Robust Covariance Matrix Estimation 43
    - 4.5.4 Inference About Partial Effects 44
  - 4.6 Prediction – Computing Probabilities 45
  - 4.7 Measuring Fit 48
  - 4.8 Estimation Issues 52
    - 4.8.1 Grouped Data 53
    - 4.8.2 Perfect Prediction 53
    - 4.8.3 Different Normalizations 54
    - 4.8.4 Censoring of the Dependent Variable 54
    - 4.8.5 Maximum Likelihood Estimation of the Ordered Choice Model 56
    - 4.8.6 Bayesian (MCMC) Estimation of Ordered Choice Models 57
    - 4.8.7 Software For Estimation of Ordered Choice Models 61
5. Specification Issues and Generalized Models 63
  - 5.1 Functional Form Issues and the Generalized Ordered Choice Model (1) 63
    - 5.1.1 Parallel Regressions 64
    - 5.1.2 Testing the Parallel Regressions Assumption – The Brant (1990) Test 65
    - 5.1.3 Generalized Ordered Logit Model (1) 70
    - 5.1.4 The Single Crossing Feature of the Ordered Choice Model 74
    - 5.1.5 Choice Invariant Ratios of Partial Effects 77
    - 5.1.6 Methodological Issues 77
  - 5.2 Accommodating Heterogeneity 78
    - 5.2.1 Threshold Models – The Generalized Ordered Probit Model (2) 79

5.2.2	Nonlinear Specifications – A Hierarchical Ordered Probit Model	81
5.2.3	Heterogeneous Scaling (Heteroscedasticity) of Random Utility	86
5.2.4	Individually Heterogeneous Marginal Utilities	89
5.2.5	Random Parameters Models	90
	Implied Heteroscedasticity	90
	Maximum Simulated Likelihood Estimation	91
	Conditional Mean Estimation in the Random Parameters Model	92
5.2.6	Latent Class and Finite Mixture Modeling	97
	The Latent Class Ordered Choice Model	97
	Estimation by Maximum Likelihood	98
	The EM Algorithm	99
	Estimating the Class Assignments	101
	A Latent Class Model Extension	101
	Application	102
	Endogenous Class Assignment and A Generalized Ordered Choice Model	106
5.2.7	Generalized Ordered Choice Model (3)	109
5.3	Specification Tests for Ordered Choice Models	114
5.3.1	Model Specifications – Missing Variables and Heteroscedasticity	114
5.3.2	Testing Against the Logistic and Normal Distribution	117
5.3.3	Unspecified Alternatives	119
6.	Ordered Choice Modeling with Panel Data	122
6.1	Ordered Choice Models with Fixed Effects	122
6.2	Ordered Choice Models with Random Effects	126
6.3	Testing for Random or Fixed Effects	129
6.4	Extending Parameter Heterogeneity Models to Ordered Choices	133
7	Extensions	137
7.1	Dynamic Models	137
7.2	Inflation Models	140
7.3	Multiple Equations	143
7.3.1	Bivariate Ordered Probit Models	143
7.3.2	Polychoric Correlation	145
7.3.3	Semi-Ordered Bivariate Probit Model	146
7.3.4	Applications of the Bivariate Ordered Probit Model	146
7.3.5	A Panel Data Version of the Bivariate Ordered Probit Model	147
7.3.6	Trivariate Ordered Probit Model	149
7.3.7	Models of Sample Selection with an Ordered Probit Selection Rule	149
7.3.8	A Sample Selected Ordered Probit Model	154
7.3.9	An Ordered Probit Model with Endogenous Treatment Effects	157
8	Semiparametric Estimators and Analyses	158
8.1	Heteroscedasticity	159
8.2	A Distribution Free Estimator with Unknown Heteroscedasticity	160
8.3	A Semi-nonparametric Approach	161
8.4	A Partially Linear Model	164
8.5	Semiparametric Analysis	164
9	Conclusions	166
	References	167
	Appendix	

# 1 Introduction

The model of ordered choice pioneered by Aitchison and Silvey (1957) and Snell (1964) and articulated in its modern form by McElvey and Zavoina (1975) has become a widely used tool in many fields. The number of applications in the current literature is large and increasing rapidly. A quick search of just the “ordered probit” model identified applications on:

- academic grades [Butler et al. (1994), Li and Tobias (2006a)],
- bond ratings [Terza (1985)],
- Congressional voting on a Medicare bill [McElvey and Zavoina (1975)],
- credit ratings [Cheung (1996)],
- driver injury severity in car accidents [Eluru et al. (2008)],
- drug reactions [Fu et al.(2004)],
- duration [Ridder (1990)],
- education [Machin and Vignoles (2005), Carneiro et al. (2001, 2003), Cameron and Heckman (1998)],
- eye disease severity [Biswas and Das (2002)],
- happiness [Winkelmann (2005), Zigarette (2007)],
- health status [Greene (2008a) based on Riphahn et. al (2003)],
- job classification in the military [Marcus and Greene (1983)],
- labor supply [Heckman and MaCurdy (1981)],
- life satisfaction [Clark et al. (2001)],
- monetary policy [Eichengreen, Watson and Grossman (1985)],
- nursing labor supply [Brewer et al. (2008)],
- obesity [Greene, Harris, Hollingsworth and Maitra (2008)],
- perceptions of difficulty making left turns [Zhang (2007)],
- pet ownership [Butler and Chatterjee(1997)],
- product quality [Bresnahan (1987), Prescott and Visscher (1977), Shaked and Sutton (1982)],
- promotion and rank in nursing [Pudney and Shields (2000)],
- stock price movements [Tsay (2005)],
- tobacco use [Harris and Zhao (2007), Kasteridis, Munkin and Yen (2008)],

and hundreds more.

Social science oriented introductions to the ordered choice model appear in journal articles such as Winship and Mare (1984), Becker and Kennedy (1992), Daykin and Moffatt (2002) and Boes and Winkelmann (2006a), and in textbook treatments including Maddala (1983), DeMaris (2004), Long (1997), Long and Freese (2006) and Greene (2008a). The practitioner who desires a quick entry level primer on the model can choose among numerous sources for a satisfactory introduction to the ordered choice model and its uses. There are also scores of surveys and primers for bioassay, including, e.g., Greenland (1994), Agresti (1999) and Ananth and Kleinbaum (1997). This survey is offered as an addition to this list for a number of purposes.

- A number of interesting extensions of the model already appearing in the literature are not mentioned in the surveys listed above.
- Recent analyses of the ordered choice model have uncovered some interesting avenues of generalization.

- The model formulation rests on a number of subtle underlying aspects that are not developed as completely as are the mechanics of using the “technique.” Only a few of the surveys devote substantial space to interpreting the model’s components once they are estimated. As made clear here and elsewhere, the coefficients in an ordered choice provide, in isolation, almost no useful information about the phenomenon under study. Yet, estimation of coefficients and tests of statistical significance are the central (sometimes, only) issue in many of the surveys listed above, and in some of the received applications.
- We will offer our own generalizations of the ordered choice model.
- With the creative development of easy to use contemporary software, many model features and devices are served up because they *can* be computed without much (or any) discussion of *why* they would be computed, or, in some cases, even *how* they are computed. To cite an example, Long and Freese (2006, pp. 195-196) state “several different measures [of fit] can be computed...” [using Stata] for the ordered probit model. Their table that follows lists 20 values, seven of which are statistics whose name contains “R squared.” The values range from 0.047 to 0.432. No discussion of what the measures are, what they mean, or how they are computed follows; the section provides the reader with a single statement that two Monte Carlo studies have found that one of the measures “closely approximates the  $R^2$  obtained by fitting the linear regression model on the underlying latent variable.” Obviously researchers differ on what information they wish to extract from the data. We will attempt to draw the focus to a manageable few aspects of the model that appear to have attained some degree of consensus.

The review proceeds as follows. The fundamental ordered choice model is developed in some detail in Section 2. The historical antecedents to the basic model are documented in Section 3. In section 4, we return to the modern form of the model, and develop the different aspects of its use, such as interpreting the model, statistical inference and fit measures. Some recent generalizations and extensions are presented in Sections 5 - 7. Semiparametric models that reach beyond the mainstream of research are discussed in Section 8. An application based on a recent study [Riphahn, Wambach and Million (2003)] will be dispersed through the discussion to provide an illustration of the points being presented.

There is an equally large literature parallel to the social science applications in the areas of biometrics and psychometrics. The distinction is not perfectly clean, but there is a tangible difference in orientation, as will be evident below. From the beginning with Bliss’s (1934a) invention of probit modeling, many of the methodological and statistical developments in the area have taken place in this setting. It will be equally evident that these two areas of application have developed in parallel, but by no means in concert. Our survey to follow is largely directed toward social science applications. However, the extensions and related features of the models and techniques in biometrics will be integrated into the presentation.

## 2. An Ordered Choice Model for Social Science Applications

The ordered probit model in its modern form was proposed by McElvey and Zavoina (1975) for the analysis of ordered, categorical, nonquantitative choices, outcomes and responses. [But, see the discussion of Gurland et al. (1960) in Section 3.5.] Familiar examples include bond ratings, discrete opinion surveys such as those on political questions, obesity measures, preferences in consumption, and satisfaction and health status surveys such as those analyzed by Boes and Winkelmann (2006a, 2006b) and other applications mentioned in the introduction. The model is used to describe the data generating process for a random outcome that takes one of a set of discrete, *ordered* outcomes. The health satisfaction or opinion survey provide clear examples.

### 2.1 A Latent Regression Model for a Continuous Measure

The model platform is an underlying random utility model or latent regression model,

$$y_i^* = \boldsymbol{\beta}'\mathbf{x}_i + \varepsilon_i, i = 1, \dots, N,$$

in which the continuous latent utility or ‘measure,’  $y_i^*$  is observed in discrete form through a censoring mechanism;

$$\begin{aligned} y_i &= 0 \text{ if } \mu_{-1} < y_i^* \leq \mu_0, \\ &= 1 \text{ if } \mu_0 < y_i^* \leq \mu_1, \\ &= 2 \text{ if } \mu_1 < y_i^* \leq \mu_2 \\ &= \dots \\ &= J \text{ if } \mu_{J-1} < y_i^* \leq \mu_J. \end{aligned}$$

The vector  $\mathbf{x}_i$  is a set of  $K$  covariates that are assumed to be strictly independent of  $\varepsilon_i$ ;  $\boldsymbol{\beta}$  is a vector of  $K$  parameters that is the object of estimation and inference. The  $N$  sample observations are labeled  $i = 1, \dots, N$ . Long and Freese (2006, p. 183) caution that one ought to insure that the model to be considered here really is appropriate for the variable of interest before embarking on the analysis. In their case, the question is whether the measured outcome really is ordered. They cite an application of ordering of occupations. Indeed, it is easy to see the validity of their conclusion; the ranking based on, say, some prestige scale is likely to be completely different from a ranking of the same set of outcomes based on expected income. The interpretation of the ordered outcome as a censoring of an underlying continuously measured preference or other measure will provide a reliable guide as to the appropriateness of the model. The thrust of the model is that the observed outcome is not simply a set of discrete outcomes that by some criterion *can* be ordered; the observed outcome is a monotonic (many to one) transformation of a single continuous outcome that naturally *must* be ordered. The further example that Long and Freese pursue, in which the response variable is one of “Strongly Disagree,” “Disagree,” “Agree,” and “Strongly Agree” is a clear example of a censoring of an underlying preference scale.

The use of models for ordered outcomes arises in many literatures, as suggested in the introduction. The literatures do have focal points in two centers, social sciences including sociology, political science, economics and psychology and in bioassay, as discussed at length below. A reading of the literature in both places suggests that social scientists are broadly comfortable with the idea of the censoring mechanism as the data generating process behind their samples of, usually, individual observations. Their counterparts in bioassay occasionally express some ambivalence about the underlying regression. In Aitchison and Silvey’s (1957) canonical application (developed at length below), there is no clear regression based data generating process at work; if anything the only stimulus in the model is the passage of time, and there are no

“coefficients” or “responses” in the equation. Nonetheless, as we explore below, there is a clear, if not perfect correspondence between their analysis and the ordered choice model. Snell (1964) in contrast, begins development of his model with “We assume there to be an underlying continuous scale of measurement along which the scale categories represent intervals.” Once again, however, the analysis to follow has nothing to do with regression; the model relates to discovery of the threshold values in the presence of an individual “effect.” But, the applications in the study clearly apply to continuous preference scales, in one case a taste test and in another an opinion survey with answers *terrible, poor, fair, good, excellent*.

The use of the latent regression to represent an underlying preference, or utility scale, and the translation of the utility into a discrete indicator has critics in many quarters. A lengthy discussion of the relevance (or irrelevance) of economics to the formulations appears in Hammermesch (2004). On the question, for example, of “how happy does your income make you?” – the question analyzed at some length by Boes and Winkelmann (2006b – see, esp., pp. 4-5) and illustrated below – Hammermesch asks whether it is meaningful to equate this “happiness” with utility. We will then associate the measured outcomes with the supposed utility. For better or worse, this is the position reached by many of the social science applications where the models of ordered choice are applied. They rest crucially on the notion of the underlying regression and the censoring process that produces the measured outcome. Ferrer-i-Carbonell and Frijters (2004) take the discussion yet another level deeper, and consider the underlying assumptions that must be at work in order to use satisfaction measures to reflect underlying welfare measures. [See, as well, Winkelmann and Winkelmann (1998).]

McCullagh (1980) is widely regarded as a codiscoverer of the ordered choice model. (Curiously, he makes no mention of McElvey and Zavoina (1975).) He states (on page 109)

Motivation for the proposed models is provided by appeal to the existence of an underlying continuous and perhaps unobservable random variable. In bioassay this latent variable usually corresponds to a “tolerance” which is assumed to have a continuous distribution in the population. Tolerances, themselves, are not directly observable but increasing tolerance as manifest through an increase in the probability of survival. The categories are envisaged as contiguous intervals on the continuous scale.... Ordinality is therefore an integral feature of such models and the imposition of an arbitrary scoring system for the categories is thereby avoided.

At least to some extent, Anderson and Philips (1981, p. 22) seem unpersuaded;

It is often possible to argue that an ordered categorical variable is a coarsely measured version of a continuous variable not itself observable. Thus, it is reasonable to assume that the ordered categories correspond to non-overlapping and exhaustive intervals of the real line. ... Although the existence of a latent continuous variable is not crucial for our arguments, it makes interpretation easier and clearer.

They do suggest that in at least one application, a method of predicting the values of the unobservable variable will be developed. Nonetheless, the development of their model begins (on p. 23) with

Suppose that individuals are grouped into  $k$  ordered groups which are identified by an ordered categorical variable  $y$  with arbitrarily assigned value  $s$  for the  $s$ th ordered group;  $s = 1, \dots, k$ . The variable  $y$  is a convenient identifier for some of the arguments presented later. *The ordering of the groups is not, in general, based on any numerical measurement.*

(Emphasis added.) Anderson (1984, p. 1) in something of a tour de force on ordered outcomes, seems to move in both directions at once:

Particular emphasis is placed on the case where  $y$  is an ordered categorical variable and the category with  $y = y_i$  is taken to be “lower” than the category with  $y = y_j$  if  $i < j$ . ... In principle, there is a single unobservable, continuous variable related to this ordered scale, but in practice, the doctor making the assessment will use several pieces of information in making his judgment on the observed category.

The notions of the latent continuous variable and the existence of the latent regression are not mere semantics. At least this is the point behind some of the preceding discussion. Superficially, the same model will arise in any case. However, the underlying platform turns out to be a crucial element of making sense of parameters that are estimated and of interpretations of the empirical model once obtained from the data. Consider, for example, also from Anderson (1984, p.2).

The dimensionality of the regression relationship between  $y$  and  $\mathbf{x}$  is determined by the number of linear functions required to describe the relationship. If only one linear function is required, the relationship is one dimensional; otherwise it is multidimensional. For example, in predicting  $k$  categories of pain relief from predictors  $\mathbf{x}$ , suppose that different functions  $\beta_1'\mathbf{x}$  and  $\beta_2'\mathbf{x}$  are required to distinguish between the pairs of categories (*worse, same*) and (*same, better*), respectively. *Then the relationship is neither one-dimensional nor ordered with respect to  $\mathbf{x}$ .*

(Emphasis added.) Essentially, the observation is about curve fitting and functional form. One might ask in this instance, “what are the coefficients?” For the current purpose, however, the question would seem to be “what if the simple regression model seems to be inadequate in terms of predicting (by an as yet unspecified procedure) the outcome?” However, the observation raises a vexing question. What if the outcomes, themselves, *are* manifestly ordered. Precisely what does the last sentence imply about the model that is generalized in such a way as to purposely be adequate to handle the full dimensionality of the outcome, as if it were not ordered at all? We will return to this issue below in the context of one of the “generalized” ordered choice models.

## 2.2 The Observed Discrete Outcome

The model contains the unknown marginal utilities,  $\beta$ , as well as  $J+2$  unknown threshold parameters,  $\mu_j$ , all to be estimated using a sample of  $n$  observations, indexed by  $i = 1, \dots, N$ . The data consist of the covariates,  $\mathbf{x}_i$  and the observed discrete outcome,  $y_i = 0, 1, \dots, J$ . The assumption of the properties of the “disturbance,”  $\varepsilon_i$ , completes the model specification. The conventional assumptions are that  $\varepsilon_i$  is a continuous random disturbance with conventional cdf,  $F(\varepsilon_i|\mathbf{x}_i) = F(\varepsilon_i)$  with support equal to the real line, and that the density,  $f(\varepsilon_i) = F'(\varepsilon_i)$  is likewise defined over the real line. The assumption of the distribution of  $\varepsilon_i$  includes independence from (or exogeneity of)  $\mathbf{x}_i$ .

By the laws of probability, the probabilities associated with the observed outcomes are

$$\text{Prob}[y_i = j | \mathbf{x}_i] = \text{Prob}[\varepsilon_i \leq \mu_j - \beta'\mathbf{x}_i] - \text{Prob}[\mu_{j-1} - \beta'\mathbf{x}_i], j = 0, 1, \dots, J.$$

Several normalizations are needed to identify the model parameters. First, in order to preserve the positive signs of all of the probabilities, we require  $\mu_j > \mu_{j-1}$ . Second, if the support is to be the entire real line, then  $\mu_{-1} = -\infty$  and  $\mu_J = +\infty$ . Since the data contain no unconditional information on scaling of the underlying variable – if  $y_i^*$  is scaled by any positive value, then

scaling the unknown  $\mu_j$  and  $\beta$  by the same value preserves the observed outcomes – an unconditional, free variance parameter,  $\text{Var}[\varepsilon_i] = \sigma_\varepsilon^2$ , is not identified (estimable). It is convenient to make the identifying restriction  $\sigma_\varepsilon = \bar{\sigma}$ , a constant. The usual approach to this normalization is to assume that  $\text{Var}[\varepsilon_i|\mathbf{x}_i] = 1$  in the probit case and  $\pi^2/3$  in the logit model – in either case to eliminate the free structural scaling parameter. Finally, assuming (as we will) that  $\mathbf{x}_i$  contains a constant term, we will require  $\mu_0 = 0$ . (If, with the other normalizations, and with a constant term present, this normalization is not imposed, then adding a constant to  $\mu_0$  and the same constant to the intercept term in  $\beta$  will leave the probability unchanged.)

We note at this point a minor ambiguity in the received literature. Some treatments omit the overall constant term in  $\beta$  and, in turn, omit the now unnecessary normalization  $\mu_0 = 0$ . The counterpart in these treatments is  $\beta_0 = 0$ , where  $\beta_0$  is the overall constant term. In related fashion, some treatments (e.g., the *Stata* and *SAS* software packages) translate the outcome variable to  $y_i = 1, 2, \dots, J$ , which produces a different count of possible outcomes. We have maintained the formulation above for two reasons. First, most empirical applications in our experience are based on data that actually contain zero as the origin – e.g., the GSOEP data analyzed by Boes and Winkelmann (2006a, 2006b). Second, as we have formulated the model, the familiar binary choice (probit and logit) models are useful parametric special cases that do not require a reformulation of the entire model. This feature is noted elsewhere by some of the authors discussed below.

The standard treatment in the received literature completes the ordered choice model by assuming either a standard normal distribution for  $\varepsilon_i$ , producing the “ordered probit” model or a standardized logistic distribution (mean zero, variance  $\pi^2/3$ , which produces the “ordered logit” model). Applications appear to be well divided between the two. A compelling case for one distribution or the other remains to be put forth – historically, a preference for the logistic distribution has been based on mathematical convenience and because of its ready revelation of “odds ratios” in a convenient closed form. [But, see Berkson (1951) who “prefers logits to probits” in a direct response to Finney. Unfortunately, Berkson’s arguments will not help to resolve the issue in the setting of this review.] Contemporary software such as *Stata* and *NLOGIT* have automated menus of other distributional choices, for example, the asymmetric Gompertz and extreme value distributions. However the motivation for these distributions is even less persuasive than that for a preference for probits over logits. These two overwhelmingly dominate the received applications; the others seem more than anything else to be gadgets that are straightforward to program in the software. [An exception is Han and Hausman (1986), who present a model in which an ordered extreme value model emerges naturally. A similar example of duration modeling by Formisano et al. (2001) is described by Simonoff (2003, pp. 435-448.)]

### 2.3 Probabilities and the Log Likelihood

With the full set of normalizations in place, the likelihood function for estimation of the model parameters is based on the implied probabilities,

$$\text{Prob}[y_i = j | \mathbf{x}_i] = [F(\mu_j - \beta' \mathbf{x}_i) - F(\mu_{j-1} - \beta' \mathbf{x}_i)] > 0, j = 0, 1, \dots, J.$$

Estimation of the parameters is a straightforward problem in maximum likelihood estimation. [See, e.g., Pratt (1981) and Greene (2007a, 2008a).] The log likelihood function is

$$\log L = \sum_{i=1}^n \sum_{j=0}^J m_{ij} \log [F(\mu_j - \beta' \mathbf{x}_i) - F(\mu_{j-1} - \beta' \mathbf{x}_i)]$$

where  $m_{ij} = 1$  if  $y_i = j$  and 0 otherwise. Maximization is done subject to the constraints  $\mu_{.1} = -\infty$ ,  $\mu_0 = 0$  and  $\mu_{.J} = +\infty$ . The remaining constraints,  $\mu_{j-1} < \mu_j$  can, in principle, be imposed by a reparameterization in terms of some underlying structural parameters, such as

$$\mu_j = \sum_{m=1}^j \exp(\alpha_m),$$

however, this is typically unnecessary. (It is necessary in the generalization suggested in Section 5.2.7 below.) Expressions for the derivatives of the log likelihood can be found in McElvey and Zavoina (1975), Maddala (1983), Long (1997), Stata (2008) and Econometric Software (2007).

The most recent literature (since 2005) includes several applications that use Bayesian methods to analyze ordered choices. Being heavily parametric in nature, they have focused exclusively on the ordered probit model. Some commentary on methods and methodology may be found in Koop and Tobias (2006). Applications to the univariate ordered probit model include Kadam and Lenk (2008), Ando (2006), Zhang et al. (2007) and Tomoyuki and Akira (2006). In the most basic cases, with diffuse priors, the “Bayesian” methods merely reproduce (with some sampling variability) the maximum likelihood estimator. [See Train (2003) for discussion of the Bernstein – von Mises result.] The MCMC methodology is often useful in settings which extend beyond the basic model. We will describe below, for example, applications to a bivariate ordered probit model [Biswas and Das (2002)], a model with autocorrelation [Czado et al. (2005) and Girard and Parent (2001)] and a model that contains a set of endogenous dummy variables in the latent regression [Munkin and Trivedi (2008).]

## 2.4 Analysis of Data on Ordered Choices

Analysis of ordered outcomes appears at many points in the literature since the (apparent) emergence with Aitchison and Silvey (1957). As discussed below, what sets McElvey and Zavoina apart is their adaptation to social science applications – the analysis of individual data. The central focus of the applications in bioassay was and is on grouped data and the analysis of proportions. The analysis of individual data, in a regression-like setting was relatively new at this point in the literature. Cox (1970), Finney (1971), Theil (1969, 1970, 1971) among others make mention of analysis of individual binary data, but McElvey and Zavoina (1975) were the first to extend the ideas of the ordered choice analysis to a model that was closely akin to regression modeling in cross sections of social science data. We will pursue this dichotomy in the next section, on the antecedents to the ordered probit (and logit) models.

### 3. Antecedents and Contemporary Counterparts

McElvey and Zavoina's proposal is preceded by several earlier developments in the statistical literature. The chronology to follow does suggest, however, that their development produced a discrete jump in the received body of techniques. The obvious starting point was the early work on probit methods in toxicology, beginning with Bliss (1934a) and made famous by Finney's (1947b) classic monograph on the subject. The ordered choice model that we are interested in here appears in three clearly discernible steps in the literature, Aitchison and Silvey's (1957) treatment of stages in the life cycle of a certain insect, Snell's (1964) analysis of ordered outcomes (without a regression interpretation) and McElvey and Zavoina's (1975) proposal of the modern form of the "ordered probit regression model." Some later papers, e.g., Anderson (1984) expanded on the basic models.

#### 3.1 The Origin of Probit Analysis: Bliss (1934), Finney (1947)

Bliss (1934a) tabulated graphically the results of a laboratory study of the effectiveness of an insecticide. He plotted the relationship between the "Percent of Aphids Killed" on the ordinate and "Milligrams of Nicotine Per 100 ML of Spray" on the abscissa of a simple figure, reproduced below as Figure 1. The figure loosely traces out the familiar sigmoid shape of the normal cdf, and in a natural fashion provides data on what kill rate can be expected for a given concentration of nicotine.

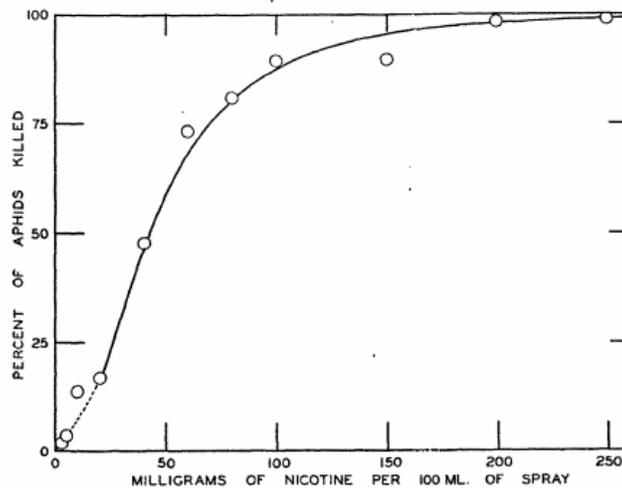


FIG. 1. Net mortality of *Aphis rumicis* L. sprayed in laboratory with different solutions of nicotine; summary of results over 3-year period. Tattersfield and Gimingham.<sup>4</sup> Heavy curve is same as that in Fig. 2 transposed back to original units.

Figure 1 Insecticide Experiment

The inverse question – “what concentration is necessary to achieve a given kill rate?” – is answered by inverting the function in the figure. Writing

$$p_i = F(c_i)$$

for the former, Bliss suggested that the latter could be answered by analyzing

$$c_i = F^{-1}(p_i).$$

The “Method of Probits” is carried out simply by referring the percent kill,  $p_i$  to a table to determine the value of  $c_i$  of interest. Of course, the question can be answered from the figure by moving eastward from the kill rate of interest to the figure then downward to the concentration. A common application involved elicitation of the lethal dose needed to achieve a 50% kill rate, denoted LD50. [See Finney (1944a,b,1947a) or (1971), for examples.]

The obvious flaw in the method just described (by the authors, not Bliss) is that different situations would provide different shaped curves, and the preceding provides no accommodation of that. His search of the then current literature suggested to Bliss that analysts had used a variety of freehand drawing methods to accommodate this kind of heterogeneity, methods that were subject to errors and approximations. Bliss (1934a, p. 38) goes on to suggest “It is believed that these and other difficulties can be minimized if percentage kill and dosage are transformed to units which may be plotted as straight lines on ordinary cross section paper and hence permit fitting by the customary technique of least squares or of the straight line regression equation.”

Superficially, Bliss suggests that the preceding model be modified to accommodate the heterogeneity

$$p_i = F(\alpha + \beta c_i).$$

What is needed for the “transformation to units...” is a definition of the specific function,  $F(\cdot)$ , for which he chose the normal distribution. The inverse transformation is

$$\alpha + \beta c_i = F^{-1}(p_i) = \Phi^{-1}(p_i) = \text{normit}(p_i) = y_i.$$

This being 1934, computation of the normits is another difficult hurdle. Bliss relied on a table published by Pearson (1914, “Tables of the Normal Probability Integral” in *Pearson’s Tables for Statisticians and Biometricians* which is reproduced below). Dealing with negative numbers was a complication of some substance in 1934, so Bliss suggests the “probability unit” or “probit”

$$\text{probit}(p_i) = \text{normit}(p_i) + 5.$$

Probits for a number of values of  $p_i$  are given in Bliss’s Table I reproduced below in Figure 2.

These are Bliss’s probits. Note that the value associated with 50% is 5.00, not 0.00. A remaining problem is how to handle the extreme tail values. The author assigned the value 0.00 to 0.01% and 10.00 to 99.99%. The level of inaccuracy for the intervening values was taken as tolerable. It is intriguing to note, the Pearson Tables (volumes of them) were themselves computed by hand (around 1910). Indeed, though the accuracy of the figures in Bliss’s table is noteworthy given when and how they were computed, it is, in fact, quite lacking in absolute terms. Figure 3 shows the percentage error in Bliss’s (Pearson’s) probits (computed using a modern computer and the INP(.) function in *NLOGIT*). It is intriguing to see that the errors are quite large at the tails and clearly not random. An approximation was being used that systematically degrades as the probability moves away from .5 in either direction.

TABLE I

Per cent. kill	Probits						
1.0	1.87	50.0	5.00	80.0	6.13	95.0	7.21
5.0	2.79	52.0	5.07	81.0	6.18	96.0	7.35
10.0	3.28	54.0	5.14	82.0	6.23	97.0	7.53
15.0	3.61	56.0	5.20	83.0	6.28	98.0	7.76
20.0	3.87	58.0	5.27	84.0	6.34	98.5	7.92
25.0	4.09	60.0	5.34	85.0	6.39	99.0	8.13
30.0	4.30	62.0	5.41	86.0	6.45	99.1	8.18
34.0	4.44	64.0	5.48	87.0	6.51	99.2	8.24
36.0	4.52	66.0	5.56	88.0	6.58	99.3	8.30
38.0	4.59	68.0	5.63	89.0	6.65	99.4	8.38
40.0	4.66	70.0	5.70	90.0	6.72	99.5	8.46
42.0	4.73	72.0	5.78	91.0	6.80	99.6	8.57
44.0	4.80	74.0	5.86	92.0	6.89	99.7	8.69
46.0	4.86	76.0	5.95	93.0	6.98	99.8	8.87
48.0	4.93	78.0	6.04	94.0	7.09	99.9	9.16

Figure 2. Table of Probits for Values of  $p_i$ .

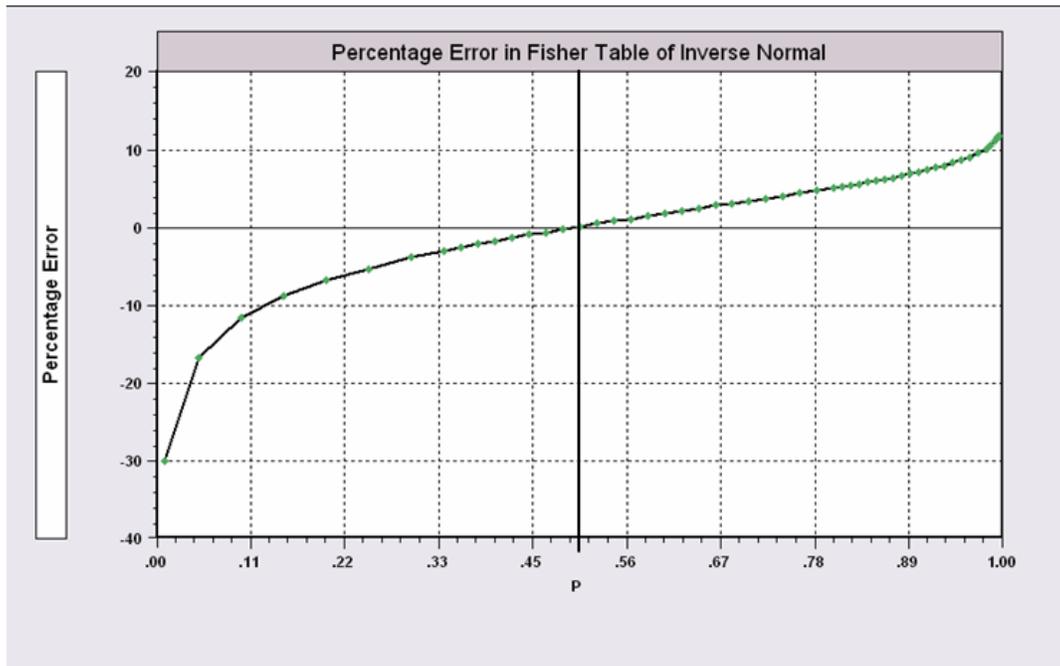


Figure 3 Percentage Errors in Pearson Table of Probability Integrals

As the model is stated above, any two points suffice to determine  $\alpha$  and  $\beta$ . To accommodate the inevitable sampling variability, the (implied) model must be modified to

$$p_i = \Phi(\alpha + \beta c_i + \varepsilon_i).$$

No assumption about the distribution of  $\varepsilon_i$  is necessary;  $\varepsilon_i$  is just sampling variability. A mean or median of zero would be a convenient normalization. Bliss then suggests the method of least squares to estimate  $\alpha$  and  $\beta$ , which might suggest that he relied (again implicitly) on symmetry of the random errors,  $\varepsilon_i$ . This would be the evident origin of probit analysis. Other authors had been doing similar analyses for years. But, this was the first point at which the technique was formalized using the inverse probability function (and the normal distribution.) [In Bliss (1934b), the author notes that two other researchers, Hemmingsen (1933) and Gaddum (1933) had used essentially the same method in a study of toxicity in mice.]

Bliss cites several advantages of his method:

- (1) It provides a test of normality (of  $\varepsilon$ ). (One could examine the variation of  $F^{-1}(p_i)$  around the fitted regression line.)
- (2) It includes the ability to do the analysis using logarithms. [See Greene, Knapp and Seaks (1993).] (At least it makes it simpler.)
- (3) It suggests a method of determining whether organisms exposed to each dosage were equivalent and the amounts administered experimentally were uniformly proportional to the effective dosage over the range covered by the experiment. (This is examined by exploring the regression relationship.)
- (4) It allows the analyst to see “the disclosure of change in the mode of lethal action with certain poisons over different sections of the dosage range indicated by an abrupt change in the slope of the regression.” The figure that is shown for this case in the article (shown below as Figure 4) is equivalent to the introduction of a linear spline in the function based on the log of the dosage, i.e.,

$$p_i = \Phi\{\alpha + \beta \log Dosage_i + \gamma[1(\log Dosage_i > 1.35) \times (\log Dosage_i - 1.35)] + \varepsilon_i\},$$

which is strikingly modern. [See Greene (2008a, pp. 111-112).]

- (5) It allows a simple method of expressing in the slope of a straight line, the relative uniformity or diversity between individuals in their susceptibility to a poison. (This seems to relate to the inherent variability of freehand methods used previously.)

In three editions of his celebrated book on the subject of probit analysis, Finney (1947b, 1952, 1971) refined Bliss’s methods and applied them to a wide array of experiments. The major practical development in the progression of this work was the advent of software and computers for maximum likelihood methods, including Finney’s own contribution to this market, a program that he named BLISS in recognition of his predecessor. [See ISI (1982).]

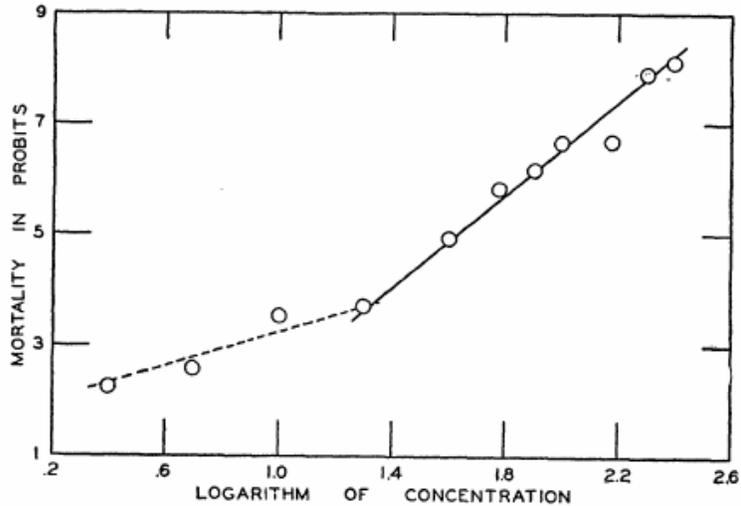


Fig. 2. Data in Fig. 1 converted to rectilinear form by use of logarithms and probits as explained.

Figure 4. Implied Spline Regression in Bliss's Probit Model

### 3.2 Social Science Data and Regression Analysis for Binary Outcomes

To this point, and in the studies noted below, the collection of methods is applied to sampling situations involving grouped data, that is proportions. The samples involved in the analyses described here consisted of observations  $(n_i, p_i, x_i)$ ,  $i = 1, \dots, N$ . That is, a group size, a proportion of “responders” and a level of the stimulus. The literature was into the 1970s before researchers began to extend the techniques to individual data. See, for example, the “Frontiers” section of Theil (1971). The formal treatment of individual data for ordered choices – the sort of data observed by social scientists – begins with Walker and Duncan (1967) in the bioassay literature and appeared first in the social sciences in 1975 with McElvey and Zavoina.

The development of the “minimum chi squared” approach to estimation, and the development of estimation methods as something closer than before to regression analysis might be seen as a bridge between these literatures. Berkson (1944, 1953, 1955a,b, 1957, 1980) and Amemiya (1975, 1980, 1985) suggest an approach to estimation along the lines of

$$p_i = F(\alpha + \beta c_i) + \varepsilon_i.$$

[Walker and Duncan (1967), drawing on Gurland and Dahm (1960), also took precisely this approach to modeling probabilities (see p. 169). However, they were concerned with individual data, not sample proportions. We will examine Walker and Duncan’s analysis in Section 3.5.] That is, the sampling variability in estimation is laid on the sample proportion,  $p_i$ , as an estimator of the population quantity,  $F(\alpha + \beta c_i)$ . Under this interpretation, the logit,  $\log p_i/(1-p_i)$  or normit transformation,  $\Phi^{-1}(p_i)$  would seem to be less useful, since now the sampling variability is moved inside the function. A two step or iterative application of weighted least squares, the *minimum chi squared* estimator provides an approach that accounts for the nonlinearity of the function and the heteroscedasticity in  $p_i$ . [See, e.g., Greene (2003, Section 21.4.6).]

The analysis of the population probability,  $F(\alpha + \beta c_i)$ , as the conditional mean in a regression relationship can be carried over to a setting of individual data. This line of approach comes to fruition in the class of “Generalized Linear Models,” [McCullagh and Nelder (1983).]

The GLIM approach to modeling binary data embodies the regression interpretation of the probability function and extends easily to the analysis of individual data.

### 3.3 Analysis of Binary Choice

By 1975, analysis of binary data by social scientists, in grouped or individual form, using maximum likelihood, or minimum chi squared estimators had come to full bloom. The GLIM approach [Grizzle et al. (1969), Nelder and Wedderburn (1972) and Wedderburn (1974), and, see, McCullagh and Nelder (1983) and Pregibon (1984)] had likewise become in bioassay. Surveys of varying length of estimation involving binary choices are given in Cox (1970), Finney (1971), Amemiya (1981), Long (1997), Greene (2007a, 2008a) and dozens of other primers and introductions.

### 3.4 Ordered Outcomes: Aitchison and Silvey (1957), Snell (1964),

Analysis of a dichotomous response (always in grouped form, however), is well developed by the 1940s. Analysis of ordered responses that are of interest in this study, begins in 1957 with an extension to Finney by Aitchison and Silvey (1957). The other relevant antecedent is Snell's (1964) parallel development of an (only apparently) different treatment of ordered outcomes. In what follows, we will use the authors own notation, though contemporary treatments use a uniformly different flavor of notation.

The modeling exercise considered by Aitchison and Silvey (1957) is as follows: Sample observations are made on a species of insect *Petrobius* Leash (Thysanura, Machilidae) that passes through  $s+1$  stages in its life cycle. An insect is necessarily observed in one stage at any point in time. The last stage is always reached. Observations are made at  $m$  different times, denoted  $x_\alpha$ ,  $\alpha = 1, \dots, m$ .

The amount of time spent by an insect in stage  $i$ , ( $i=1, \dots, s$ ) is an observation on a nonnegative random variable,  $\xi_i$ . Interest is in estimation of  $\lambda_i = E[\xi_i]$  = the average amount of time that will be spent in stage  $i$ . The total time spent in stages  $1, \dots, r$  is  $\eta_r = \sum_{i=1}^r \xi_i$ , also a nonnegative random variable. Interest might be in estimation of  $\mu_r = E[\eta_r]$  as well. Since  $\lambda_i = \mu_i - \mu_{i-1}$ ,  $\lambda_i$  is estimable from  $\mu_i$ .

Total time spent in stages up to the observation,  $\eta_r$ , is a continuous random variable with cdf  $\Pr(\eta_r \leq x) = G_r(x)$ . Probabilities of observation of an insect in the  $s+1$  stages at time  $x$  are

$$\begin{aligned} \pi_1(x) &= \Pr(\eta_1 > x) = 1 - G_1(x), \\ \pi_2(x) &= \Pr(\eta_1 \leq x \text{ and } \eta_2 > x) \\ &= \Pr(\eta_1 \leq x) - \Pr(\eta_1 \leq x \text{ and } \eta_2 \leq x) \\ &= \Pr(\eta_1 \leq x) - \Pr(\eta_2 \leq x) \\ &= G_1(x) - G_2(x). \end{aligned}$$

(This makes use of the result that if  $\eta_r < x$ , then  $\eta_{r-1} < x$ .)

$$\begin{aligned} \pi_s(x) &= G_s(x) - G_{s-1}(x), \\ \pi_{s+1}(x) &= G_s(x). \end{aligned}$$

The proportions of insects (subjects) observed in stage  $s$  at time  $x$ ,  $p_s(x)$  are moment estimators of  $\pi_s(x)$ . Estimation of the means is based on the model assumption that the random variables  $\eta_r$  are normally distributed with mean  $\mu_r$  and standard deviation  $\theta_r$ , so

$$G_r(x) = \Phi[(x - \mu_r)/\theta_r]$$

The authors consider method of moments estimation of  $\mu_r$  and  $\theta_r$ . Let  $p_{ar}$  denote the sample estimate of  $\pi_r(x_\alpha)$ . That is,  $p_{ar}$  is the proportion of subjects in stage  $r$  at time  $x_\alpha$ . Then the relationship above suggests

$$\Phi^{-1}(p_{ar}) = Y_{ar} = x_\alpha/\theta_r - \mu_r/\theta_r.$$

They observe, then, “for given  $r$  a straight line fitted to the points  $(x_\alpha, Y_{ar})$  will cross the  $x$ -axis near the maximum likelihood estimate of  $\mu_r$ , while the gradient will approximate to the maximum-likelihood estimate of  $-\theta_r^{-1}$ .” By this device, all the parameters of this model may be estimated. Some obvious problems will arise with data sets in which  $p_{ar}$  is near zero or one. Moreover, estimation of the scale parameters was complicated (this being 1955), so they considered model simplifications, arriving at  $\theta_r^2 = \sigma^2\mu_r$  and then using, instead, maximum likelihood based on the method of scoring. The authors, noting the connection to Finney’s work, label this a “generalized probit model.” Although the preceding does not involve the same sort of estimation problem as Finney’s (in short, the coefficient on  $x$  in this model is  $1/\theta_r$  and we are, in principle, only estimating the threshold values), there is an obvious relationship. They state (p. 139)

Clearly a situation might arise where in place of a simple dichotomy, [Finney’s case] subjects are divided into more than two classes by any dose of the stimulus. Accordingly, we envisage an experiment where random samples of subjects are subjected to  $m$  doses  $x_\alpha$  ( $\alpha = 1, 2, \dots, m$ ) of a stimulus and as a result of the application of the dose  $x_\alpha$  each subject is placed in one of  $s+1$  classes. A straightforward illustration of such an experiment is given by Tattersfield, Gimmingham and Morris (1925) who classified insects subject to a poison as *unaffected*, *slightly affected*, *moribund* or *dead*. The particular problem discussed above is another illustration if, in this case, time is regarded as the stimulus.

Thus, Aitchison and Silvey have clearly laid the foundation for the ordered probit model as we now understand it, albeit, the application described does not resemble it very closely at all. They go on to suggest conditions that “must be satisfied in this general experiment in order that the method of analysis used in our particular case should be applicable”

- (i) The classes must be ordered, mutually exclusive and exhaustive.
- (ii) The reactions of a subject to increasing doses must be systematic in the sense that if dose  $x$  places a subject in the  $i$ th class, then a dose greater than  $x$  is required to place this subject in the  $j$ th class where  $j$  is greater than  $i$ .

Point (i) is obvious – the model is designed for ordered outcomes. The second point seems to relate to the latent regression interpretation of the modern view of the model. The authors discuss a “tolerance” for the given classes defined in the model which the surrounding discussion associates with levels of a latent variable that is observed by the analyst only through the class observed. Finally, the authors note that “if  $s = 1$  then the present analysis becomes an ordinary probit analysis and it is in this sense that we have generalized probit analysis.”

Before leaving Aitchison and Silvey, it is interesting to note that although their application did not actually generalize probit analysis, the speculation in the paragraph noted above, in fact, did. The application that they pursued is extended by Feinberg (1980) in what he calls the *continuation ratio model*. [See, as well, Long and Freese (2006, pp. 221-222).] The model is a regression style model that is designed for sequential (so, by implication, ordered) outcomes. The example given by Long and Freese is faculty rank, which would typically include assistant, then associate, then full professor (and perhaps instructor at the left and chaired professor at the right). The functional form is written for  $m$  stages in the progression in which the

probability that an observed individual is in stage  $m$  given  $\mathbf{x}$  is  $\Pr(y = m|\mathbf{x})$  and the probability that they are in a higher stage is  $\Pr(y > m|\mathbf{x})$ . Then, the “continuation model” for the log odds is

$$\log \left[ \frac{\Pr(y = m | x)}{\Pr(y > m | x)} \right] = \theta_m - \boldsymbol{\beta}' \mathbf{x}.$$

It is not obvious how the ordering aspect of the outcomes enters this model. The requirement in the model (and in university life) that for a given individual,

$$\Pr(y \leq m|\mathbf{x}) < \Pr(y \leq m+1|\mathbf{x})$$

is induced by the fact that  $m+1$  means there are more ranks at or below  $m$  than  $m+1$ , not that the next rank has a higher order than the previous one. For the scenario described, the flaw in the model would seem to be that it is a static model being used to describe a dynamic phenomenon. Although one must pass through the stages in order (though individuals have been known to skip stages), the probabilities in the model do not have any intrinsic relationship to the ordering of the stages but rather arise the same way if we merely count ranks.

Snell (1964) considers specifically analyzing a set of scores for a ranked set of outcomes such as *Excellent, Very Good, Good, Not Very Good, Poor, Very Poor*, recorded, perhaps, 6,5,4,3,2,1 or the like. Conventional analysis of such data (Aitchison and Silvey (1957) notwithstanding) was done using analysis of variance techniques, e.g., regression methods assuming (a) normally distributed disturbances and (b) homogeneous variances.

Their departure point is “[w]e assume there to be an underlying continuous scale of measurement along which the scale categories represent intervals.” The scale is divided into intervals labeled  $k = 0, 1, \dots, k$  by  $k+2$  points,  $x_{-1}, x_0, x_1, \dots, x_k$ . Observations in the data, indexed by  $i$ , consist of group size,  $n_i$  and proportions,  $p_{ij}$ ,  $j = 0, 1, \dots, k$ . The underlying continuous distribution function is denoted  $P_i(x_j)$ . It is unclear what continuous random outcome this is meant to refer to, in connection to the “ $i$ .” However, it is obvious from the context that in fact what is implied is that we describe the realization of a random variable,  $X_i$  which is the unobserved aforementioned “measurement.” Thus, by the construction above, the probability of observing an individual that is in group  $i$  will be in category  $s_j$  is equal to

$$P_i(x_j) - P_i(x_{j-1}), \quad i = 1, \dots, m; \quad j = 0, \dots, k.$$

Once again, the reference to “ $i$ ” above refers to a group, so it can only be inferred that what the author has in mind is that group “ $i$ ” consists of  $n_i$  realizations of  $X_i$ , and the preceding gives the probabilities associated with each member of the group. (Note that there is nothing so far in the data other than the observation subscript,  $i$ , to distinguish the groups, e.g., no stimulus  $x_i$ .) To continue, “We take the distribution function to be of the form”

$$P_i(x_j) = [1 + \exp(-f_{ij})]^{-1} = \Lambda(f_{ij}) \quad (\text{using a contemporary notation}).$$

Finally,  $f_{ij}$  is defined to be the “logit” of the “proportion”  $P_i(x_j)$ ,

$$f_{ij} = \log[P_i(x_j)/(1 - P_i(x_j))] = a_i + b_i x_j.$$

The model now has for each  $i$ , a location parameter  $a_i$  and a spread parameter  $b_i$ . To impose homoscedasticity on the data, they assume  $b_i = 1$ . The log likelihood for the observed data is

$$\log L(a_1, \dots, a_m, x_{-1}, x_0, \dots, x_k) = \sum_{i=1}^m n_i \sum_{j=0}^k p_{ij} \log [P_{ij} - P_{i,j-1}].$$

It is apparent that a normalization is required to use the entire real line, so  $x_0 = -\infty$  and  $x_k = +\infty$ . He also notes “since the choice of origin is arbitrary, we take  $x_1 = 0$ .” (In fact, since there is no other invariant constant term in the model, this last normalization is not necessary – it now constitutes a substantive restriction.) The remainder of the analysis focuses on methods of estimating  $m$  fixed effects  $a_i$  and  $k-2$  threshold values,  $x_j$ .

The parameters of the model can be loosely estimated by a method of moments type of calculation. Approximate estimates of the threshold values  $x_k$  are based on group size weighted averages of the group proportions. Initial estimates of the fixed effects are computed using

$$a_i = -\sum_{j=1}^k p_{ij} s_j$$

where  $s_j = (x_j - x_{j-1})/2$ ,  $j = 2, 3, \dots, k-1$ . The two end points corresponding to the lower and upper tails are problematic, and a solution, ultimately,  $s_1 = x_1 - 1$  and  $s_k = x_{k-1} + 1$ , is suggested. Newton’s method is ultimately used to complete the estimation.

Snell’s model is functionally equivalent to  $P_i(x_j) = \Lambda[x_j - (-a_i)]$  so that the log likelihood function is

$$\log L(a_1, \dots, a_m, x_{-1}, x_0, \dots, x_k) = \sum_{i=1}^m n_i \sum_{j=0}^k p_{ij} \log [\Lambda(x_j + a_i) - \Lambda(x_{j-1} + a_i)]$$

This corresponds essentially to a modern form of the ordered choice model, though it should be noted that the assumption of a different “effect,”  $a_i$  for each cross section observation does not appear in the recent literature. (It is estimable, perhaps counter to intuition, because there is more than a single observation for each  $i$ ; there is a whole set of  $p_{ij}$ s for each  $i$ .)

It is worth noting as well, that the terms in the log likelihood function above are only positive if the  $x_j$  terms are strictly ordered. The initial, “approximate” values will certainly be, because they are functions of the cumulative group proportions. But, the application of Newton’s method that follows makes no mention of this restriction, and could break down numerically. The method was only suggested in the text; the author used the approximate, method of moments estimators in the applications.

Some of the closing remarks in the paper are intriguing.

“The aim throughout this paper has been to present a method based upon a theoretical model and yet to keep the procedure as simple as possible. For this reason, attention has been directed very much towards an approximate solution.”

The method of solution is the method of moments; in principle it could have been done with a hand calculator. (In 1964, Texas Instruments had just begun production of their first four function calculators, so that might have been optimistic. However, IBMs 7090 series of mainframe computers was already well established and the 360 series was on the near horizon. There would have been no shortage of computing power. A computing language, Fortran (Formula Translation), had been invented in the 1950s.) Snell does note that the iterative method “can easily be carried out on a desk machine, and one iteration should be sufficient.”

“The model upon which the method is based takes no account of the experimental design behind the data.”

We read this to state that there is no data generating process assumed to be at work here (though, in fact, there must be one in the background – the data arise through some kind of stochastic process; we have attached probabilities to the outcomes.) In fact, the method is semiparametric – the fixed effects approach does stop short of regression. However, of course, the choice of logistic distribution was not entirely innocent. It was made for mathematical convenience, however the numerical results depend on it. The same set of computations could have been done, at considerable cost in complexity, using the normal distribution.

“Finally, there is no reason why the use of this technique should be restricted to subjective measurement.”

Indeed, the recent history has demonstrated the versatility of the method.

### 3.5 Minimum Chi Squared Estimation of an Ordered Response Model: Gurland et al. (1960)

Gurland, Lee and Dahm (1960) considered the following analysis in bioassay (p. 383): [We will modify their notation slightly so that their model will fit more neatly into the discussion used herein.]

Suppose  $N$  groups consisting of  $n_1, \dots, n_N$  houseflies are exposed to dosages  $x_1, \dots, x_N$ , respectively. Out of the  $n_i$  flies exposed at dosage  $x_i$ , suppose that at the given time of observation,

$r_{i1}$  are dead,  $r_{i2}$  are moribund,  $r_{i3}$  are alive.

Write the observed proportions as

$$p_{i1} = r_{i1}/n_i, p_{i2} = r_{i2}/n_i, p_{i3} = r_{i3}/n_i = 1 - p_{i1} - p_{i2}.$$

Let

$$P_{i1} = E[p_{i1}], P_{i2} = E[p_{i2}], P_{i3} = 1 - P_{i1} - P_{i2}$$

be the corresponding expected proportions or true probabilities. Then, ...

$$P_{i1} = \Phi(\alpha_1 + \beta x_i) \quad (1)$$

$$P_{i1} + P_{i2} = \Phi(\alpha_2 + \beta x_i), i = 1, \dots, N \quad (2)$$

where

$$\beta = 1/\sigma, \alpha_1 = -\mu_1/\sigma, \alpha_2 = -\mu_2/\sigma.$$

... This assumes a normal tolerance distribution  $N[\mu_1, \sigma^2]$  of lethal dosages and a normal tolerance distribution  $N[\mu_2, \sigma^2]$  of moribund dosages. Furthermore,  $\mu_1 > \mu_2$ . Since a fly becomes moribund before it dies, the expression in (2), which is the probability a fly is moribund or dead, must involve the same parameter,  $\beta$  as in (1). If the  $\beta$  were not common, the two curves would cross, but this is obviously not permissible since  $P_{i1} + P_{i2} > P_{i1}$ .

Note, first, the interpretation of  $P_{ij}$  as  $E[p_{ij}]$  implies  $p_{ij} = P_{ij} + \varepsilon_{ij}$ , precisely as in Section 3.2. The authors propose a regression approach to estimation of the model parameters, as opposed to maximum likelihood estimation. They proceed to develop a weighted least squares (minimum

chi squared) estimator. Second, presumably, the normal distributions assumed above apply to the distributions of tolerances across individual flies. It follows from their analysis, then, that for any particular housefly,  $t = 1, \dots, n_i$ ,

$$\begin{aligned} \text{Prob}(dead_{it}|x_i) &= \Phi[-\mu_1/\sigma + (1/\sigma)x_i] \\ &= \text{Prob}[T^* \leq (1/\sigma)x_i - \mu_1/\sigma] \\ \text{Prob}(dead_{it}|x_i) + \text{Prob}(moribund_{it}|x_i) &= \Phi[-\mu_2/\sigma + (1/\sigma)x_i] \\ &= \text{Prob}[T^* \leq (1/\sigma)x_i - \mu_2/\sigma] \end{aligned}$$

Where  $T^*$  is the tolerance across flies in the experiment. This would appear to be precisely the model ultimately analyzed by McElvey and Zavoina (1975). There is a loose end in the preceding which makes the model an imperfect precursor, however. The authors have avoided the latent regression – they make no mention of it. They state specifically that there are different tolerance distributions with the same variance but different means. But, they do force the same  $\beta$  to appear in both probabilities, arguing that without this restriction, we will be able, for some dosage,  $x_i$  to have the probability of dead or moribund be less than that the probability of dead, which is a contradiction of the axioms of probability. It does follow, however, that there are different *prior* distributions for flies that will die after dosage  $x_i$  and flies that will be moribund – i.e., the different tolerance distributions. Thus, there is an ambiguity in the formulation as to what random variable the assumed normal distributions are meant to describe. By a reasonable construction, for example, we might infer that the distribution describes the observed flies only after the reaction to the dosage.

The ambiguities notwithstanding, Gurland et al. (1960) have laid the platform for analysis of ordered outcomes with something resembling a regression approach. The approach is still, however, focused on the analysis of sample proportions. The minimum chi squared (iterated weighted least squares) estimator that they develop is proposed because it “is simpler to apply.”

### 3.6 Individual Data and Polychotomous Outcomes: Walker and Duncan (1967)

Walker and Duncan (1967) were concerned with the problem of using a large number of covariates to analyze the probabilities of outcomes. The experiment in the study involved four large surveys of individuals who were free of heart disease at entry to the study and who were examined long after for presence of (1) myocardial infarction (*MI*), (2) angina pectoris (*AP*) and (3) no coronary heart disease ( $\overline{CHD}$ ). After considering whether the first two categories might be unordered or ordered, the authors opted to build a model for the latter. Previous analyses had studied crosstabulated data based on one or two factors and by age and sex. The use of numerous other factors – the application involved 8 in addition to age and sex – mandated a different approach.

The three outcome model follows along the lines of Gurland et al. (1960) with two major exceptions. First, the large number of factors compels analysis of the individual data, rather than the sample proportions. Second, though only in passing, they note a natural characterization of the data generating process as “Considered jointly they involve the further assumption that the state of an individual described by the vector  $\mathbf{x}$ , which is sufficient to entail the more severe form *MI*, is certainly sufficient to entail the less severe form *AP*. *If MI and AP are in reality grades of severity of coronary disease, this assumption will hold at least approximately.* If on the other hand these are distinct, even though closely related diseases, it is not likely to hold.” [Emphasis added.] (p. 173.) Coupled with the assumption of the strict ordering of the outcomes, this does sound like the rudiments of an “underlying regression” interpretation. If so, then the authors’

assumption of the logistic distribution as shown below completes the formulation of the ordered logit model. Continuing, “The mathematical reflexion of this assumption is seen in the fact that  $P_1 + P_2 \geq P_1$ , which holds if and only if the ‘slope’ coefficient  $\beta$  is identical in (6.1) and (6.2), as is easily shown. (In fact, this is only the case if  $\alpha_2 > \alpha_1$ . Otherwise, it is neither necessary nor sufficient.)

Their three outcome model (where, as before, we have changed their notation a bit for clarity) is,

$$\begin{aligned}
z_{i1} &= 1 \text{ if } MI_i \text{ and 0 otherwise} \\
z_{i2} &= 1 \text{ if } AP_i \text{ and 0 otherwise} \\
z_{i3} &= 1 \text{ if } \overline{CHD}_i \text{ and 0 otherwise} \\
P_1 &= E[z_{i1}|\mathbf{x}_i] \quad (6.1) \\
P_2 &= E[z_{i2}|\mathbf{x}_i] \\
P_3 &= 1 - P_1 - P_2. \\
E[z_{i1}|\mathbf{x}_i] &= P_1 = \Lambda(\alpha_1 + \beta'\mathbf{x}_i) \quad (6.1) \\
E[z_{i1} + z_{i2}|\mathbf{x}_i] &= P_1 + P_2 = \Lambda(\alpha_2 + \beta'\mathbf{x}_i) \quad (6.2)
\end{aligned}$$

To preserve the result  $P_1 + P_2 \geq P_1$ , it must also be true that  $\alpha_2 > \alpha_1$ . The implied model structure is

$$\begin{aligned}
\text{Prob}(MI_i|\mathbf{x}_i) &= \Lambda(\alpha_1 + \beta'\mathbf{x}_i) \\
\text{Prob}(AP_i|\mathbf{x}_i) &= \text{Prob}(Heart Disease|\mathbf{x}_i) - \text{Prob}(MI_i|\mathbf{x}_i) = \Lambda(\alpha_2 + \beta'\mathbf{x}_i) - \Lambda(\alpha_1 + \beta'\mathbf{x}_i) \\
\text{Prob}(\overline{CHD}_i|\mathbf{x}_i) &= 1 - \Lambda(\alpha_2 + \beta'\mathbf{x}_i).
\end{aligned}$$

Walker and Duncan are the first to pursue the analysis of ordered probabilities with individual data. In fact, the latent regression model is not necessary to reach their model formulation; we have superimposed our own interpretation on their model to obtain it. They, in turn, did not appear quite ready to make the assumption. Their model is only consistent with that specification. Indeed, what they have proposed is a mathematical model of a set of probabilities that preserve the supposed (severity) ordering of the first and second outcomes. No appeal to a latent regression is needed. On the other hand, quite clearly, it is a small extension to broaden this model to include the formal ordered probit regression model proposed by McElvey and Zavoina (1975).

### 3.7 McElvey and Zavoina (1975)

McElvey and Zavoina’s (1975) proposed model is described at length above. Based on the preceding very short chronology, it would seem that their model was a significant jump forward, not an increment to the existing machinery. In fact, neither Aitchison and Silvey (1957) nor Snell (1964) proposed anything resembling a regression approach to the analysis of ordered outcomes. There is an obvious hint in this direction at the end of the former, but no direct modification of their proposed model would produce a regression style formulation. Certainly, Walker and Duncan’s model can easily be made consistent with the structure of McElvey and Zavoina. But, McElvey and Zavoina were the first to formalize the model in terms of an individual choice setting based on a theory of regression, and to develop an effective iterative method of estimation. Walker and Duncan were in similar territory, but they relied on a weighted least squares procedure and an algorithm based on a Kalman filter [Kalman (1960)] that has not reappeared in the literature. McElvey and Zavoina and Walker and Duncan were the also the first analysts to propose using individual data. The predecessors relied entirely on grouped data (proportions), essentially on the method of moments (or maximum likelihood in a few cases).

### 3.8 Developments Since McElvey and Zavoina

As noted earlier, McCullagh (1977, 1979, 1980) is credited with codiscovering the ordered choice model. The proposed model, shown below, is precisely a counterpart to the ordered probit model. However, McCullagh stopped short of hanging the framework on a latent regression. Though he departs from “Motivation for the proposed model is provided by appeal to the existence of an underlying continuous random variable,” he goes on to state (page 110)

All the models advocated in this paper share the property that the categories can be thought of as contiguous intervals on some continuous scale. They differ in their assumptions concerning the distributions of the latent variable (e.g. normality (after suitable transformation), homoscedasticity etc.). It may be objected, in a particular example, that there is no sensible latent variable and that these models are therefore irrelevant or unrealistic. However, the models as introduced in Sections 2.1 and 3.1 make no reference to the existence of such a latent variable and its existence is not required for model interpretation. If such a continuous underlying variable exists, interpretation of the model with reference to this scale is direct and incisive. If no such continuum exists the parameters of the models are still interpretable in terms of the particular categories recorded and not those which might have obtained had the defining criteria  $\{\theta_j\}$  been different. Quantitative statements of conclusions are therefore possible in both cases although more succinct and incisive statements are usually possible when direct appeal to a latent variable is acceptable.

McCullagh seems to be holding back from a commitment to an underlying regression. As he notes, however, it will emerge ultimately that interpretation of the coefficients of the model without such an assumption becomes a bit ambiguous.

Though the idea of the ordered logit model shown below is sometimes attributed to McCullagh, elements of it appear earlier in Andrich (1979) and Plackett (1974), and McCullagh cites Plackett for some of his results. The model proposed is based on a discrete random variable with “ $k$  ordered categories of the response” with probabilities  $\pi_1(\mathbf{x})$ ,  $\pi_2(\mathbf{x})$ , ...,  $\pi_k(\mathbf{x})$ . (“In the case of two groups,  $\mathbf{x}$  is an indicator variable or two level factor indicating the appropriate group.” This appears to suggest a contingency table sort of analysis, for which, of course, the “ordering” would be superfluous.) The response variable,  $Y$ , takes values  $y = 1, \dots, k$  with the listed probabilities. Define  $\kappa_j(\mathbf{x})$  to be the odds that  $Y \leq j$  given  $\mathbf{x}$ . Then, the “proportional odds model” specifies that

$$\kappa_j(\mathbf{x}) = \kappa_j \times \exp(-\boldsymbol{\beta}'\mathbf{x}), j = 1, \dots, k.$$

The ratio of corresponding odds is

$$\kappa_j(\mathbf{x}_1)/\kappa_j(\mathbf{x}_2) = \exp[-\boldsymbol{\beta}'(\mathbf{x}_1 - \mathbf{x}_2)],$$

which is independent of  $j$  and depends only on the difference between the covariate vectors. Given the odds ratio stated as above and defining  $\gamma_j(\mathbf{x}) = \pi_1(\mathbf{x}) + \dots + \pi_j(\mathbf{x})$ , the *proportional odds model* becomes equivalent to

$$\log[\gamma_j(\mathbf{x})/(1-\gamma_j(\mathbf{x}))] = \theta_j - \boldsymbol{\beta}'\mathbf{x}, j = 1, \dots, k.$$

This is, of course, mathematically identical to the familiar ordered choice model discussed earlier. Formally, using a more recent notation,

$$\text{Prob}[y \leq j] = \Lambda(\theta_j - \boldsymbol{\beta}'\mathbf{x}),$$

which is the ordered logit model. As the author notes, no appeal to an underlying regression model is necessary to achieve this result. Left to be determined is the mechanism by which the observed discrete random variable is assigned to  $k$  exhaustive, exclusive and *ordered* categories. The model is meant to apply to proportions, as shown in a series of applications that follows. The application that follows immediately, however, does fall naturally into the latent continuous measure framework, a study of tonsil sizes in a sample of 1398 children [Holmes and Williams (1954)], shown in Figure 5.

TABLE 1  
*Tonsil size of carriers and non-carriers of Streptococcus pyogenes*

	<i>Present but not enlarged</i>	<i>Enlarged</i>	<i>Greatly enlarged</i>	<i>Total</i>
Carriers	19	29	24	72
Non-carriers	497	560	269	1326
Total	516	589	293	1398

**Figure 5 McCullagh Application of Ordered Outcomes Model**

For the simple case shown above, interpretation of the  $\beta$  in the “regression” will be simple, as it will highlight the differences in the probabilities or odds for the outcomes in the two groups. For more complicated kinds of regressors, for example, if age, height, or weight appeared in the data set above, then interpretation of the coefficients would be much more complicated without resort to a regression model of some sort and a notion of “holding other things constant.” In his analysis of this data set, Tutz (1990, 1991) argues that the higher outcomes (more to the right) can only be reached by passing through the lower ones. This calls for a different approach, which he labels the *sequential model*. The simplest case would be Agresti’s (1984) *continuation ratio model*,

$$\text{Prob}(y = r \mid y \geq r, \mathbf{x}) = D(\theta_r - \beta' \mathbf{x})$$

where  $D(\cdot)$  is a transformation of the index. This yields the unconditional probabilities

$$\text{Prob}(y = r \mid \mathbf{x}) = D(\theta_r - \beta' \mathbf{x}) \prod_{i=1}^{r-1} [1 - D(\theta_i - \beta' \mathbf{x})].$$

A variety of extensions are suggested. [For another survey of this and related models, see Barnhart and Sampson (1994).]

Anderson and Philips (1981) continue McCullagh’s development in two directions. Researchers in this area work back and forth around the assumption of the latent continuous variable and latent regression. Second, they introduced some results related to functional form. As noted earlier, their departure point is “... an ordered categorical variable is a coarsely measured version of a continuous variable not itself observable.” The model proposed is as follows: “[I]ndividuals are grouped into  $k$  ordered groups which are identified by an ordered categorical variable  $y$  with arbitrarily assigned value  $s$  for the  $s$ th ordered group;  $s = 1, \dots, k$ . ... The ordering of groups is not, in general, based on any numerical measurement.” (The authors are holding back from the assumption. However, one might ask, on what basis is the ordering of groups assigned if not some underlying quantitative measure?) A regressor vector,  $\mathbf{x}$ , is defined. The Plackett (1974, 1981) and McCullagh (1980) functional form is

$$\text{Prob}(y \leq s | \mathbf{x}) = \frac{\exp(\theta_s - \boldsymbol{\beta}'\mathbf{x})}{1 + \exp(\theta_s - \boldsymbol{\beta}'\mathbf{x})}, s = 0, 1, \dots, k,$$

where  $\theta_1 \leq \theta_2 \leq \dots \leq \theta_{k-1}$ ,  $\theta_0 = -\infty$ ,  $\theta_k = +\infty$ . (The author uses weak inequalities, though in order to prevent zero probabilities for nonnull events, strong inequalities are required.) It follows, as we observed earlier, that

$$\text{Prob}(y = s | \mathbf{x}) = \Lambda(\theta_s - \boldsymbol{\beta}'\mathbf{x}) - \Lambda(\theta_{s-1} - \boldsymbol{\beta}'\mathbf{x}),$$

which is the “logistic model.” This is also labeled the “cumulative odds model” by McCullagh (1980). The authors suggest, instead, that we write

$$\text{Prob}(y \leq s | \mathbf{x}) = \Psi(\theta_s - \boldsymbol{\beta}'\mathbf{x})$$

where  $\Psi(\cdot)$  is a “completely specified cumulative distribution function. This is a generalized “linear” model, but “nonlinear” versions are possible and are referred to in the discussion. The above models will be called *ordered regression models*. (Emphasis added. This is the first occurrence of the term that we have encountered in this literature search.)

The authors justify the model in terms of a latent unobservable,  $z$ , where, conditioned on  $\mathbf{x}$ ,  $z$  has a logistic distribution. Although  $z$  is not observed, a related, grouped version of  $z$ ,  $y$ , is observable. Of course, this is precisely the interpretation that McElvey and Zavoina have provided for the model. (Once again, however, there is no mention of McElvey and Zavoina or their model.) We have on the suggested basis,

$$y = s \text{ if } \theta_{s-1} \leq z < \theta_s \text{ (} s = 1, \dots, k \text{)}.$$

Note that assumptions are made only about the conditional distribution of  $z$  given  $\mathbf{x}$  and  $y$  given  $\mathbf{x}$ . No assumption is made about the marginal distribution of  $\mathbf{x}$ , which prompts the claim that these models make only moderate distributional assumptions.

“Other assumptions are possible for the form of the distribution of  $z$  given  $\mathbf{x}$ . One obvious choice is that this should be the normal distribution,  $N(\boldsymbol{\beta}'\mathbf{x}, 1)$ , leading to the probit model,

$$\text{Prob}(y \leq s | \mathbf{x}) = \Phi(\theta_s - \boldsymbol{\beta}'\mathbf{x}).$$

Here,  $\Phi(\cdot)$  represents the usual probit function. For practical purposes, the logistic and probit models are virtually indistinguishable, but the logistic model of (1) and (2) is often preferred for its computational convenience.” [Anderson and Philips (1981).] Thus, the ordered probit model is (re)born, here in 1981.

Aitchison and Bennett (1970) is occasionally cited as another antecedent to the ordered choice models considered here. In fact, they were concerned with a different setting altogether, though it is intriguing to note that their formulation is precisely that used to motivate McFadden’s conditional logit model (1974). Since they did not consider ordered outcomes, we will forego a detailed discussion of their results.

### 3.9 Other Related Models

Many authors have modified these models at various edges for different situations and types of data. Some major references to examine for details are Agresti (1984, 1990), Clogg and

Shihadeh (1994) and Greenwood and Farewell (1988). Before closing this review, we note two that have particular relevance for our discussion.

### 3.9.1 Known Thresholds

Stewart (1983), Terza (1985) and Bhat (1994) examine a setting in which essentially the conditions of the ordered probit model emerge, save that there is more information about the censoring than merely the categories. An obvious example considered by these authors is given by bracketed income data. When income data are censored into known ranges, then the resulting data generating process is precisely that of the ordered choice model except that the threshold values are known. Suppose, for example, that  $y^*$  = log of income is normally distributed with mean  $\mu = \beta'x$  and variance  $\sigma^2$ , so

$$y^* = \beta'x + \varepsilon,$$

and the censoring mechanism is

$$y = j \text{ if } A_{j-1} < y^* \leq A_j,$$

where  $A_{j-1}$  and  $A_j$  are known. Then, the log likelihood is built up from the probabilities for the observed outcomes;

$$\log \text{Prob}(y = j | x) = \log \left[ \Phi \left( \frac{A_j - \beta'x}{\sigma} \right) - \Phi \left( \frac{A_{j-1} - \beta'x}{\sigma} \right) \right].$$

For this model, the parameters  $\beta$  and  $\sigma$  are both identified (estimable). The ordering of the outcomes is enforced a fortiori by the ordering of the known brackets. This model is, in fact, not a discrete choice model in the spirit of the others that are considered here. Rather, it is a less complicated censoring model more closely resembling the tobit model. [Tobin (1958), Amemiya (1985a, 1985b), Greene (2008a).] There is a temptation to treat this model using linear regression analysis, substituting, e.g., the midpoints of the brackets for intermediate values and some reasonable value for the upper and lower ranges. The temptation should be resisted, since (1) the likelihood for the data and the structural parameters is well defined (and the estimator is available as a preprogrammed procedure in modern software) and (2) least squares in this setting will be inconsistent. The OLS estimator will suffer from truncation bias. The overall result is that because there is variation in  $x$  that is not associated with variation in  $y$ , the OLS slopes will tend to be biased toward zero. The maximum likelihood estimator, which does not display this feature, is easily obtained. We do note, however, if, instead of midpoints, one uses for the substituted values

$$E[y^* | A_{j-1} < y^* \leq A_j, x] = \beta'x + \sigma \left[ \frac{\phi[(A_{j-1} - \beta'x)/\sigma] - \phi[(A_j - \beta'x)/\sigma]}{\Phi[(A_j - \beta'x)/\sigma] - \Phi[(A_{j-1} - \beta'x)/\sigma]} \right]$$

then, with an appropriate iterate for  $\sigma$  as well as this implicit estimator for  $\beta$ , this is equivalent to the EM algorithm [see Dempster, Laird and Rubin (1977)], and is an effective, albeit inefficient way to compute the maximum likelihood estimators of  $\sigma$  and  $\beta$ . (It will be slow to converge compared to other gradient methods such as Newton's method.)

### 3.9.2 Nonparallel Regressions

A second modification of the model, due to Anderson (1984) is of interest here. He notes (p. 4) “The ordering of the categories, or subsets of them, with respect to the regression variables is open to question in some cases. Hence, we start with the logistic regression model suitable for a qualitative, categorical response variable [Cox (1970), Anderson (1972)].” This is

$$\text{Prob}(y = y_s | \mathbf{x}) = \frac{\exp(\beta_{0s} - \boldsymbol{\beta}'_s \mathbf{x})}{\sum_{t=1}^k \exp(\beta_{0t} - \boldsymbol{\beta}'_t \mathbf{x})}$$

where  $\beta_{0k} = 0$  and  $\boldsymbol{\beta}_k = \mathbf{0}$  are introduced to simplify the notation. In fact, the function listed is homogeneous of degree zero, and the “simplifications” are normalizations needed for identification. This is precisely the multinomial logit model developed by McFadden, (1974) and Nerlove and Press (1972). Characteristically (apparently), there is no connection across the branches of the literature. (This being before the Internet, perhaps the lack of connection across disparate literatures is an understandable consequence of the difficulty of a detailed search. We take that sort of thing for granted now.) Anderson proposes this model for unordered categorical outcomes. He notes, in passing, however, that this model often “gives a good fit” even when the  $\boldsymbol{\beta}$ s are “restricted to be parallel.” “This is particularly true when the categories are ordered.” That is to suggest, the ordered choice model considered thus far embodies the *restriction* that the  $\boldsymbol{\beta}$ s are the same. By a simple transformation of the ordered logit model, we find

$$\text{logit}(j) = \log[\text{Prob}(y \leq j | \mathbf{x}) / \text{Pr}(y > j | \mathbf{x})] = \mu_j - \boldsymbol{\beta}' \mathbf{x}$$

which means that  $\partial \text{logit}(j) / \partial \mathbf{x} = \boldsymbol{\beta}$  for all  $j$ . This has come to be known as the “parallel regressions assumption.” [See, e.g., Long (1997, p. 141).] This feature of the model has motivated one form of the “generalized ordered logit” (and probit) model. We will reconsider this generalization of the model in some detail below.

## 4. Estimation, Inference and Analysis Using the Ordered Choice Model

In this section, we will survey the elements of estimation, inference and analysis with the ordered choice model. It will prove useful to develop an application as part of the discussion.

### 4.1 Application of the Ordered Choice Model to Self Assessed Health Status

Riphahn, Wambach and Million (RWM, 2003) analyzed individual data on health care utilization (doctor visits and hospital visits) using various models for counts. The data set is a large panel extracted from the German Socioeconomic Panel (GSOEP). [See RWM (2003) and Greene (2008a) for discussion of the data set in detail.] The data set is an unbalanced panel including 7,293 German households observed from 1 to 7 times and a total of 27,326 observations. (We will visit the panel data aspects of the data and models later.) Among the several interesting variables in this data set is HSAT, a self reported health assessment that is recorded with values 0,1,...,10 (so,  $J = 10$ ). Figure 6 shows the distribution of outcomes for the full sample: The figure reports the variable NewHSAT, not the original variable. Forty of the 27,326 observations on HSAT in the original data were coded with noninteger values between 6.5 and 6.95. We have changed these 40 observations to 7s. In order to construct a compact example that is sufficiently general to illustrate the technique, we will aggregate the categories shown as follows: (0-2)=0, (3-5)=1, (6-8)=2, (9)=3, (10)=4. [One might expect collapsing the data in this fashion to sacrifice some information and, in turn, produce a less efficient estimator of the model parameters. See Murad et al. (2003) for some analysis of this issue.] Figure 7 shows the result, once again for the full sample, stratified by gender. The families were observed in 1984-1988, 1991 and 1995. For purposes of the application, to maintain as closely as possible the assumptions of the model, at this point, we have selected the most frequently observed year, 1988, for which there are a total of 4,483 observations, 2313 males and 2170 females. We will use the variables in the regression part of the model,

$$\mathbf{x} = (\text{constant}, \text{Age}, \text{Income}, \text{Education}, \text{Married}, \text{Kids}).$$

In the original data set, *Income* is HHNINC (household income) and *Kids* is HHKIDS (household kids). *Married* and *Kids* are binary variables, the latter indicating whether or not there are children in the household. Descriptive statistics for the data used in the application are shown in Table 1.

### 4.2 Distributional Assumptions

As suggested earlier, one of the ambiguities in the set of procedures for ordered choice modeling is the distributional assumption. There seems to be little to determine whether the logit, probit, or some other distribution is to be preferred. The logistic model has some mathematical features to recommend it, but any of these, such as the computation of odds ratios can be replicated under other assumptions, perhaps at some minor inconvenience (depending on one's software). The deeper question of how the distributional assumption relates to the model structure remains unresolved. Stewart (2003) proposes, beyond the familiar choices a "seminonparametric generalized ordered probit" that is considerably more complicated than the logit and probit models examined here. The model is automated in a *Stata* command however. Stewart's and other semiparametric approaches are developed in Section 8. We do note, the offered procedure produces coefficient estimates, but it is unclear how these can be translated into

partial effects or other useful quantities. It remains true in this (and all parametric and semiparametric forms) that the vector of partial effects is a scalar multiple of  $\beta$ . On this basis, Stewart argues that ratios of coefficients are useful substitutes for partial effects.

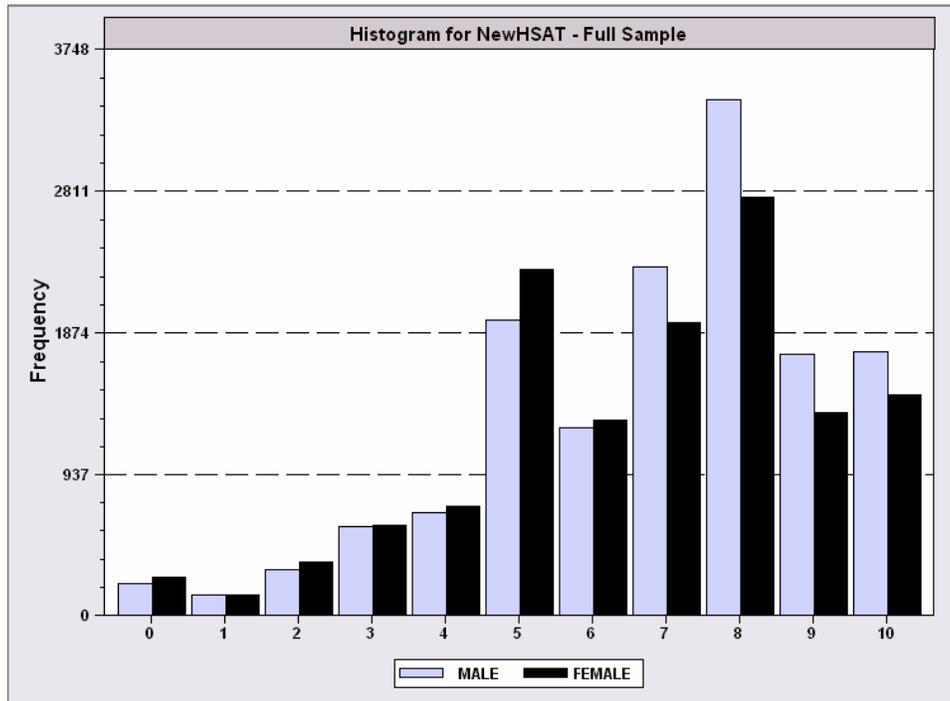


Figure 6 Self Reported Health Satisfaction

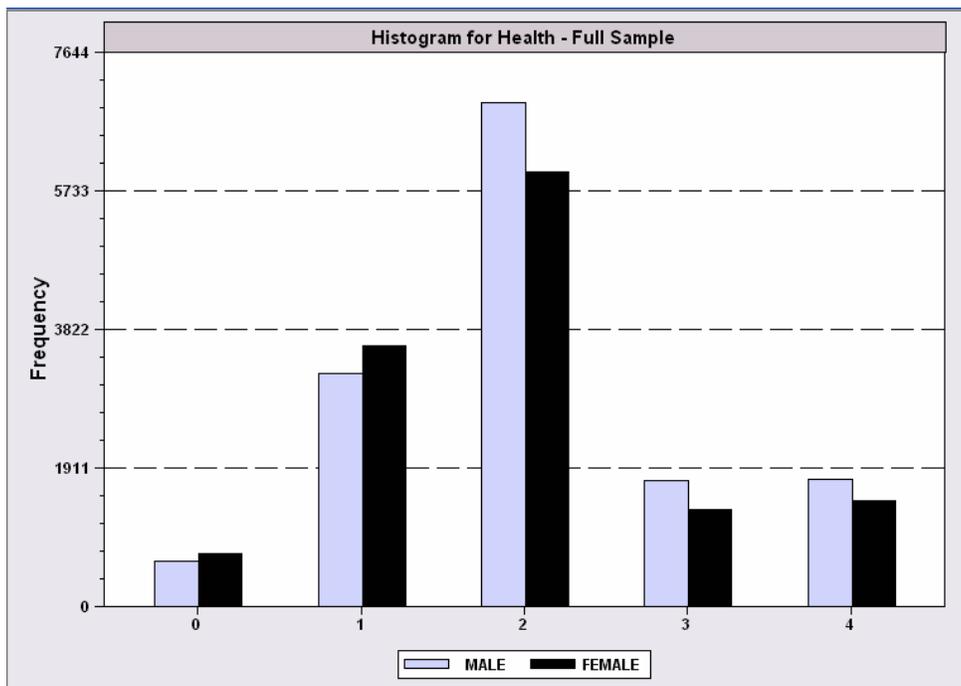


Figure 7 Health Satisfaction with Combined Categories

**Table 1 Data Used in Ordered Choice Application**

Variable	Mean	Std.Dev.	Minimum	Maximum	Cases	Missing
-----+-----						
Stratum is FEMALE =		.000.	Obs.=	2313.000		
AGE	42.7289	11.2966	25.0000	64.0000	2313	2170
EDUC	11.8269	2.49357	7.00000	18.0000	2313	2170
INCOME	.355342	.164814	.000000	2.00000	2313	2170
MARRIED	.756161	.429490	.000000	1.00000	2313	2170
KIDS	.386511	.487055	.000000	1.00000	2313	2170
-----+-----						
Stratum is FEMALE =		1.000.	Obs.=	2170.000		
AGE	44.1982	11.2320	25.0000	64.0000	2170	2313
EDUC	10.9824	2.14195	7.00000	18.0000	2170	2313
INCOME	.341703	.163252	.500000E-02	2.00000	2170	2313
MARRIED	.747926	.434303	.000000	1.00000	2170	2313
KIDS	.371889	.483420	.000000	1.00000	2170	2313
-----+-----						
All observations in current sample						
AGE	43.4401	11.2880	25.0000	64.0000	4483	0
EDUC	11.4181	2.36767	7.00000	18.0000	4483	0
INCOME	.348740	.164183	.000000	2.00000	4483	0
MARRIED	.752175	.431798	.000000	1.00000	4483	0
KIDS	.379433	.485300	.000000	1.00000	4483	0

### 4.3 The Estimated Ordered Probit (Logit) Model

Table 2 presents estimates of the ordered probit and logit models for the pooled data set. (Results from the computer program have been extracted and blended to display the estimates. All computations were carried out using *NLOGIT*. They can all be replicated with equal convenience with *Stata* and, perhaps with a bit more programming, with *EViews*, *TSP*, *SAS* and most other commercial programs.) The tabulated results include diagnostic statistics such as the log likelihood function, a description of the observed data on the outcome, followed by standard presentations of the coefficients, standard errors, etc. These will be examined in detail in the sections to follow.

The estimates for the probit model imply

$$y^* = 1.97882 - .01806Age + .03556Educ + .25869Income \\ - .03100Married + .06065Kids + \varepsilon$$

$$y = 0 \text{ if } y^* \leq 0 \\ y = 1 \text{ if } 0 < y^* \leq 1.14835 \\ y = 2 \text{ if } 1.14835 < y^* \leq 2.54781 \\ y = 3 \text{ if } 2.54781 < y^* \leq 3.05639 \\ y = 4 \text{ if } y^* > 3.05639.$$

Figure 8 shows the implied model for a person of average age (43.44 years), education (11.418 years) and income (0.3487) who is married (1) with children (1). The figure shows the implied probability distribution in the population for individuals with these characteristics. As we will examine in the next section, the force of the regression model is that the probabilities change as the characteristics (**x**) change. In terms of the figure, changes in the characteristics induce changes in the placement of the partitions in the distribution and, in turn, in the probabilities of the outcomes.

**Table 2 Estimated Ordered Choice Models: Probit and Logit**

Ordered Probability Model (PROBIT)		Ordered Probability Model (LOGIT)	
Dependent variable	HEALTH	Dependent variable	HEALTH
Number of observations	4483	Number of observations	4483
Log likelihood function	-5752.985	Log likelihood function	-5749.157
Number of parameters	9	Number of parameters	9
Info. Criterion: AIC =	2.57059	Info. Criterion: AIC =	2.56889
Info. Criterion: BIC =	2.58346	Info. Criterion: BIC =	2.58175
Info. Criterion:HQIC =	2.57513	Info. Criterion:HQIC =	2.57342
Restricted log likelihood	-5875.096	Restricted log likelihood	-5875.096
McFadden Pseudo R-squared	.0207847	McFadden Pseudo R-squared	.0214362
Chi squared	244.2238	Chi squared	251.8798
Degrees of freedom	5	Degrees of freedom	5
Prob[ChiSq > value] =	.0000000	Prob[ChiSq > value] =	.0000000

TABLE OF CELL FREQUENCIES FOR ORDERED PROBABILITY MODEL

Outcome	Frequency		Cumulative < =		Cumulative > =	
	Count	Percent	Count	Percent	Count	Percent
HEALTH=00	230	5.1305	230	5.1305	4483	100.0000
HEALTH=01	1113	24.8271	1343	29.9576	4253	94.8695
HEALTH=02	2226	49.6542	3569	79.6119	3140	70.0424
HEALTH=03	500	11.1532	4069	90.7651	914	20.3881
HEALTH=04	413	9.2349	4483	100.0000	414	9.2349

Variable	Coefficient	Standard Error	b/St.Er.	P[ Z >z]	Mean of X
----------	-------------	----------------	----------	----------	-----------

-----+Index function for probability

Constant	1.97882***	.11616998	17.034	.0000	
AGE	-.01808***	.00161885	-11.166	.0000	43.440107
EDUC	.03556***	.00713213	4.986	.0000	11.418086
INCOME	.25869**	.10387504	2.490	.0128	.3487401
MARRIED	-.03100	.04203080	-.737	.4608	.7521749
KIDS	.06065	.03823694	1.586	.1127	.3794334

-----+Threshold parameters for index

Mu(1)	1.14835***	.02115847	54.274	.0000	
Mu(2)	2.54781***	.02161803	117.856	.0000	
Mu(3)	3.05639***	.02646225	115.500	.0000	

**PROBIT**

-----+Index function for probability

Constant	3.51787***	.20382097	17.260	.0000	
AGE	-.03214***	.00287516	-11.178	.0000	43.440107
EDUC	.06454***	.01247422	5.174	.0000	11.418086
INCOME	.42626**	.18649143	2.286	.0223	.3487401
MARRIED	-.06452	.07455619	-.865	.3868	.7521749
KIDS	.11477*	.06685997	1.717	.0861	.3794334

-----+Threshold parameters for index

Mu(1)	2.12132***	.03705395	57.249	.0000	
Mu(2)	4.43457***	.03902131	113.645	.0000	
Mu(3)	5.37772***	.05199833	103.421	.0000	

**LOGIT**

Note: \*\*\*, \*\*, \* = Significance at 1%, 5%, 10% level.

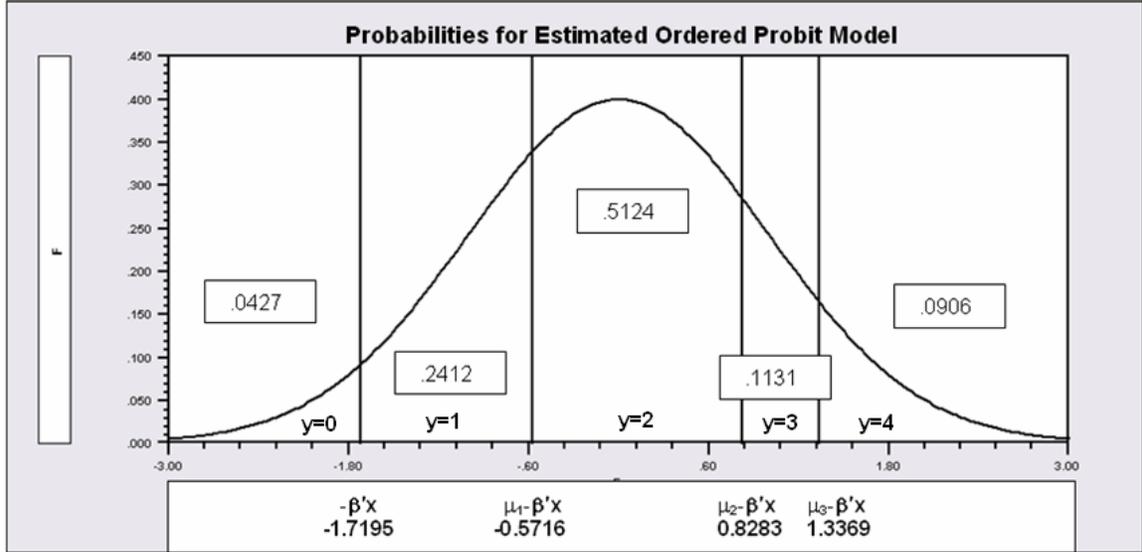


Figure 8 Estimated Ordered Probit Model

#### 4.4 Interpretation of the Model – Partial Effects and Scaled Coefficients

Interpretation of the coefficients in the ordered probit model is more complicated than in the ordinary regression setting. [See, e.g., Daykin and Moffatt (2002).] There is no natural conditional mean function in the model. The outcome variable,  $y$ , is merely a label for the unordered, non-quantitative outcomes. As such, there is no conditional mean function,  $E[y|x]$  to analyze. (This is characteristic of discrete choice models.) In order to attach meaning to the parameters, one typically refers to the probabilities themselves. The partial effects in the ordered choice model are

$$\delta_j(\mathbf{x}_i) = \frac{\partial \text{Prob}(y = j | \mathbf{x}_i)}{\partial \mathbf{x}_i} = [f(\mu_{j-1} - \beta' \mathbf{x}_i) - f(\mu_j - \beta' \mathbf{x}_i)] \beta.$$

A moment's inspection shows that neither the sign nor the magnitude of the coefficient is informative about the result above, so the direct interpretation of the coefficients is fundamentally ambiguous. [A counterpart result for a dummy variable in the model would be obtained by using a difference of probabilities, rather than a derivative. [See Boes and Winkelmann (2006a) and Greene (2007a, Chapter E22).] That is, suppose  $D$  is a dummy variable in the model (such as *Married*) and  $\gamma$  is the coefficient on  $D$ . We would measure the effect of a change in  $D$  from 0 to 1 with all other variables held at the values of interest (perhaps their means) using

$$\Delta_j(D) = [F(\mu_j - \beta' \mathbf{x}_i + \gamma) - F(\mu_{j-1} - \beta' \mathbf{x}_i + \gamma)] - [F(\mu_j - \beta' \mathbf{x}_i) - F(\mu_{j-1} - \beta' \mathbf{x}_i)]$$

(One might on occasion compute the partial effect for a dummy variable by differentiating as if it were a continuous variable. The results will typically resemble the finite change computation, sometimes surprisingly closely – the finite change is a discrete approximation to the derivative. Nonetheless, the latter computation is the more appropriate one.)

The implication of the result above is that the effect of a change in one of the variables in the model depends on all the model parameters, the data, and which probability (cell) is of interest. It can be negative or positive. To illustrate, we consider a change in the education

variable on the implied probabilities in Figure 8. Since the changes in a probability model are typically “marginal” (small), we will exaggerate the effect a bit so that it will show up in a figure. Consider, then, the same individual shown in Figure 8, except now, with a Ph.D. (college plus four years of postgraduate work). That is, 20 years of education, instead of the average 11.4 used earlier. The effect of an additional 8.6 years of education is shown in Figure 9. All five probabilities have changed. The two at the right end of the distribution have increased while the three at the left have decreased.

The partial effects, however computed, give the impacts on the specific probabilities per unit change in the stimulus or regressor. For example, for continuous variable *Educ*, we find partial effects for the five cells of -.0034, -.00885, .00244, .00424, .00557, respectively, which give the expected change on the probabilities per additional year of education. For the income variable, for the highest cell, the estimated partial effect is .04055. However, some care is needed in interpreting this in terms of a unit change. The income variable has a mean of 0.3417 and a standard deviation of 0.1632. A full unit change in income would put the average individual nearly six standard deviations above the mean. Thus, for the marginal impact of income, one might want to measure a change in standard deviation units. Thus, an assessment of the impact of a change in income on the probability of the highest cell probability might be  $0.3417 \times 0.1632 = 0.0558$ . Precisely how this computation should be done will vary from one application to another.

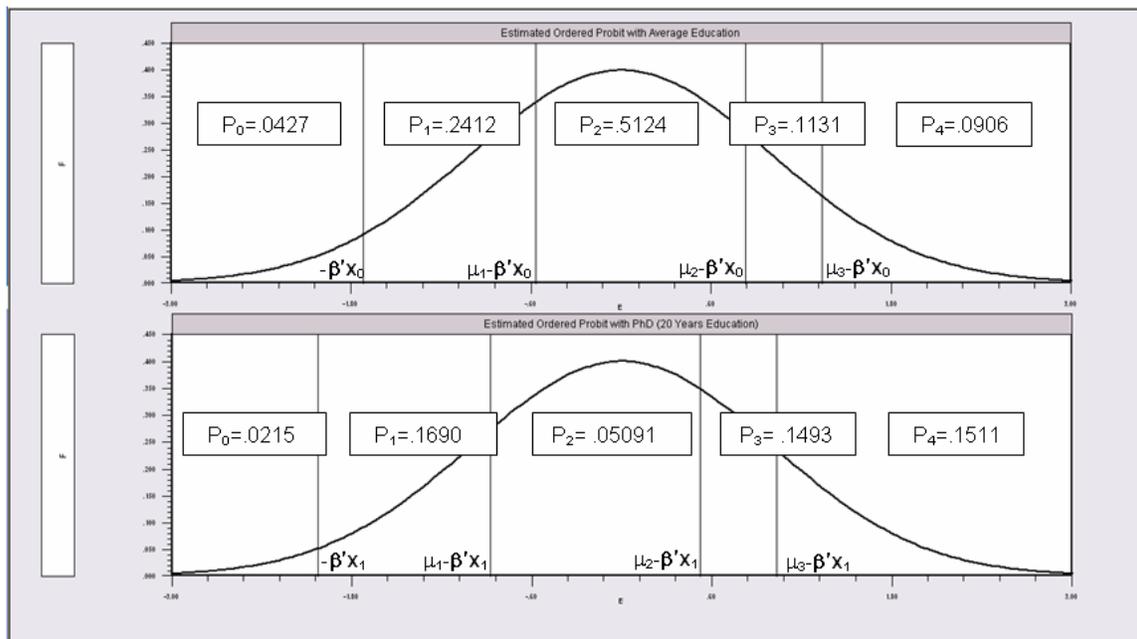


Figure 9 Partial Effect in Ordered Probit Model

Neither the signs nor the magnitudes of the coefficients are directly interpretable in the ordered choice model. It is necessary to compute partial effects of something similar to interpret the model meaningfully. In this computation, the only certainties in the signs of the partial effects in this model are as follows, where we consider a variable with a positive coefficient:

- Increases in that variable will increase the probability in the highest cell and decrease the probability in the lowest cell.
- The sum of all the changes will be zero. (The new probabilities must still sum to one.)
- The effects will begin at  $\text{Pr}(0)$  with one or more negative values, then change to a set of positive values; there will be one sign change. (This is the “single crossing” feature of the model. We will reconsider this aspect below.)

These are reversed for a variable with a negative coefficient.

One might also be interested in cumulative values of the partial effects, such as

$$\frac{\partial \text{Prob}(y \leq j | \mathbf{x}_i)}{\partial \mathbf{x}_i} = \sum_{m=0}^j [f(\mu_{m-1} - \boldsymbol{\beta}' \mathbf{x}_i) - f(\mu_m - \boldsymbol{\beta}' \mathbf{x}_i)] \boldsymbol{\beta}.$$

See, e.g., Brewer et al. (2008). (Note that the last term in this set is zero by construction.) An example appears below in Table 3.

Note in Table 2 there is a large difference in the coefficients obtained for the probit and logit models. The logit coefficients are roughly 1.8 times as large (not uniformly). This difference, which will always be observed, points up one of the risks in attempting to interpret directly the coefficients in the model. This difference reflects an inherent difference in the scaling of the underlying variable and in the shape of the distributions. The difference can be traced back (at least in part) to the different underlying variances in the two models. In the probit model,  $\sigma_\varepsilon = 1$ ; in the logit model  $\sigma_\varepsilon = \pi/\sqrt{3} = 1.81$ . The models are roughly preserving the ratio  $\boldsymbol{\beta}/\sigma_\varepsilon$  in the estimates. Note that the difference is greatly diminished (though not quite eliminated) in the partial effects reported in Table 3. That is the virtue of the scaling done to compute the partial effects. The inherent characteristics of the model are essentially the same for the two functional forms.

**Table 3 Estimated Partial Effects for Ordered Choice Models**

```

=====
|| Summary of Marginal Effects for Ordered Probability Model (PROBIT) ||
|| Effects computed at means. Effects for binary variables are ||
|| computed as differences of probabilities, other variables at means. ||
=====
||
|| Continuous Variable AGE
Outcome Effect dPy<=nn/dX dPy>=nn/dX
Y = 00 .00173 .00173 .00000
Y = 01 .00450 .00623 -.00173
Y = 02 -.00124 .00499 -.00623
Y = 03 -.00216 .00283 -.00499
Y = 04 -.00283 .00000 -.00283
=====
||
|| Continuous Variable EDUC Continuous Variable INCOME
Outcome Effect dPy<=nn/dX dPy>=nn/dX Effect dPy<=nn/dX dPy>=nn/dX
Y = 00 -.00340 -.00340 .00000 -.02476 -.02476 .00000
Y = 01 -.00885 -.01225 .00340 -.06438 -.08914 .02476
Y = 02 .00244 -.00982 .01225 .01774 -.07141 .08914
Y = 03 .00424 -.00557 .00982 .03085 -.04055 .07141
Y = 04 .00557 .00000 .00557 .04055 .00000 .04055
=====
||
|| Binary(0/1) Variable MARRIED Binary(0/1) Variable KIDS
Outcome Effect dPy<=nn/dX dPy>=nn/dX Effect dPy<=nn/dX dPy>=nn/dX
Y = 00 .00293 .00293 .00000 -.00574 -.00574 .00000
Y = 01 .00771 .01064 -.00293 -.01508 -.02081 .00574
Y = 02 -.00202 .00861 -.01064 .00397 -.01684 .02081
Y = 03 -.00370 .00491 -.00861 .00724 -.00960 .01684
Y = 04 -.00491 .00000 -.00491 .00960 .00000 .00960
=====
|| Summary of Marginal Effects for Ordered Probability Model (LOGIT) ||
=====
||
|| Continuous Variable AGE
Outcome Effect dPy<=nn/dX dPy>=nn/dX
Y = 00 .00145 .00145 .00000
Y = 01 .00521 .00666 -.00145
Y = 02 -.00166 .00500 -.00666
Y = 03 -.00250 .00250 -.00500
Y = 04 -.00250 .00000 -.00250
=====
||
|| Continuous Variable EDUC Continuous Variable INCOME
Outcome Effect dPy<=nn/dX dPy>=nn/dX Effect dPy<=nn/dX dPy>=nn/dX
Y = 00 -.00291 -.00291 .00000 -.01922 -.01922 .00000
Y = 01 -.01046 -.01337 .00291 -.06908 -.08830 .01922
Y = 02 .00333 -.01004 .01337 .02197 -.06632 .08830
Y = 03 .00502 -.00502 .01004 .03315 -.03318 .06632
Y = 04 .00502 .00000 .00502 .03318 .00000 .03318
=====
||
|| Binary(0/1) Variable MARRIED Binary(0/1) Variable KIDS
Outcome Effect dPy<=nn/dX dPy>=nn/dX Effect dPy<=nn/dX dPy>=nn/dX
Y = 00 .00287 .00287 .00000 -.00511 -.00511 .00000
Y = 01 .01041 .01327 -.00287 -.01852 -.02363 .00511
Y = 02 -.00313 .01014 -.01327 .00562 -.01801 .02363
Y = 03 -.00505 .00509 -.01014 .00897 -.00904 .01801
Y = 04 -.00509 .00000 -.00509 .00904 .00000 .00904
=====

```

#### 4.4.1 Nonlinearities in the Variables

In the computation of partial effects, it is assumed that the independent variables can vary independently. When the model contains interactions of variables, or nonlinear functions of variables, the computation of partial effects becomes problematic, though more so in practice than in theory. [See Norton and Ai (2003) for extensive analysis of this issue.] Consider, for example, in our model if we added variables  $EducSq$  and  $Educ*Age$ . The estimated model is shown in Table 4 with some of the partial effects. Separate partial effects are shown for  $Educ$ ,  $Age$ ,  $EducSq$  and  $EducAge$ , as if they were independent variables. In fact, in this model, the partial effect for education would be

$$\delta_j(Educ) = \frac{\partial \text{Prob}(y = j | \mathbf{x})}{\partial Educ} = [f(\mu_{j-1} - \beta' \mathbf{x}) - f(\mu_j - \beta' \mathbf{x})] (\beta_{Educ} + 2\beta_{EducSq} Educ + \beta_{EducAge} Age)$$

As Norton and Ai argued, none of the widely used computer packages computes this sort of result automatically. (It would be impossible for the software to anticipate every possible nonlinear function that might appear in the index function or recognize that function if it were implicit in a variable such as  $EducAge$ .) The analyst would have to compute this for themselves. This can be computed using the results reported, as

$$\delta_j(Educ) = \frac{\partial \text{Prob}(y = j | \mathbf{x})}{\partial "Educ"} + (2Educ) \frac{\partial \text{Prob}(y = j | \mathbf{x})}{\partial "EducSq"} + Age \frac{\partial \text{Prob}(y = j | \mathbf{x})}{\partial "EducAge"}$$

The derivatives shown for the top cell are 0.03326, -0.00093, -0.00007, respectively, and the means of  $Age$  and  $Education$  are 43.44 and 11.41, respectively. This, the partial effect is 0.0089966. In our original model with the linear index function, the estimated effect was 0.00557. Thus, the simple estimate of 0.03326 would be quite misleading, though, in fact, the partial effect does go up; but by 61%, not 600%!

#### 4.4.2 Average Partial Effects

In computing partial effects, we have evaluated the functions by inserting the sample means of the regressors. That is, our computation for  $Educ$ , for example, is

$$\frac{\partial \text{Prob}(y = j | \bar{\mathbf{x}})}{\partial Educ} = [f(\mu_{j-1} - \beta' \bar{\mathbf{x}}) - f(\mu_j - \beta' \bar{\mathbf{x}})] \beta_{Educ}$$

The average partial effect, or APE, is computed instead by evaluating the partial effect for each individual and averaging the computed effects. thus,

$$APE_j(Educ) = \frac{1}{N} \sum_{i=1}^N [f(\mu_{j-1} - \beta' \mathbf{x}_i) - f(\mu_j - \beta' \mathbf{x}_i)] \beta_{Educ}$$

In practice, unless the sample size is very small or the data are highly skewed and affected by outliers, this will give a very similar result. For the example suggested, the first computation gives 0.005557 and the second gives 0.005723, a difference of about 2.7%. Further discussion of the computation of APEs and standard errors using the delta method appear in Greene (2008a, pp. 783-785).



### 4.4.3 Interpreting the Threshold Parameters

In most treatments, the threshold parameters,  $\mu_j$  are treated as nuisance parameters; necessary for the computations, but of no intrinsic interest on their own. Daykin and Moffatt (2002,p. 162) argue that in psychology applications with attitude scales, “If the statement is one with which most people are either in strong agreement or strong disagreement, we would expect the cut points to be tightly bunched in the middle of the distribution. If, in contrast, the statement is one on which people are not keen to be seen expressing strong views, we would expect the cut points to be more widely dispersed.” Thus, in the absence of other information, this suggests that the threshold parameters can reveal some information about the preferences of the respondents. [In contradiction, Anderson (1984, p. 4) states “The estimates of the  $\theta_s$  are strongly related to the average proportion in the corresponding categories, as recourse to any specified functional form for  $F(\cdot)$  indicates. Hence, the  $\theta_s$  parameters are not informative about the closeness of categories. As noted above, the regression relationship is based on  $\beta'x$  and is firmly one dimensional.”]

### 4.4.4 The Underlying Regression

One would typically not be interested in the underlying regression. The observed variable will always be the discrete, ordered outcome. Nonetheless, the model does imply a set of partial changes for the latent regressand,

$$\partial E[y^*|x]/\partial x = \beta$$

This differs from more familiar cases in that the scaling of the dependent variable has been lost due to the censoring. Thus, it is impossible to attach any meaning to the change in the mean. McElvey and Zavoina (1975) suggest that if one is going to base interpretation of the model on the latent regression, then the coefficients should be “standardized.” That is, changes should be measured in standard deviation units. A standardized regression coefficient for variable  $k$  would be

$$\beta_k^* = \beta[s_{kk}/s_{y^*}]$$

where  $s_{kk}$  is the standard deviation of the regressor of interest and  $s_{y^*}$  is the standard deviation of  $y^*$ . Measurement of  $s_{kk}$  is straightforward based on the observed data. For  $s_{y^*}$ , the authors suggest the computation be based on the implication of the regression;

$$y^* = \beta'x + \varepsilon$$

so

$$\text{Var}[y^*] = \beta' \Sigma_{xx} \beta + \sigma_\varepsilon^2.$$

The two components are easily computed using the observed data and the normalized value of  $\sigma_\varepsilon^2 = 1$  or  $\pi^2/3$ . For our ordered choice model, the estimate of  $s_{y^*}$  is 1.03156. The results of the computation are shown below.

Variable	$\beta$	$\beta^*$
Age	-.01808	-2.23279
Educ	.03556	.19325
Income	.25869	.00676
Married	-.03100	-.00560
Kids	.06065	.01385

Some caution is needed when interpreting these. The variable that is assumed to be changing is an underlying preference scale. The notion of a unit or standard deviation change in utility or feeling is a bit dubious. That is among the motivations for discrete choice analysis of this sort; it frees the analyst from having to attach units of measure to unmeasurable quantities while still enabling them to learn about important features of preferences.

## 4.5 Inference

This section considers hypothesis tests about model components.

### 4.5.1 Inference about Coefficients

The model has been fit by maximum likelihood. The assumptions underlying the regularity conditions for maximum likelihood estimation should be met, so inference can be based on conventional methods. Standard errors for the estimated coefficients are computed by inverting an estimator of the negative of the expected second derivatives of the log likelihood. This will either be based on the actual second derivatives

$$\begin{aligned} \mathbf{V}_H = Est.Asy.Var \begin{bmatrix} \hat{\boldsymbol{\beta}}_{MLE} \\ \hat{\boldsymbol{\mu}}_{MLE} \end{bmatrix} &= \left[ -\sum_{i=1}^N \frac{\partial^2 \log \Pr(y = y_i | \mathbf{x}_i, \hat{\boldsymbol{\beta}}_{MLE}, \hat{\boldsymbol{\mu}}_{MLE})}{\partial \begin{bmatrix} \hat{\boldsymbol{\beta}}_{MLE} \\ \hat{\boldsymbol{\mu}}_{MLE} \end{bmatrix} \partial \begin{bmatrix} \hat{\boldsymbol{\beta}}'_{MLE} & \hat{\boldsymbol{\mu}}'_{MLE} \end{bmatrix}} \right]^{-1} \\ &= \left[ -\sum_{i=1}^N \hat{\mathbf{H}}_i \right]^{-1} \\ &= \left[ -\hat{\mathbf{H}} \right]^{-1} \end{aligned}$$

or the sum of the outer products of the first derivatives (the BHHH or outer product of gradients, OPG, estimator),

$$\begin{aligned} \mathbf{V}_{OPG} &= \left[ \sum_{i=1}^N \left( \frac{\partial \log \Pr(y = y_i | \mathbf{x}_i, \hat{\boldsymbol{\beta}}_{MLE}, \hat{\boldsymbol{\mu}}_{MLE})}{\partial \begin{bmatrix} \hat{\boldsymbol{\beta}}_{MLE} \\ \hat{\boldsymbol{\mu}}_{MLE} \end{bmatrix}} \right) \left( \frac{\partial \log \Pr(y = y_i | \mathbf{x}_i, \hat{\boldsymbol{\beta}}_{MLE}, \hat{\boldsymbol{\mu}}_{MLE})}{\partial \begin{bmatrix} \hat{\boldsymbol{\beta}}_{MLE} \\ \hat{\boldsymbol{\mu}}_{MLE} \end{bmatrix}} \right)' \right]^{-1} \\ &= \left[ \sum_{i=1}^N \hat{\mathbf{g}}_i \hat{\mathbf{g}}_i' \right]^{-1} \\ &= \left[ \hat{\mathbf{G}}' \hat{\mathbf{G}} \right]^{-1} \end{aligned}$$

Generally, two procedures, the Wald test and the likelihood ratio test are used for testing hypotheses. A third, the LM test, is available, but rarely used because of the simplicity of the other two.



A test about more than one coefficient can be carried out using a Wald test. For a null hypothesis of the form

$$H_0: \mathbf{R}\boldsymbol{\beta} = \mathbf{q}$$

where  $\mathbf{R}$  is a matrix of coefficients in the linear restrictions and  $\mathbf{q}$  is a vector of constants, the statistic will be

$$W = (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{q})' [\mathbf{RVR}' ]^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{q})$$

where  $\mathbf{V}$  is the estimated asymptotic covariance matrix of the coefficients. The difficulty of this computation will vary from one program to another. Both *Stata* and *NLOGIT* have built in “*Wald*” commands that can be used to do the computation as well as matrix algebra routines that also allow the user to program the computation themselves. For example, the following tests the null hypothesis that the coefficients on *EducSq* and *EducAge* in our expanded model are simultaneously zero. As noted, the statistic is treated as chi squared statistic with degrees of freedom equal to the number of restrictions. In the results below, for example, we see that we would reject the hypothesis that both are zero, evidently because of the significance of the first one.

```
Ordered ; Lhs = Health
          ; Rhs = one,age,educ,income,married,kids,educsq,educage $
Wald     ; fn1 = b_educsq ; fn2 = b_educag $
```

```
+-----+
| WALD procedure. Estimates and standard errors |
| for nonlinear functions and joint test of    |
| nonlinear restrictions.                      |
| Wald Statistic          =          6.64372   |
| Prob. from Chi-squared[ 2] =          .03609 |
+-----+
+-----+-----+-----+-----+-----+
|Variable| Coefficient | Standard Error |b/St.Er.|P[|Z|>z]|
+-----+-----+-----+-----+-----+
|Fncn(1) |   -.00596** |   .00234587   | -2.541 | .0110 |
|Fncn(2) |   -.00042   |   .00065479   |  -.641 | .5213 |
+-----+-----+-----+-----+-----+
| Note: ***, **, * = Significance at 1%, 5%, 10% level. |
+-----+-----+-----+-----+-----+

```

The counterparts for this computation in *Stata* would be

```
. oprobit health age educ income married kids educsq educage
. test educsq educage
( 1) [health]educsq = 0
( 2) [health]educage = 0
      chi2( 2) =    6.644
      Prob > chi2 =    0.0361
```

The computation can be programmed directly using matrix algebra, e.g., with *NLOGIT* as

```
Matrix      ; b2=b(7:8);v22=varb(7:8,7:8) $
Matrix      ; list ; Wald = b2'<v22>b2 $
Matrix WALD      has 1 rows and 1 columns.
              1
              +-----+
1|          6.64372
```

and using the *Mata* package in *Stata* or *PROC MATRIX* in *SAS*. In any case, using the built in procedure has the advantage of producing the “*p*-value” for the statistic as well as the statistic itself.

The likelihood ratio test will usually be simpler than the Wald test if the hypothesis is more involved than the simple zero restrictions shown above, though it does require estimation of both the null (restricted) and alternative (unrestricted) models. The test statistic is simply twice the difference between the log likelihoods for the null and alternative models. For the earlier example, the log likelihood for the (alternative) model that includes *EducSq* and *EducAge* is -5749.664 while, as seen earlier, the log likelihood for the (null) model that omits these variables is -5752.985. The test statistic is

$$LR = 2(-5749.664 - (-5752.985)) = 6.642.$$

This is nearly the same as the Wald statistic and produces the same conclusion. The two tests can conflict for a particular significance level. This is a finite sample result – asymptotically, the two statistics have the same characteristics when the assumptions of the model are met. As a general occurrence (albeit not necessarily), the Wald statistic will usually be larger than the *LR* statistic. Purely heuristically, because it uses more information – it is based on both models – we prefer the *LR* statistic.

A common test of the sort considered here is a “test of the model” in the spirit of the overall *F* statistic in the linear regression model that is used to test the null hypothesis that all coefficients in the model save the constant term are zero. The counterpart for the ordered choice model would be likelihood ratio test against the null hypothesis that the model contains only a constant term and the threshold parameters. This test statistic is routinely reported with the standard results for the estimated model by all commercial packages. For the preceding, we have

```

+-----+
| Ordered Probability Model
| Underlying probabilities based on Normal
| Maximum Likelihood Estimates
| Dependent variable           HEALTH
| Number of observations       4483
| Number of parameters         9
| Log likelihood function      -5752.985
| Restricted log likelihood     -5875.096
| Chi squared                  244.2238
| Degrees of freedom           5
| Prob[ChiSqd > value] =      .0000000
+-----+

```

Note it is not necessary to estimate the null model to carry out this test. The maximum likelihood estimates of the parameters of the model when it contains only a constant term are equivalent to method of moments estimators based on the following moment equations involving the raw sample proportions:

$$\begin{aligned}
 P_0 &= \Pr(y = 0) = F(-\alpha) \\
 P_1 &= \Pr(y \leq 1) = F(\mu_1 - \alpha) \\
 P_j &= \Pr(y \leq j) = F(\mu_j - \alpha) \\
 &\text{and so on.}
 \end{aligned}$$

These can be solved directly, in the logit case using a hand calculator (e.g.,  $a = \log(P_0/(1-P_0))$ ). These (with  $\beta = 0$ ) are the usual starting values for the iterations, so the log likelihood computed at entry to the iterative procedure provides the needed value for the null model.

## 4.5.2 Testing for Structural Change or Homogeneity of Strata

The likelihood ratio test provides a more convenient approach for testing homogeneity of strata in the data. For example, our data are separated by men and women in the introduction, and one might be interested in testing whether the same model should be used to describe the two groups. The counterpart to a “Chow test” in linear regression would be a test of group homogeneity in the choice model. The test statistic is easily computed using

$$LR = 2[\sum_{g=groups} \log L_g - \log L_{pooled}].$$

The statistic has a limiting chi squared distribution with degrees of freedom equal to  $G-1$  times the number of parameters in the model (slopes and thresholds). Our data are segmented by gender in the introduction. For a test of the null hypothesis that the same ordered choice model applies to the two groups, we find  $\log L_{Male} = -29.52.05$ ,  $\log L_{Female} = -2798.03$  and  $\log L_{Pooled} = -5752.98$ . Applying the preceding result gives a chi squared value of 5.83 with 9 degrees of freedom. The  $p$ -value is 0.7569 (the 95% critical value is 16.92). On this basis we conclude that is appropriate to pool these two subsamples. (In RWM’s analysis, they maintained the sample division.)

## 4.5.3 Robust Covariance Matrix Estimation

As noted earlier, there are two candidates available for the estimated asymptotic covariance matrix of the parameter estimators,  $-\mathbf{H}^{-1}$  based on the Hessian and  $(\mathbf{G}'\mathbf{G})^{-1}$  based on the first derivatives. The implication of the *Information Matrix Equality* [see Greene (2008a, Ch. 16)] is that these two matrices estimate the same covariance matrix and are, for practical purposes, interchangeable. A third matrix, the “robust” covariance matrix is often computed in recent applications, that being

$$\mathbf{V}_R = [-\mathbf{H}^{-1}] (\mathbf{G}'\mathbf{G}) [-\mathbf{H}^{-1}].$$

The logic of the computation can be seen by assuming that Netwon’s method is used to estimate the parameters. The maximum likelihood estimator at the maximum will produce

$$\sqrt{N}(\hat{\boldsymbol{\theta}}_{MLE} - \boldsymbol{\theta}^0) = \left[ -\hat{\mathbf{H}} / N \right]^{-1} \left( \sqrt{N} \sum_{i=1}^N \mathbf{g}_i \right) + o(1)$$

where  $\boldsymbol{\theta}^0$  is the vector of parameters that the MLE converges to and  $o(1)$  denotes a trailing term that converges to zero as  $N \rightarrow \infty$ . The asymptotic variance of the MLE is obtained by multiplying the limiting variance of the right hand side by  $1/N$ . The trailing terms will disappear. The leading matrix in brackets converges (we assume) to its expectation – a constant matrix. For the vector in parentheses, if the model assumptions are correct, then by the information matrix equality, its limiting variance will be  $-\mathbf{H}/N$ . Two occurrences of  $\mathbf{H}$  will cancel and we are left with  $\mathbf{V}_H$  as the usual estimator. But, ignoring the information matrix equality, whether it is met or not, the asymptotic variance of the MLE will be estimable by using  $(1/N)\mathbf{G}'\mathbf{G}$  as an estimator of the variance matrix of the quantity in parentheses. Then, the “robust” covariance matrix estimator becomes the sandwich estimator given above.

This produces two cases: If the model assumptions are correct, then the robust estimator is the same as either of the conventional estimators. If the model assumptions are incorrect, then the robust estimator still produces the asymptotic covariance matrix for the MLE. (A familiar application of this result is the “White” (1980) estimator for the asymptotic covariance matrix of

the OLS estimator in the presence of heteroscedasticity.) But, a new question arises in the second case. If the model assumptions are not correct, then what is  $\theta^0$ ? In order for this computation to be useful, it must be the case that in spite of the failure of the model assumptions,  $\hat{\theta}_{MLE}$  must still be a consistent estimator of the parameters of interest, in the present case,  $(\beta', \mu')$ . Once again, the case of OLS in the presence of heteroscedasticity provides a useful benchmark. On the other hand, for the ordered probit model, any of the following will render the estimator of the parameters inconsistent: (i) omitted variables even if they are orthogonal to included variables, (ii) heteroscedasticity in  $\varepsilon$ , (iii) incorrect distributional assumption – e.g., using the logit model when the probit model is the correct one, (iv) correlation across observations, (v) endogeneity of any of the regressors, (vi) omission of latent heterogeneity – this is equivalent to an omitted variable. Indeed, it is difficult to produce a model failure that the estimator is robust to. The upshot is that either the “robust covariance matrix” estimator is the same as the other two already considered, or it is a “robust” covariance matrix for an inconsistent estimator of the parameters. [Additional commentary on this result appears in Freedman (2006).]

#### 4.5.4 Inference About Partial Effects

Partial effects are computed using either the derivatives or first differences for discrete variables;

$$\delta_j(\mathbf{x}_i) = \frac{\partial \text{Prob}(y = j | \mathbf{x}_i)}{\partial \mathbf{x}_i} = [f(\mu_{j-1} - \beta' \mathbf{x}_i) - f(\mu_j - \beta' \mathbf{x}_i)] \beta$$

$$\Delta_j(d, \mathbf{x}_i) = [F(\mu_j - \beta' \mathbf{x}_i + \gamma) - F(\mu_{j-1} - \beta' \mathbf{x}_i + \gamma)] - [F(\mu_j - \beta' \mathbf{x}_i) - F(\mu_{j-1} - \beta' \mathbf{x}_i)].$$

Since these are functions of the estimated parameters, they are subject to sampling variability and one might desire to obtain appropriate asymptotic covariance matrices and/or confidence intervals. For this purpose, the partial effects are typically computed at the sample means. [See Greene (2008a, pp. 780-785) for analysis of this computation for average partial effects.] The delta method is used to obtain the standard errors. Let  $\mathbf{V}$  denote the estimated asymptotic covariance matrix for the  $(K+J-2) \times 1$  parameter vector  $(\hat{\beta}', \hat{\mu}')$ . Then, the estimator of the asymptotic covariance matrix for each vector of partial effects is

$$\mathbf{Q} = \hat{\mathbf{C}} \mathbf{V} \hat{\mathbf{C}}'$$

where

$$\hat{\mathbf{C}} = \begin{bmatrix} \frac{\partial \hat{\delta}_j(\bar{\mathbf{x}})}{\partial \hat{\beta}'} & \frac{\partial \hat{\delta}_j(\bar{\mathbf{x}})}{\partial \hat{\mu}'} \end{bmatrix}.$$

The appropriate row of  $\hat{\mathbf{C}}$  is replaced with the derivatives of  $\Delta_j(d, \bar{\mathbf{x}})$  when the effect is being computed for a discrete variable.

Patterns of statistical significance for the partial effects will usually echo those for the coefficients themselves. This will follow from the fact that  $\mathbf{C}$  is of the form

$$\mathbf{C} = [a_{ij} \mathbf{I}, \mathbf{0}] + [\mathbf{C}_{\beta,2}, \mathbf{C}_{\mu}]$$

where  $a_{ij}$  is the bracketed scalar term in  $\hat{\delta}_j(\bar{\mathbf{x}})$ . The second matrix is typically much smaller than the first. Thus, the estimated asymptotic covariance matrix for  $\hat{\delta}_j(\bar{\mathbf{x}}) = a_{ij}\beta$  typically resembles  $a_{ij}^2\mathbf{V}$ . The scale factor would cancel out of a “z value” leaving the typical result. It is clearly visible in the results in Table 6 below. This result does raise a vexing question. It is conceivable for the significance tests of  $\delta_j(x_k)$  to conflict with each other, that is, with  $\delta_m(x_k)$  for an  $m \neq j$ , and/or with a test about the associated coefficient,  $\beta_k$ . Since  $\delta_j(x_k) = a_{ij}\beta_k$ , the tests would seem to be in direct contradiction. The natural question for the practitioner, then, is where should the appropriate test of significance be carried out. Opinions differ and there is no single answer. It might logically be argued that the overall purpose of the regression analysis is to compute the partial effects, so that is where the tests should be carried out. On the other hand, the meaning of the test with respect to the partial effects is ambiguous, since they are functions of all the parameters as well as the data. The number of possible contradictions is large. Our preference on the methodological basis is for the structural coefficients, not the partial effects.

#### 4.6 Prediction – Computing Probabilities

One might want to use the model for prediction as well as inference. The natural predictor would seem to be  $\hat{y}^* = \hat{\beta}'\mathbf{x}$ . However, the underlying variable is typically unobservable, and often of no intrinsic interest in its own right. (E.g., in the bioassay case, the “tolerance” of a particular insect would probably be of little interest. In the preference scale case such as in our health satisfaction example, the underlying utility is inherently unmeasurable.) The more natural exercise would be to predict the observed outcome. Since it is discrete, the linear predictor is of little use. The starting point would be the predicted probabilities. The model provides predictors

$$\begin{aligned}\hat{P}_j(x_i) &= F(\hat{\mu}_j - \hat{\beta}'\mathbf{x}_i) - F(\hat{\mu}_{j-1} - \hat{\beta}'\mathbf{x}_i) \\ &= \hat{F}_{j,i} - \hat{F}_{j-1,i}, j = 0, 1, \dots, J\end{aligned}$$

If the sample is small enough and particular observations are of interest, a simple listing might be useful. For our sample of 4,483 observations, this would probably not be helpful. One might, instead, tabulate predicted probabilities against variables of interest. For example, for reasons unknown to us, the presence of children in the household appears to have a substantial (increasing) impact on whether one reports the lowest value of health satisfaction.

Standard errors and confidence intervals can be computed using the delta method. These are a bit simpler than for the partial effects, as there is no need to make a distinction between discrete and continuous variables. The matrix of derivatives has a row for each outcome, containing

$$\frac{\partial \hat{P}_j(\mathbf{x}_i)}{\partial (\hat{\beta}' \quad \hat{\mu}')} = \left[ (\hat{f}_{j-1}(\mathbf{x}_i) - \hat{f}_j(\mathbf{x}_i))\mathbf{x}_i' \quad (0, \dots, -\hat{f}_{j-1}, \hat{f}_j, 0, \dots) \right]$$

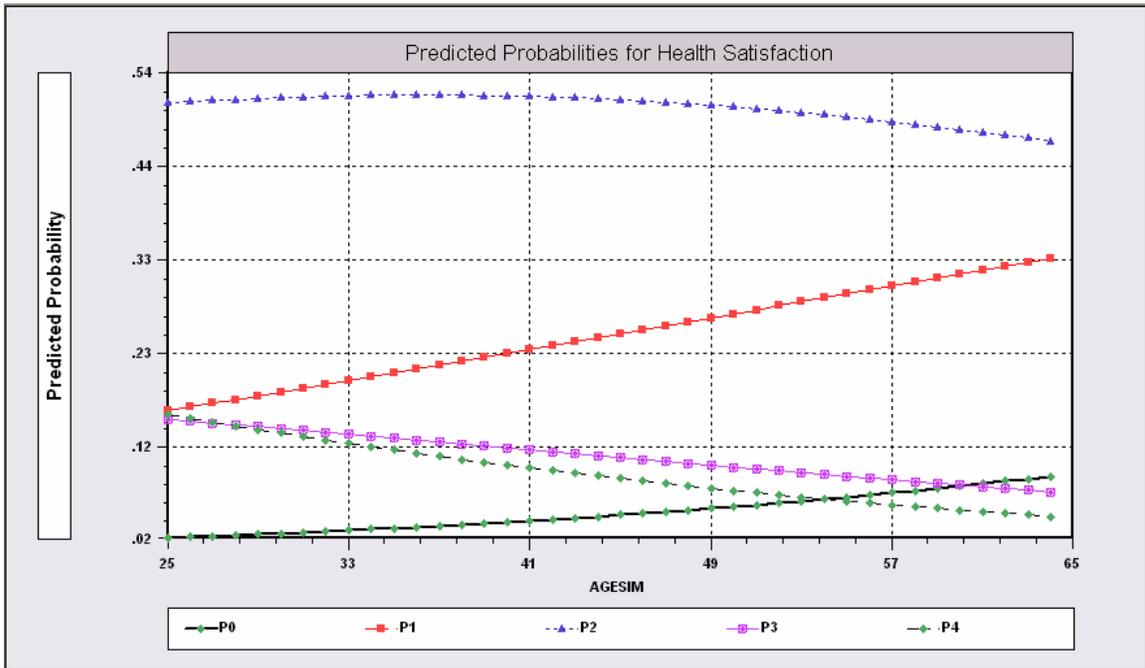
For certain variables of interest, a plot of the predicted probabilities against the values of the variable might be useful. In our application, *Age* seems to be an important determinant of self assessed health satisfaction. A plot of the predicted probabilities for this model for the values of *Age* in the sample, 25 to 64, for a person who has average income and education, and is married with children appears in Figure 10.

**Table 6 Estimated Partial Effects with Asymptotic Standard Errors**

Marginal effects for ordered probability model					
M.E.s for dummy variables are $\Pr[y x=1]-\Pr[y x=0]$					
Names for dummy variables are marked by *.					
Variable	Coefficient	Standard Error	b/St.Er.	P[ Z >z]	Mean of X
These are the effects on Prob[Y=00] at means.					
AGE	.00173***	.00016500	10.488	.0000	43.440107
EDUC	-.00340***	.00069211	-4.919	.0000	11.418086
INCOME	-.02476**	.00997292	-2.483	.0130	.3487401
*MARRIED	.00293	.00392048	.747	.4551	.7521749
*KIDS	-.00574	.00357811	-1.603	.1089	.3794334
Variable	Coefficient	Standard Error	b/St.Er.	P[ Z >z]	Mean of X
These are the effects on Prob[Y=01] at means.					
AGE	.00450***	.00040305	11.161	.0000	43.440107
EDUC	-.00885***	.00177514	-4.986	.0000	11.418086
INCOME	-.06438**	.02585078	-2.490	.0128	.3487401
*MARRIED	.00771	.01044010	.738	.4604	.7521749
*KIDS	-.01508	.00949339	-1.588	.1122	.3794334
Variable	Coefficient	Standard Error	b/St.Er.	P[ Z >z]	Mean of X
These are the effects on Prob[Y=02] at means.					
AGE	-.00124***	.00016956	-7.310	.0000	43.440107
EDUC	.00244***	.00054946	4.438	.0000	11.418086
INCOME	.01774**	.00735553	2.411	.0159	.3487401
*MARRIED	-.00202	.00261143	-.774	.4387	.7521749
*KIDS	.00397	.00241917	1.641	.1009	.3794334
Variable	Coefficient	Standard Error	b/St.Er.	P[ Z >z]	Mean of X
These are the effects on Prob[Y=03] at means.					
AGE	-.00216***	.00024067	-8.958	.0000	43.440107
EDUC	.00424***	.00090065	4.709	.0000	11.418086
INCOME	.03085**	.01255878	2.457	.0140	.3487401
*MARRIED	-.00370	.00503280	-.736	.4620	.7521749
*KIDS	.00724	.00459950	1.574	.1154	.3794334
Variable	Coefficient	Standard Error	b/St.Er.	P[ Z >z]	Mean of X
These are the effects on Prob[Y=04] at means.					
AGE	-.00283***	.00027111	-10.452	.0000	43.440107
EDUC	.00557***	.00113041	4.931	.0000	11.418086
INCOME	.04055**	.01633487	2.482	.0130	.3487401
*MARRIED	-.00491	.00673253	-.729	.4657	.7521749
*KIDS	.00960	.00612040	1.569	.1166	.3794334
Note: ***, **, * = Significance at 1%, 5%, 10% level.					

**Table 7 Mean Predicted Probabilities by Kids**

Variable	Mean	Std.Dev.	Minimum	Maximum	Cases Missing	
=====						
Stratum is KIDS	=	.000.	Obs.= 2782.000,	Sum of wts. =	2782.000	
-----						
P0	.595860E-01	.281820E-01	.956145E-02	.125545	2782	1701
P1	.268398	.634147E-01	.106526	.374712	2782	1701
P2	.489603	.243695E-01	.419003	.515906	2782	1701
P3	.101163	.301566E-01	.525888E-01	.181065	2782	1701
P4	.812503E-01	.412504E-01	.281517E-01	.237842	2782	1701
Stratum is KIDS	=	1.000.	Obs.= 1701.000,	Sum of wts. =	1701.000	
-----						
P0	.363923E-01	.139256E-01	.109540E-01	.105794	1701	2782
P1	.217619	.396625E-01	.115439	.354036	1701	2782
P2	.509830	.904826E-02	.443130	.515906	1701	2782
P3	.125049	.194545E-01	.616726E-01	.176725	1701	2782
P4	.111111	.304129E-01	.353675E-01	.222307	1701	2782
All observations in current sample						
-----						
P0	.507855E-01	.263258E-01	.956145E-02	.125545	4483	0
P1	.249130	.608208E-01	.106526	.374712	4483	0
P2	.497278	.222687E-01	.419003	.515906	4483	0
P3	.110226	.290206E-01	.525888E-01	.181065	4483	0
P4	.925804E-01	.402074E-01	.281517E-01	.237842	4483	0



**Figure 10 Predicted Probabilities for Different Ages**

## 4.7 Measuring Fit

The search for a scalar measure of model fit for discrete choice models must be among the least satisfying of the exercises in the modeling effort. Superficially, the search is for a counterpart to the  $R^2$  = “proportion of the variation in the dependent variable that is explained by variation in the independent variables.” The search is frustrated in this (and other discrete choice models) for two reasons:

- There is no “dependent variable.” In the ordered choice model, there are  $J+1$  explained variables that are defined by  $m_{ij} = 1$  if  $y_i = j$  and 0 otherwise and which satisfy the constraints  $m_{ij} = 0$  or 1 and  $\sum_j m_{ij} = 1$ . (This is true for the bioassay case as well; the observed proportions for each  $i$  consist of the sample means of  $m_{ij}$  for  $n_i$  observations with a common  $\mathbf{x}_i$ .) The observed variable  $y_i$  is nothing more than a labeling convention for the regions of the real line defined by the partitioning in the model specification.
- There is no “variation” (around the mean) to be explained. The outcome is not a measure of a quantity; it is a label. There is no conditional mean, as such, either.

For these reasons, one needs to exert a considerable amount of caution in computing and reporting “measures of fit” in this setting.

A “fit measure” that one computes can be used for two purposes: (i) to assess the fit of the predictions by the model to the observed data, compared to no model and (ii) to compare the model one estimates to a different model. For the first of these, we (and a generation of others) have suggested the overall model chi squared,

$$\chi^2[K+J-2] = 2[\log L_{Model} - \log L_{No Model}].$$

A transformation of this statistic that is (very) often reported in the contemporary literature is McFadden’s (1977) “pseudo  $R^2$ ” which is computed as

$$R_{Pseudo}^2 = 1 - \log L_{Model} / \log L_{No Model}.$$

A degrees of freedom adjusted version is sometimes reported,

$$Adjusted R_{Pseudo}^2 = 1 - [\log L_{No Model} - M] / \log L_{Model}.$$

where  $M$  is the number of parameters in the model. This fit measure has the virtues that it is bounded by 0 and 1, and increases whenever the model increases in size – that is, the pseudo  $R^2$  is larger for any model compared to a model that is nested within it. It is important to emphasize, as is clear from the definition, it is not a measure of model fit to the data and it is not a measure of the proportion of variation explained in any sense. (It is also worth noting that it is not necessarily bounded by zero and one unless the model in question is a discrete choice model for which the log likelihood function is necessarily negative. For example, it is a simple exercise to show that the log likelihood for a linear normal regression model can be positive or negative, depending on the value of  $\sigma_\varepsilon$ , which could produce values outside the unit interval.) Lastly, the *Pseudo  $R^2$*  cannot reach one, though it can equal zero.

The value of the *Pseudo  $R^2$*  in the model we have analyzed above can be found in Table 2 for the basic model (0.02075) and in Table 4 for the expanded model (0.02134). The low values

might seem a bit surprising given the several highly significant coefficient estimates in the reported results. However, as with the counterpart in linear regression, highly significant coefficients need not attend a high fit measure.

A second commonly reported measure for the ordered choice model was suggested by McKelvey and Zavoina (1975). The logic of their measure is based on predicting the underlying latent variable,  $y^*$ . As noted in Section 4.4.4, the total variance *in the underlying variable* in the choice model is

$$\text{Var}[y^*] = \boldsymbol{\beta}' \boldsymbol{\Sigma}_{xx} \boldsymbol{\beta} + \sigma_\varepsilon^2.$$

where  $\boldsymbol{\Sigma}_{xx}$  is the theoretical covariance matrix of  $\mathbf{x}_i$ . The first part of this is estimable using the maximum likelihood estimates of  $\boldsymbol{\beta}$  and the sample covariance matrix for the data, and the second part is known to be one or  $\pi^2/3$  for the probit and logit models, respectively. Thus, the authors suggested

$$R_{MZ}^2 = 1 - \frac{\sigma_\varepsilon^2}{\hat{\boldsymbol{\beta}}' \mathbf{S}_{xx} \hat{\boldsymbol{\beta}} + \sigma_\varepsilon^2}$$

They defined the “explained” part of this computation in terms of deviations from a prediction,  $e_i = \hat{y}_i - \hat{\bar{y}}$  where  $\hat{y}_i = \hat{\boldsymbol{\beta}}' \mathbf{x}_i$ , producing

$$R_{MZ}^2 = \frac{\sum_{i=1}^N (\hat{y}_i - \hat{\bar{y}})^2}{\sum_{i=1}^N (\hat{y}_i - \hat{\bar{y}})^2 + N}$$

With this computation, we obtain an improvement over the *PseudoR*<sup>2</sup>; for our model,  $R_{MZ}^2 = 0.06024$ .

Long and Freese (2006) list a variety of other measures that are computed for the ordered choice models. (We note, this set of results is produced by a *Stata* program called `FitStat` written by one of the authors. We mention it at this juncture to illustrate the problem of searching for a fit measure in a particular discrete choice model, not to recommend that analysts either do or do not use it or these results. The formulas below do not appear in Long and Freese or in the documentation for *Stata*; they are described in long detail by UCLA/ATS (2008) among others and, of course, piecemeal by the original designers.) These include

$$R_{Cox,Snell}^2 = 1 - \left[ \frac{\log L_{No\ Model}}{\log L_{Model}} \right]^{2/N}$$

$$R_{Cragg,Uhler/Nagelkerke}^2 = \frac{1 - \left[ \frac{\log L_{No\ Model}}{\log L_{Model}} \right]^{2/N}}{1 - [\log L_{No\ Model}]^{2/N}}$$

In UCLA/ATS (2008), it is noted that “pseudo *R*-squareds” for categorical variables serve three functions:

- Measures of explained variability,
- Measures of improvement from null model to fitted model,
- Square of the correlation.

None of the already suggested fit measures bear any relation to the first and third of these. All are connected to the improvement in the log likelihood by the addition of the variables in the model to a constants only model. Of course, the log likelihood functions, themselves, do that, and what these statistics add to the two values is a transformation that is between zero and one. It is worth noting, the measures are strictly between zero and one. None can achieve one even if the model predicts perfectly (somehow – we have not defined what would be meant by “predict”). Nonetheless, what they do all share is that they increase as the model grows and they are bounded by zero and one. (However, the “adjusted pseudo  $R^2$ ” can decline as variables are added, in the same fashion as  $\bar{R}^2$  for linear regression.)

UCLA/ATS (2008) observe (with reference to a binary logit model),

When analyzing data with a logistic regression, an *equivalent statistic to R-squared does not exist*. [Emphasis added.] The model estimates from a logistic regression are maximum likelihood estimates arrived at through an iterative process. They are not calculated to minimize variance, so the OLS approach to goodness-of-fit does not apply. However, to evaluate the goodness-of-fit of logistic models, several pseudo  $R$ -squareds have been developed. These are “pseudo”  $R$ -squareds because they look like  $R$ -squared in the sense that they are on a similar scale, ranging from 0 to 1 (though some pseudo  $R$ -squareds never achieve 0 or 1) with higher values indicating better model fit, but they cannot be interpreted as one would interpret an OLS  $R$ -squared and different pseudo  $R$ -squareds can arrive at very different values

We note, the notion of “model fit” in this and elsewhere relates to the log likelihood for the model, not to an assessment of how well the model predicts the outcome variable, as it does in regression analysis.

It seems appropriate to add a fourth item to the list above; fit measures are used to compare models to each other, not only to baseline, “null” models. For this purpose, a handful of other fit measures that are not normalized to the unit interval, but are based on the log likelihood function, are often used:

$$\textit{Akaike Information Criterion} = \textit{AIC} = [-2\log L + 2M]/N$$

$$\textit{Finite Sample AIC} = \textit{AIC}_{FS} = \textit{AIC} + 2M(M+1)/(N - M - 1)$$

$$\textit{Bayes Information Criterion} = \textit{BIC} = [-2\log L + M/\log N]/N$$

$$\textit{Hannan-Quinn IC} = \textit{HQIC} = [-2\log L + 2 M \log \log N]/N$$

The information measures are all created in the spirit of adjusted  $R^2$  – they reward a model for “fit” with few parameters and small samples. A better model is one with a smaller information criterion. (Long and Freese mention two others, “*AIC used by Stata*” and “*BIC used by Stata*” that we have been unable to decipher.)

Long and Freese (p. 196) and UCAL/ATS (2008) mention two other measures that seem (to these authors) to have received far less attention than these likelihood based measures. These are

$$\textit{Count } R^2 = \frac{\textit{Number of Correct Predictions}}{N}$$

and

$$\text{Adjusted Count } R^2 = \frac{\text{Number of Correct Predictions} - N_j^*}{N - N_j^*}.$$

Where  $N_j^*$  is the count of the most frequent outcome. The discussion is about binary choice models, so we have to extend the idea to our ordered choice model. There is a long catalog of fit measures for binary choice models based on this sort of computation. [See, e.g., Greene (2008a, pp. 790-793).] The central feature is a fitting mechanism: Predict  $y = j$  if the model states that  $j$  is the most likely outcome. In the binary choice case, the rule is to use as the prediction, the outcome which has probability exceeding 0.5. For the ordered choice case, this would suggest using the rule

$$\hat{y}_i = j^* \text{ such that estimated } \Pr(y_i = j^* | \mathbf{x}_i) > \Pr(y_i = j | \mathbf{x}_i) \forall j \neq j^*$$

That is, put the predicted  $y$  in the cell with the highest probability. This rule has an aesthetic appeal, and in the absence of priors (as in a Bayesian setting) we have not found a preferable approach. Nonetheless, this can lead to an unexpected outcome. For our example, this rule produces the following table.

```

+-----+
|      Cross tabulation of predictions. Row is actual, column is predicted.      |
|      Model = Probit      .      Prediction is number of the most probable cell.  |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| Actual|Row Sum| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|      0|    230| 0 | 0 | 230| 0 | 0 |   |   |   |   |   |
|      1|   1113| 0 | 0 | 1113| 0 | 0 |   |   |   |   |   |
|      2|   2226| 0 | 0 | 2226| 0 | 0 |   |   |   |   |   |
|      3|    500| 0 | 0 | 500| 0 | 0 |   |   |   |   |   |
|      4|    414| 0 | 0 | 414| 0 | 0 |   |   |   |   |   |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| Col Sum|  4483| 0 | 0 | 4483| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+

```

By this method, our model, with its highly significant overall fit and several highly significant variables seems, nonetheless, to fail utterly on this criterion. It always predicts  $y = 2$ . By the *Count*  $R^2$  measure, our model achieves a fit of 0.4965, which is looks like a substantial improvement over the *Pseudo*  $R^2$  of 0.02075. Lest we become too enthusiastic about the result, however, note that the *Adjusted Count*  $R^2$  is zero! The reason is that the model does not improve on the model free “always predict 2.”

The situation in which the model always predicts the same value is not uncommon. It takes a high correlation (in some general sense) between the covariates and the outcome and a large amount of variation in the covariates within the sample to spread the predictions across the outcomes. Briefly, another example is provided by a standard data set used by the authors of *Stata* to demonstrate the ordered choice model in their documentation. The “automobile data,” (<http://www.stata-press.com/data/r8/fullauto.dta>) is used in [R] `oprobit` to model the 1977 repair records of 66 foreign and domestic cars. The variable *rep77* takes values *poor*, *fair*, *average*, *good* and *excellent*. The explanatory variables in the model are *foreign* (origin of manufacture), *length* (a proxy for size) and *mpg*. (The computations below were obtained with both *Stata* and *NLOGIT*, which obtained identical results.) The predictions produced by this model are listed below. The McFadden *Pseudo*  $R^2$  is 0.1321. The *Count*  $R^2$  is  $(1+0+21+7+1)/66 = 0.454$ . The adjusted value is  $(30 - 27)/(66-27) = 0.077$ .

Cross tabulation of predictions. Row is actual, column is predicted. Model = Probit . Prediction is number of the most probable cell.											
Actual	Row Sum	0	1	2	3	4	5	6	7	8	9
0	3	1	0	2	0	0					
1	11	0	0	9	2	0					
2	27	0	1	21	5	0					
3	20	0	0	11	7	2					
4	5	0	0	2	2	1					
Col Sum	66	1	1	45	16	3	0	0	0	0	0

This survey does not conclude with a proposal for the appropriate or optimal fit measure. The search for a scalar counterpart to the  $R^2$  in a linear regression does seem unproductive. Fit measures based on the log likelihood can be used for comparing models. For this purpose, the log likelihood itself or one of the information criteria seems sensible; the AIC dominates the received applications. For assessing the predictions of the model, it would seem that the scalar measures based on the log likelihood would be useless. The maximum likelihood estimator is not computed so as to maximise the number of correction predictions – in the linear normal regression model, the MLE of  $\beta$  is computed to maximize  $R^2$ , but that is coincidental; minimizing  $e'e$  does maximize  $R^2$ . Indeed, there may be (as yet not proposed) other estimators that improve on the MLE for predicting the outcome variable, as the Maximum Score Estimator [see Manski (1975, 1985, 1986, 1988)] improves on the MLE of the logit or probit model for binary choice. In any event, it does seem appropriate, if one seeks a “measure of fit” one should first decide upon a procedure (rule) for producing the predictions, then assess, against a benchmark, how well that method does. The *Count*  $R^2$  measures shown above seem better suited to that specific purpose than pseudo  $R^2$  measures based on the log likelihood.

#### 4.8 Estimation Issues

McKelvey and Zavoina (1975) provide expressions for the first and second derivatives of the log likelihood function for the ordered probit model, and suggest Newton’s method as an algorithm for estimation. They do conjecture, however, about the possible problem of multiple roots of the log likelihood. Pratt (1981), was able to show that the ordered probit model was a member of a class of discrete choice models in which the log likelihood functions are globally concave. Thus, estimation of the model can be counted on to converge (when it does at all), to the single root of the log likelihood function. We note at this point a few other aspects of estimation of the ordered choice model.

### 4.8.1 Grouped Data

The adaptation of the maximum likelihood estimator to the grouped data (bioassay) treatment is a trivial modification. The log likelihood for a sample in which the stimulus,  $\mathbf{x}_i$  is repeated  $n_i$  times is

$$\begin{aligned} \log L &= \sum_{i=1}^N \sum_{t=1}^{n_i} \sum_{j=0}^J m_{it,j} \log [F(\mu_j - \beta' \mathbf{x}_i) - F(\mu_{j-1} - \beta' \mathbf{x}_i)] \\ &= \sum_{i=1}^N \sum_{j=0}^J \log [F(\mu_j - \beta' \mathbf{x}_i) - F(\mu_{j-1} - \beta' \mathbf{x}_i)] \sum_{t=1}^{n_i} m_{it,j} \\ &= \sum_{i=1}^N \sum_{j=0}^J n_{ij} \log [F(\mu_j - \beta' \mathbf{x}_i) - F(\mu_{j-1} - \beta' \mathbf{x}_i)] \\ &= \sum_{i=1}^N n_i \sum_{j=0}^J p_{ij} \log [F(\mu_j - \beta' \mathbf{x}_i) - F(\mu_{j-1} - \beta' \mathbf{x}_i)] \end{aligned}$$

Mechanically, in the log likelihood for a cross section of individual data, the terms  $m_{ij}$  are replaced with the group proportions,  $p_{ij}$ , and the observations in the log likelihood and its derivatives are weighted by the group size.

### 4.8.2 Perfect Prediction

A problem of nonconvergence can be caused by a condition in the data that Long and Freese (2006, p. 192) label “Predicting Perfectly.” If a variable in the data set predicts perfectly one of the implicit dependent variables,  $m_{ij} = 1$  if and only if  $y_i = j$ , then it will not be possible to fit the coefficients of the model – in this instance, the corresponding threshold parameter becomes inestimable. The suggested case is a dummy variable that takes only one value within a particular cell – it may also take that value in other cells. Within our example, suppose married people (*Married* = 1) always responded with *Health* = 4; i.e., married people always report the highest health satisfaction. Then, knowing someone is married allows a perfect prediction of *Health* = 4 for them. In such as case, it is necessary to drop such observations from the sample. *Stata* detects this condition automatically and reports a diagnostic “Note: nn observations completely determined. Standard errors are questionable.” As it is, the diagnostic is correct. But, it is incomplete. Because the offending variable enjoys such a relationship with the outcome variable, it is almost certainly endogenous in the model, and not only are the standard errors questionable, the parameter estimates themselves are as well. In a vague way, this is a cousin to a problem of sample selection. The observations that have been discarded have not been done so randomly. They have been discarded by a criterion that is specifically related to the dependent variable. This particular feature of the model is as of this writing an obscure corner of the model development, but there would seem to be scope for further analysis of the issue.

It is tempting in this instance just to drop the offending variable. Whether this is advisable or not is unclear. If one is certain that but for the (perhaps unexpected) data problem the variable is an important feature of the data generating process, then the resulting model when the variable is dropped now has an omitted regressor. One problem has been traded for another. On the other hand, if the problem considered here involves more than just a handful of observations, one might question the overall structure of the model. Treating such a variable as if it were exogenous might be inappropriate.

### 4.8.3 Different Normalizations

We have noted at a few points that the normalization of the thresholds is a crucial feature of the model. However, it is not the case that different normalizations produce different results. Whether one assumes  $\mu_0 = 0$  and includes an overall constant in the model, or allows  $\mu_0$  to be a free parameter and drops the constant, will have no implications for the log likelihood, the other parameters, or the predictions of the model. An example to illustrate the point is useful. Consider, once again, the car repair data discussed in the previous section. We have fit the model using *NLOGIT*, which uses the first normalization and *Stata* which uses the second. The two sets of results are given in Table 8. Note that the log likelihoods and estimates of the coefficients in  $\beta$  are identical. (The differences in the standard errors result from *Stata*'s use of the Hessian for the standard errors vs. *NLOGIT*'s use of the outer products estimator.) The first "cut point" in the *Stata* results is precisely the negative of *NLOGIT*'s overall constant. For the remaining threshold parameters, we can see that "cut point  $j$ " equals *NLOGIT*'s  $(\mu_j - \alpha)$ . As expected, then, the results are identical.

### 4.8.4 Censoring of the Dependent Variable

In some applications, there can be a second layer of censoring of the variable of interest in the ordered choice model. (The first level of censoring is the translation of  $y_i^*$  to  $y_i$  by measuring only the interval in which  $y_i^*$  appears.) Consider a model of educational attainment in which the variable of interest is "education" and in which the recorded value is only 0 for primary school, 1 for secondary school (high school), 2 for college, 3 for masters and 4 for Ph.D. If an observation is recorded as "at least high school," for example, then values 2, 3 and 4 are censored. This case is easily handled using the laws of probability. The appropriate log likelihood for the ordered choice model is

$$\log L = \sum_{i=1}^N \sum_{j=0}^J m_{ij} \log(P_{ij} - P_{i,j-1})$$

where heretofore  $m_{ij}$  indicated the one cell that applies to observation  $i$ , and now indicates all of the cells that apply. For the example given, we would have  $m_{i0} = 0$  and  $m_{ij} = 1$  for  $j = 1, 2, 3, 4$ . The change in the computations of the model parameters is trivial. It should be noted, one must know the upper bound,  $J$ , and for an observation, of course, it must be known that it is or is not censored. Censoring of the dependent variable in an ordered choice context has appeared in models of schooling attainment by Lillard and King (1987), Glewwe (1997) and Glewwe and Jacoby (1994, 1995) and in duration models, where the observed outcome is the length of time between transitions, sometimes coded as "short," medium or long, or similarly. See, e.g., Tsay (2005), Han and Hausman (1988) and Buckle and Carlson (2000).



#### 4.8.5 Maximum Likelihood Estimation of the Ordered Choice Model

The log likelihood function for the basic ordered choice model is

$$\begin{aligned}\log L &= \sum_{i=1}^N n_i \sum_{j=0}^J w_{ij} \log [F(\mu_j - \beta' \mathbf{x}_i) - F(\mu_{j-1} - \beta' \mathbf{x}_i)] \\ &= \sum_{i=1}^N n_i \sum_{j=0}^J w_{ij} \log (F_{i,j} - F_{i,j-1}) \\ &= \sum_{i=1}^N n_i \sum_{j=0}^J w_{ij} \log P_{i,j}\end{aligned}$$

where

$n_i$  = the group size in the grouped data (typical bioassay) case or  
 $n_i = 1$  in the individual data case,

and

$w_{ij} = p_{ij}$  = the proportion of group  $i$  that responds with outcome  $j$ , or  
 $w_{ij} = m_{ij} = 1$  if individual  $i$  chooses outcome  $j$  in the individual data case

$F(t)$  is the functional form in use, typically  $\Lambda(t)$  for the ordered logit model or  $\Phi(t)$  for the ordered probit model. For the moment, we will leave the functional form indeterminate. For obtaining the log likelihood and its derivatives, only the term  $\log P_{i,j}$  is of consequence. The relevant derivatives are

$$\begin{aligned}\frac{\partial \log P_{i,j}}{\partial \beta} &= \frac{f_{i,j} - f_{i,j-1}}{P_{i,j}} (-\mathbf{x}_i) \\ \frac{\partial \log P_{i,j}}{\partial \mu_j} &= \frac{f_{i,j}}{P_{i,j}} \\ \frac{\partial \log P_{i,j}}{\partial \mu_{j-1}} &= \frac{-f_{i,j-1}}{P_{i,j}}\end{aligned}$$

where  $f_{i,j}$  is the density corresponding to  $F_{i,j}$ . For the moment, we are carrying  $\mu_{-1}$ ,  $\mu_0$  and  $\mu_J$  as if they were unconstrained. The constraints are imposed later. Thus, the parameter vector contains  $\beta$  and  $\mu$ , which has  $J+2$  elements only  $J-1$  of which are free to vary. The derivative vector  $\partial \log P_{i,j} / \partial \mu$  has  $J+2$  elements, but only two are nonzero. The second derivatives are as follows:

$$\begin{aligned}\frac{\partial^2 \log P_{i,j}}{\partial \beta \partial \beta'} &= \left[ \left( \frac{f'_{i,j} - f'_{i,j-1}}{P_{i,j}} \right) - \left( \frac{f_{i,j} - f_{i,j-1}}{P_{i,j}} \right)^2 \right] \mathbf{x}_i \mathbf{x}_i' \\ \frac{\partial^2 \log P_{i,j}}{\partial \beta \partial \mu_j} &= \left[ \frac{f'_{i,j}}{P_{i,j}} - \frac{(f_{i,j} - f_{i,j-1}) f_{i,j}}{P_{i,j}^2} \right] (-\mathbf{x}_i) \\ \frac{\partial^2 \log P_{i,j}}{\partial \beta \partial \mu_{j-1}} &= \left[ \frac{-f'_{i,j-1}}{P_{i,j}} - \frac{(f_{i,j} - f_{i,j-1})(-f_{i,j-1})}{P_{i,j}^2} \right] (-\mathbf{x}_i)\end{aligned}$$

$$\frac{\partial^2 \log P_{i,j}}{\partial \mu_j^2} = \left[ \frac{f'_{i,j}}{P_{i,j}} - \left( \frac{f_{i,j}}{P_{i,j}} \right)^2 \right]$$

$$\frac{\partial^2 \log P_{i,j}}{\partial \mu_{j-1}^2} = \left[ \frac{-f'_{i,j-1}}{P_{i,j}} - \left( \frac{-f_{i,j-1}}{P_{i,j}} \right)^2 \right]$$

$$\frac{\partial^2 \log P_{i,j}}{\partial \mu_j \partial \mu_{j-1}} = \left[ \frac{(-f'_{i,j})(-f'_{i,j-1})}{P_{i,j}^2} \right]$$

The Hessian has a nonzero  $2 \times 2$  block within the full  $(J+2) \times (J+2)$  submatrix for  $\boldsymbol{\mu}$ . The relevant constraints on the terms for the fixed elements of  $\boldsymbol{\mu}$  are

$$\begin{aligned} \mu_{-1} &= -\infty, \quad \mu_0 = 0, \quad \mu_J = \infty \\ F_{i,-1} &= 0, \quad f_{i,-1} = 0, \quad f'_{i,-1} = 0 \\ F_{i,J} &= 1, \quad f_{i,J} = 0, \quad f'_{i,J} = 0. \end{aligned}$$

Finally, for the two most commonly used functional forms,

$$\text{logit: } F(t) = \Lambda(t), f(t) = \Lambda(t)[1 - \Lambda(t)], f'(t) = \Lambda(t)[(1 - \Lambda(t)) [1 - 2\Lambda(t)]]$$

$$\text{probit: } F(t) = \Phi(t), f(t) = \phi(t), f'(t) = -t \phi(t).$$

As Pratt (1981) showed, the second derivatives matrix is negative definite, so common gradient methods such as Newton or BFGS should be effective for maximizing the log likelihood function. Occasionally (rarely in our experience, however), the threshold parameters can become unordered during optimization. This points to the utility of a line search and a careful iteration. It is possible to force the threshold parameters to be ordered by reparameterizing them. For the model proposed in Section 5.2.7, we used the formulation

$$\mu_j = \mu_{j-1} + \exp(\alpha_j).$$

starting with  $\mu_0 = 0$ .

#### 4.8.6 Bayesian (MCMC) Estimation of Ordered Choice Models

Bayesian estimation of ordered choice models builds on the method pioneered by Albert and Chib (1993). The Gibbs sampler is constructed using a crucial device labeled “data augmentation.” [See Tanner and Wong (1987).] The binary choice case departs from

$$y_i^* = \boldsymbol{\beta}' \mathbf{x}_i + \varepsilon_i, \quad \varepsilon_i \sim \text{with mean 0 and known variance, } 1 \text{ (probit) or } \pi^2/3 \text{ (logit).}$$

$$y_i = 1 \text{ if } y_i^* > 0.$$

Let the prior for  $\boldsymbol{\beta}$  be denoted  $p(\boldsymbol{\beta})$ . Then, the posterior density for the probit or logit (symmetric

distribution) models is

$$p(\boldsymbol{\beta} | \mathbf{y}, \mathbf{X}) = \frac{p(\boldsymbol{\beta}) \prod_{i=1}^N F[(2y_i - 1)\boldsymbol{\beta}'\mathbf{x}_i]}{\int_{\boldsymbol{\beta}} p(\boldsymbol{\beta}) \prod_{i=1}^N F[(2y_i - 1)\boldsymbol{\beta}'\mathbf{x}_i] d\boldsymbol{\beta}},$$

where we use  $\mathbf{y}$  and  $\mathbf{X}$  (and later,  $\mathbf{y}^*$ ) to denote the full set of  $N$  observations. Estimation of the posterior mean is done by setting up a Gibbs sampler in which the unknown values  $y_i^*$  are treated as nuisance parameters to be estimated. For convenience at this point, we will assume the probit model is of interest. Conditioned on  $\boldsymbol{\beta}$  and  $\mathbf{x}_i$ ,  $y_i^*$  has a normal distribution with mean  $\boldsymbol{\beta}'\mathbf{x}_i$  and variance 1. However, when conditioned on  $y_i$  (observed), as well, the sign of  $y_i^*$  is known;

$$p(y_i^* | \boldsymbol{\beta}, \mathbf{y}, \mathbf{X}) = \text{normal with mean } \boldsymbol{\beta}'\mathbf{x}_i \text{ and variance 1, truncated at zero;} \\ \text{truncated from below if } y_i = 1 \text{ and from above if } y_i = 0.$$

Using basic results for Bayesian analysis of the linear model with known disturbance [see Greene (2008a, p. 605)] and a diffuse prior, the posterior for  $\boldsymbol{\beta}$  conditioned on  $\mathbf{y}^*$ ,  $\mathbf{y}$  and  $\mathbf{X}$  would be

$$p(\boldsymbol{\beta} | \mathbf{y}^*, \mathbf{y}, \mathbf{X}) = N_{\mathbf{K}}[\mathbf{b}, (\mathbf{X}'\mathbf{X})^{-1}] \text{ where } \mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}^*.$$

If, instead, the prior for  $\boldsymbol{\beta}$  is normal with mean  $\boldsymbol{\beta}^0$  and covariance matrix,  $\boldsymbol{\Sigma}$ , then the posterior density is normal with mean

$$E[\boldsymbol{\beta} | \mathbf{y}^*, \mathbf{y}, \mathbf{X}] = [\boldsymbol{\Sigma}^{-1} + (\mathbf{X}'\mathbf{X})]^{-1} (\boldsymbol{\Sigma}^{-1} \boldsymbol{\beta}^0 + \mathbf{X}'\mathbf{y}^*)$$

and

$$\text{Var}[\boldsymbol{\beta} | \mathbf{y}^*, \mathbf{y}, \mathbf{X}] = [\boldsymbol{\Sigma}^{-1} + (\mathbf{X}'\mathbf{X})]^{-1}$$

This sets up a strikingly simple Gibbs sampler for drawing from the joint posterior,  $p(\boldsymbol{\beta}, \mathbf{y}^* | \mathbf{y}, \mathbf{X})$ . It is customary to use a diffuse prior for  $\boldsymbol{\beta}$ . Then, compute initially,  $(\mathbf{X}'\mathbf{X})^{-1}$  and the lower triangular Cholesky matrix,  $\mathbf{L}$  such that  $\mathbf{L}\mathbf{L}' = (\mathbf{X}'\mathbf{X})^{-1}$ . (The matrix  $\mathbf{L}$  can be computed only once at the outset for the informative prior as well.) To initialize the iterations, any reasonable value of  $\boldsymbol{\beta}$  may be used. Albert and Chib suggest the classical MLE. The iterations are then given by

1. Compute the  $N$  draws from  $p(\mathbf{y}^* | \boldsymbol{\beta}, \mathbf{y}, \mathbf{X})$ .

Draws from the appropriate truncated normal can be obtained using

$$y_i^*(r) = \boldsymbol{\beta}'\mathbf{x}_i + \Phi^{-1}[\Phi(-\boldsymbol{\beta}'\mathbf{x}_i) + U(1 - \Phi(-\boldsymbol{\beta}'\mathbf{x}_i))] \text{ if } y_i = 1 \text{ and}$$

$$y_i^*(r) = \boldsymbol{\beta}'\mathbf{x}_i + \Phi^{-1}[U \Phi(-\boldsymbol{\beta}'\mathbf{x}_i)] \text{ if } y_i = 0$$

where  $U$  is a single draw from a standard uniform population.

2. Draw an observation on  $\boldsymbol{\beta}$  from the posterior  $p(\boldsymbol{\beta} | \mathbf{y}^*, \mathbf{y}, \mathbf{X})$  by first computing the mean

$$\mathbf{b}(r) = (\mathbf{X}'\mathbf{X})\mathbf{X}'\mathbf{y}^*(r).$$

Use a draw,  $\mathbf{v}$ , from the  $K$ -variate standard normal, then compute  $\boldsymbol{\beta}(r) = \mathbf{b}(r) + \mathbf{L}\mathbf{v}$ .

(We have used “ $(r)$ ” to denote the  $r$ th cycle of the iteration.) The iteration cycles between steps 1 and 2 until a satisfactory number of draws is obtained (and a burn-in number are discarded), then

the retained observations on  $\boldsymbol{\beta}$  are analyzed. With an informative prior, the draws at step 2 involving the prior mean and variance are slightly more time consuming. The matrix  $\mathbf{L}$  is only computed at the outset, but the computation of the mean adds a matrix multiplication and addition.

The extension to  $J+1$  ordered outcomes is now straightforward. We maintain the probit model, as is common. The model is, now,

$$y_i^* = \boldsymbol{\beta}'\mathbf{x}_i + \varepsilon_i, \quad \varepsilon_i \sim N[0, 1],$$

$$y_i = j \text{ if } \mu_{j-1} < y_i^* < \mu_j.$$

Diffuse priors are assumed for  $\boldsymbol{\beta}$  and  $\boldsymbol{\mu}$ , with the usual constraints on  $\mu_{-1}$  and  $\mu_0$ . Based on the same results as before, we still have

$$p(\boldsymbol{\beta} | \mathbf{y}^*, \boldsymbol{\mu}, \mathbf{y}, \mathbf{X}) = N_{\mathbf{K}}[\mathbf{b}, (\mathbf{X}'\mathbf{X})^{-1}].$$

$$p(y_i^* | \boldsymbol{\mu}, \boldsymbol{\beta}, \mathbf{y}, \mathbf{X}) = N(\boldsymbol{\beta}'\mathbf{x}_i, 1) \text{ truncated in both tails by } \mu_{j-1} \text{ and } \mu_j.$$

We will note below how to do the simulation for  $y_i^*$ . Finally, the authors provide the posterior for  $\mu_j$  ( $j = 1, \dots, J-1$ ), conditioned on the other threshold parameters,;

$$p(\mu_j | \boldsymbol{\beta}, \mathbf{y}^*, \boldsymbol{\mu}_{(j)}, \mathbf{y}, \mathbf{X}) \propto \prod_{i=1}^N \{ \mathbb{I}[y_i = j] \times \mathbb{I}[\mu_{j-1} < y_i^* < \mu_j] + \mathbb{I}[y_i = j+1] \times \mathbb{I}[\mu_j < y_i^* < \mu_{j+1}] \}$$

where the density is the posterior for  $\mu_j$  given the other threshold parameters, denoted  $\boldsymbol{\mu}_{(j)}$ , and the other parameters. The steps in the Gibbs sampler consist of initializing  $\boldsymbol{\beta}$  and  $\boldsymbol{\mu}$  as before, now with the MLE of the ordered probit model, then, in order,

1. Sample  $\mu_j$  from a uniform distribution with limits

$$Lower = \max_i \{ \max(y_i^* | y_i = j), \mu_{j-1} \} \quad (\text{i.e., the maximum over the } N \text{ observations})$$

$$Upper = \min_i \{ \min(y_i^* | y_i = j+1), \mu_{j+1} \}$$

Sampling from this uniform distribution is easily done by scaling a draw from  $U(0,1)$  by  $1/(Upper - Lower)$ .

2. Sample  $y_i^*$  from the truncated normal distribution where the underlying variable has mean  $\boldsymbol{\beta}'\mathbf{x}_i$  and standard deviation 1 and the truncation limits are  $\mu_{j-1}$  and  $\mu_j$  for the corresponding observation on  $y_i = j$ . The necessary result for this step is given in Greene (2008a, p. 575). To sample a draw from this distribution, define  $P_L = \Phi(\mu_{j-1} - \boldsymbol{\beta}'\mathbf{x}_i)$  and  $P_U = \Phi(\mu_j - \boldsymbol{\beta}'\mathbf{x}_i)$ . Note that  $P_L = 0$  if  $y_i = 0$ , and  $P_U = 1$  if  $y_i = J$ . Then, let  $U$  denote a draw from the  $U(0,1)$  population – a single uniform draw. Then, the draw for  $y_i^*$  is

$$y_i^* | y_i, \boldsymbol{\mu}, \boldsymbol{\beta}, \mathbf{x}_i = \boldsymbol{\beta}'\mathbf{x}_i + \Phi^{-1} [P_L + U \times (P_U - P_L)].$$

3. Sample  $\boldsymbol{\beta}$  from the multivariate normal population as shown earlier for the binary probit case. The only change is the data used to compute  $\mathbf{b}$ , now using the results of the doubly truncated sample in step 2 immediately above.

We then cycle through steps 1 – 3 for a large number of iterations (say tens of thousands). After discarding the first several thousand draws, the remaining draws on  $\beta$  and  $\mu$  constitute a sample from the joint posterior. The posterior mean is estimated by the average of the draws.

A convenient aspect of the MCMC approach to estimation is that often the estimator for a more complex model is easily obtained by adding layers to a simpler one. Consider the bivariate ordered probit model analyzed by Biswas and Das (2002). The model is a direct extension of the univariate model:

$$y_{i1}^* = \beta_1' \mathbf{x}_{i1} + \varepsilon_{i1}, \quad \varepsilon_{i1} \sim N[0, 1],$$

$$y_{i1} = j \text{ if } \mu_{j-1} < y_{i1}^* < \mu_j$$

$$y_{i2}^* = \beta_2' \mathbf{x}_{i2} + \varepsilon_{i2}, \quad \varepsilon_{i2} \sim N[0, 1],$$

$$y_{i2} = k \text{ if } \gamma_{k-1} < y_{i2}^* < \gamma_k.$$

$$\text{Corr}(\varepsilon_{i1}, \varepsilon_{i2}) = \rho$$

Each ordered probit is handled as before. The draws from the posterior of  $(\beta_1, \beta_2)$  are obtained by a two equation GLS regression; conditioned on the other parameters, the two latent regressions are a seemingly unrelated regressions system. The draws for  $(\mu_j, \gamma_k)$  are drawn jointly from a rectangle, with each dimension handled as in the univariate case. The draws on  $y_{i1}^*$  and  $y_{i2}^*$  are drawn from a truncated bivariate normal population. (Biswas and Das suggest to do this draw by a rejection method. It can be done in a “one draw” manner using a bivariate truncated normal analog to the method shown above. [See, e.g., Geweke (1991).]) The loose end is sampling from the posterior of  $\rho$ . Biswas and Das handle this by defining  $\Sigma$  to be an *unrestricted*  $2 \times 2$  covariance matrix of the two disturbances. The prior for  $\Sigma$  is assumed to be proportional to  $|\Sigma|^{-3/2}$ . This produces a conditional posterior for  $\Sigma$  that is an inverse Wishart population. [See Train (2003) for sampling from this population.] Note that they have introduced two new free parameters,  $\sigma_{11}$  and  $\sigma_{22}$  and are now estimating  $\sigma_{12} = \rho\sigma_1\sigma_2$ .

There is a peculiar loose end in the Biswas and Das (2002) study. In the ordered choice model, the scale parameters of the disturbances,  $\sigma_m^2 = \text{Var}[\varepsilon_{im}]$  are not identified and are normalized to 1.0. (In an alternative normalization of the model, one of the slopes is normalized at 1.0, which “identifies” the scale parameter – though not actually if that scale parameter is meant to be interpreted as the variation of  $\varepsilon$ . It merely moves the normalization off one of the parameters. See Section 8 below for applications.) Biswas and Das treated these variances as free parameters, and did not normalize one of the other parameters. As such, the model they purport to estimate is not identified. The evidence is in the reported values of the posterior means of  $\sigma_1^2 = 22.62$  and  $\sigma_2^2 = 13.33$ . These values are far outside the reasonable range for a choice model of this sort; they are supposed to be normalized at 1.0. (One might surmise that they are “identified” purely by the prior; there is no sample information about them.) This application points up a note of caution needed in MCMC estimation. The log likelihood function developed in Section 7.3.1 cannot be maximized if it is formulated in terms of an unrestricted  $\Sigma$  as used above. Ultimately, the derivatives will be collinear and the Hessian will be singular – that is the impact of a model that contains unidentified parameters. There is no counterpart control when using the Gibbs sampler. The signal that something has gone awry will arrive when the chain fails to converge, or when it arrives at a very different vector of posterior means from one run to another. It is necessary to check these failures – one run of the Gibbs sampler, regardless of how long it is, will not reveal this condition. (Redemption of the model would be obtained by formulating it in terms of a prior over  $\rho$  to begin with, and imposing the necessary normalizations on  $\sigma_1$  and  $\sigma_2$ .)

As noted earlier, the Bayesian segment of this literature is relatively compact and quite recent. Methodological contributions are offered by Albert and Chib (1993), Koop and Tobias (2006) and Imai et al. (2003) who have developed an “R” routine for some of the computations. Applications include Girard and Parent (2001), Biswas and Das (2002), Czado et al. (2005), Tomoyuki et al. (2006), Ando (2006), Zhang et al. (2007), Kadam and Lenk (2008) and Munkin and Trivedi (2008) and a handful of others. Doubtless there are more to come. Nonetheless, as of this writing, Bayesian analysis of ordered choice data is a small niche in the literature. There are, of course, a cornucopia of applications to binary data.

#### **4.8.7 Software For Estimation of Ordered Choice Models**

There are numerous commercial packages that can be used to estimate basic ordered choice models. (We mention the packages only by name here. Each of them is described in detail on their own respective website, listed below, so we will forego any detailed descriptions.) The primary ones in current use are *SAS*, *Stata*, *LIMDEP*, *NLOGIT* and *SPSS*. In addition, *Latent Gold* and a few other programs less oriented to cross section and panel data, including *RATS*, *Eviews* and *TSP*, also contain built-in estimators for the essential model. For Bayesians, there are routines in *R* provided in *ZELIG* by Imai et al. (2008). *WinBugs* also contains a routine for discrete choice models. The log likelihood is not particularly complicated, and *Gauss* and *Matlab* programs are also widely circulated.

For more advanced, exotic or obscure variants of the model, the choices are much more limited. These can, of course, be programmed by the user in the low level languages such as *Matlab*, or in many cases, even in the higher level matrix languages of the integrated packages such as *Stata*. For prepackaged routines, *Stata* and *NLOGIT/LIMDEP* contain optional features, such as heteroscedasticity and individual specific thresholds. Models with random coefficients can be fit with *PROC MIXED* in *SAS*, *GLAMM* in *Stata*, and with several of the routines in *NLOGIT*. To our knowledge, only *Latent Gold* and *NLOGIT/LIMDEP* have built in latent class treatments for ordered choice models. For panel data applications, the random effects model (Butler and Moffitt) is quite common as well and appears in all the familiar packages. Random effects models are “random constants” models. So any random parameters module can also handle random effects in a panel. That we are aware of, the fixed effects model with essentially unlimited numbers of effects (beyond the capacity to just add the dummy variables to the model) is available only in *NLOGIT* and *LIMDEP*.

The following is a list of the websites of the packages mentioned above. This is far from a complete list of software used in econometrics and statistics. For a lengthy guide that comes close to one, the econometric software resource

*Econometrics* <http://www.oswego.edu/~economic/econsoftware.htm>

is a useful reference point. The widely used packages are:

<i>Eviews</i>	<a href="http://www.eviews.com">http://www.eviews.com</a>
<i>Gauss</i>	<a href="http://www.aptech.com">http://www.aptech.com</a>
<i>Latent Gold</i>	<a href="http://www.statisticalinnovations.com/">http://www.statisticalinnovations.com/</a>
<i>LIMDEP</i>	<a href="http://www.limdep.com">http://www.limdep.com</a>
<i>Matlab</i>	<a href="http://www.mathworks.com">http://www.mathworks.com</a>
<i>NLOGIT</i>	<a href="http://www.nlogit.com">http://www.nlogit.com</a>
<i>RATS</i>	<a href="http://www.estima.com">http://www.estima.com</a>
<i>SAS</i>	<a href="http://www.sas.com">http://www.sas.com</a>
<i>SPSS</i>	<a href="http://www.spss.com">http://www.spss.com</a>
<i>Stata</i>	<a href="http://www.stata.com">http://www.stata.com</a>
<i>TSP</i>	<a href="http://www.tspintl.com">http://www.tspintl.com</a>
<i>WinBugs</i>	<a href="http://www.mrc-bsu.cam.ac.uk/bugs/winbugs/contents.shtml">http://www.mrc-bsu.cam.ac.uk/bugs/winbugs/contents.shtml</a>
<i>ZELIG</i>	<a href="http://gking.harvard.edu/zelig/">http://gking.harvard.edu/zelig/</a>

## 5. Specification Issues and Generalized Models

Anderson (1984, p. 2) discusses the inadequacy of the ordered choice model we have examined thus far. “We argue here that the class of regression models currently available for ordered categorical response variables is not wide enough to cover the range of problems that arise in practice. Factors affecting the kind of regression model required are (i) the type of ordered categorical variable, (ii) the observer error process and (iii) the “dimensionality” of the regression relationship. These factors relate to the processes giving rise to the observations and have been rather neglected in the literature.” Generalizations of the model, e.g., Williams (2006), have been predicated on Anderson’s observations, as well as some observed peculiarities in data being analyzed.

It is useful to distinguish between two directions of the contemporary development of the ordered choice model. Although it hints at some subtle aspects of the model (underlying data generating process), Anderson’s arguments, it will emerge, direct attention to the functional form of the model and its inadequacy in certain situations. Beginning with Terza (1985), a number of authors have focused, instead, on the fact that the model does not account adequately for individual heterogeneity that is likely to be present in micro- level data.

### 5.1 Functional Form Issues and the Generalized Ordered Choice Model (1)

Once again, referring to Anderson (1984, p. 2), “The dimensionality of the regression relationship between  $y$  and  $x$  is determined by the number of linear functions required to describe the relationship. If only one linear function is required, the relationship is one-dimensional; otherwise it is multi-dimensional. For example, in predicting  $k$  categories of pain relief from predictors  $x$ , suppose that different functions  $\beta_1'x$  and  $\beta_2'x$  are required to distinguish between the pairs of categories (*worse,same*) and (*same,better*), respectively. Then, the relationship is neither one-dimensional nor ordered with respect to  $x$ .” The fundamental flaw in the argument is in its opening premise. There is no regression relationship between  $y$  and  $x$ . The observed variable is merely a set of labels. What follows is curve fitting – suggesting that two equations might better fit two binary choices than a single one. (It remains to determine by what criterion different functions are *required*.) On the other hand, the author’s earlier (also p. 2) analysis of the data generating process puts a better face on the argument.

For example, Anderson and Philips (1981) refer to the “extent of pain relief after treatment.” *worse, same, slight improvement, moderate improvement, marked improvement or complete relief. In principle, there is a single, unobservable, continuous variable related to this ordered scale, [emphasis added] but in practice, the doctor making the assessment will use several pieces of information in making his judgment on the observed category. For example, he might use severity of pain, kind of pain, consistency in the time and degree of disability. We will refer to variables of the second type as “assessed” ordered categorical variables and argue that, in general, a different approach to modeling regression relationships is appropriate for the two types. Assessed ordered variables occur frequently in the biomedical, social and other social sciences.*

Thus, he argues that, at least in some situations, the dependent variable is not really ordered, or might not be. In such a case, he argues, essentially, that it makes sense to partition the outcomes, and treat them as a set of binary choices, or at least not as a single ordered choice. For the specific application considered, the issue depends crucially on whose assessment is being recorded, the doctor’s (not necessarily cleanly ordered as measured against some objective yardstick) or the patient’s (one would assume, necessarily ordered). The upshot is that, at least as

argued here, increasing the “dimensionality” of the fitting problem follows from the nature of the data generating process, not (evidently) from a need to accommodate curvature in the data.

### 5.1.1 Parallel Regressions

Anderson departs from the familiar ordered choice model that we have examined so far;

$$\text{Prob}(y \leq y_s | \mathbf{x}) = F(\theta_s - \boldsymbol{\beta}'\mathbf{x}), s = 1, \dots, k.$$

Continuing the line of argument suggested earlier, he then suggests his “new” model,

$$\text{Prob}(y = y_s | \mathbf{x}) = \frac{\exp(\beta_{0s}^* + \boldsymbol{\beta}'_s \mathbf{x})}{\sum_{t=0}^k \exp(\beta_{0t}^* + \boldsymbol{\beta}'_t \mathbf{x})}, \beta_{0k}^* = 0, \boldsymbol{\beta}_k = \mathbf{0}.$$

This is, of course, the multinomial logit model proposed by Nerlove and Press (1972) for  $k$  *unordered* choices. Later, it is observed “Model (5) [the model above] often gives a good fit to real data, even when the  $\boldsymbol{\beta}_s$  are *restricted to be parallel*. This is particularly true when the categories are ordered.” [Emphasis added.] Thus appears (apparently) the first occurrence of the “parallel regressions” notion in this literature. Note the implication is that the model is not intended for ordered data; but it seems to work well when applied to ordered outcomes. By “parallel,” the author states the restriction  $\boldsymbol{\beta}_s = -\phi_s \boldsymbol{\beta}$  where  $\phi_k \equiv 0$ . [Note that the last  $\phi_s$  is a parameter that is not identified under either the null or the alternative hypothesis because the corresponding  $\boldsymbol{\beta}_s = \mathbf{0}$ . See Andrews and Ploberger (1994).] A further identifying normalization (no longer merely for convenience) is  $\phi_1 \equiv 1$ . The resulting model,

$$\frac{\text{Pr}(y = y_s | \mathbf{x})}{\text{Pr}(y = y_k | \mathbf{x})} = \exp(\beta_{0s} - \phi_s \boldsymbol{\beta}), s = 1, \dots, k \quad [8]$$

is labeled the “*Stereotype Ordered Regression Model*.” As stated, the name is a misnomer, as the model does not enforce the ordering of the outcome; it is simply a parametric restriction on a model for *unordered* outcomes. Indeed, no linear restriction on the parameters of this model can enforce the ordering of the dependent variable, that is, the sequence

$$\text{Pr}(y \leq y_s | \mathbf{x}) < \text{Pr}(y \leq y_{s+1} | \mathbf{x}).$$

As he notes, the model “often gives a good fit to real data.” However, the ordering aspect of it would depend on the data. It is not a feature of the model. We should note, the underlying structure has been lost in this process. It is not possible to discern what underlying data generating process would give rise to such a functional form for a strictly ordered outcome that arises from an underlying continuous measure.

Anderson follows with a prescription for enforcing the ordering of the outcomes. “The next step is to order the  $\boldsymbol{\beta}_s$  to obtain a regression relationship. This is achieved by ordering the  $\phi_s$ ,

$$1 = \phi_1 > \phi_2 > \dots > \phi_k = 0. \quad [10]$$

“The ordered regression model [8] subject to constraints [10] will be termed the stereotype model.” This form is prescribed for ordered data. Unfortunately, the model is still not out of difficulty. The implied probabilities still do not enforce the ordering rule unless the constant terms are monotonically increasing;  $\beta_{01} < \beta_{02} < \dots < \beta_{0k}$ . Thus, Anderson’s remedy for the

“parallel regressions” restriction, if we enforce the ordering of the probabilities, is a progressive scaling of the parameter vector by the constants  $\phi_s$ , but it is not an internally consistent model for ordered choices without the constraint on the constant terms.

Long (1997) departs from our (now) familiar formulation of the ordered choice model.

$$\text{Prob}(y \leq j | \mathbf{x}_i) = F(\mu_j - \boldsymbol{\beta}'\mathbf{x}_i).$$

Differentiating these functions, we have

$$\partial \text{Prob}[y_i \leq j | \mathbf{x}_i] / \partial \mathbf{x}_i = -f(\mu_j - \boldsymbol{\beta}'\mathbf{x}_i) \boldsymbol{\beta}$$

This defines a set of binary choice models with different constants but common slope vector,  $\boldsymbol{\beta}$ . If we then fix the probability at, say  $P = P^*$  for any outcome, it must follow (by monotonicity of the cdf) that  $f(\mu_j - \boldsymbol{\beta}'\mathbf{x}_i)$  is fixed at  $f^*$ . It follows that *for a particular choice of probability*, we have

$$\partial \text{Prob}[y_i \leq j | \mathbf{x}_i] / \partial \mathbf{x}_i = f^* \boldsymbol{\beta} = \partial \text{Prob}[y_i \leq m | \mathbf{x}_i] / \partial \mathbf{x}_i, m = 0, \dots, J.$$

where  $f^*$  is the same for all  $j$ , that is, a multiple of the same  $\boldsymbol{\beta}$ . This is the feature of the model that has been labeled the “parallel regression assumption.” [See, e.g., Long (1997, p. 141).] This is an intrinsic feature of the ordered choice model. There is no obvious implication of the restriction for the underlying behavioral assumption – we will examine this issue in the next section. Note that the restriction cannot hold for a particular individual, since it requires the thresholds to adjust to equality. (I.e., we cannot fix all the probabilities to equal the chosen value at the same time. Rather, the “restriction” states that if  $P_1$  equals  $P^*$ , then the derivative is the same as if  $P_2$  equals the same  $P^*$ .)

### 5.1.2 Testing the Parallel Regressions Assumption – The Brant (1990) Test

Brant (1990), approaches the parallel regressions issue, but couches it in different terms. Defining

$$\gamma_j = \text{Prob}(y \leq j | \mathbf{x}) = F(\mu_j - \boldsymbol{\beta}'\mathbf{x}),$$

the logit form of the model implies (as well) that

$$\log \left( \frac{\gamma_j}{1 - \gamma_j} \right) = \mu_j - \boldsymbol{\beta}'\mathbf{x},$$

a “restriction” labeled the “proportional odds” restriction, or the “proportional odds model.” [McCullagh (1980)] Brant notes, this *is* a testable restriction, as we explore shortly. One is left to wonder, what feature of the model or of the behavior underlying it has been revealed when the null “hypothesis” of parallel regressions is rejected statistically, as it frequently is. Other than the purely mechanical observation that in a “model” with different coefficient vectors for each choice, the parallel regressions restriction is that those coefficients are the same, it is unclear in modeling terms, what the assumption means. Brant raised the same question. Before we reconsider that question, we will examine the proposed test procedure.

Several approaches to examining the parallel regressions feature have been developed.

All center on the set of implied binary choice “models” for the probit and logit cases,

$$\text{Prob}(y \geq j | \mathbf{x}) = F(\boldsymbol{\beta}'\mathbf{x} - \mu_j), j = 1, \dots, J-1.$$

Thus, one can, in principle, fit  $J-1$  such models separately. Each should produce its own constant term and a consistent estimator of the common  $\boldsymbol{\beta}$ . An “informal” examination of the differences [see Clogg and Shihadeh (1994, pp. 159-160)] should be revealing. A Lagrange multiplier test of the hypothesis is presented SAS Institute (2008). A much more straightforward (and intuitive) test is Brant’s (1990) Wald test which directly examines the restrictions

$$\boldsymbol{\beta}_1 = \boldsymbol{\beta}_2 = \dots = \boldsymbol{\beta}_{J-1}.$$

The Brant (1990) test of this hypothesis for the ordered logit model follows from the implication of the model,

$$\text{Prob}[y_i \geq j | \mathbf{x}_i] = \Lambda(\beta_{0j} + \boldsymbol{\beta}_j'\mathbf{x}_i)$$

where  $\beta_{0j} = \beta_0 - \mu_j$  and  $\Lambda(t)$  is the logistic cdf,  $1/(1+\exp(-t))$ . The slope vector  $\boldsymbol{\beta}_j$  should be the same in every equation. Thus, the specification implies  $J-1$  binary choice “models” that can be estimated one at a time, each with its own constant term and (by assumption) the same slope vector.

Expressions for the mechanics of the test appear in Long (1997, pp. 144-145.) The null hypothesis is equivalent to

$$H_0: \boldsymbol{\beta}_q - \boldsymbol{\beta}_1 = \mathbf{0}, q = 2, \dots, J-1$$

which can be summarized as

$$H_0: \mathbf{R}\boldsymbol{\beta}^* = \mathbf{0}$$

where

$$\mathbf{R} = \begin{bmatrix} \mathbf{I} & -\mathbf{I} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{I} & \mathbf{0} & -\mathbf{I} & \dots & \mathbf{0} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{I} & \mathbf{0} & \mathbf{0} & \dots & -\mathbf{I} \end{bmatrix}, \boldsymbol{\beta}^* = \begin{bmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \\ \boldsymbol{\beta}_3 \\ \vdots \\ \boldsymbol{\beta}_{J-1} \end{bmatrix}$$

The Wald statistic will be

$$\chi^2[(J-1)K] = (\mathbf{R}\hat{\boldsymbol{\beta}}^*)' \left[ \mathbf{R} \times \text{Asy. Var}[\hat{\boldsymbol{\beta}}^*] \times \mathbf{R}' \right]^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}}^*)$$

where  $\hat{\boldsymbol{\beta}}^*$  is obtained by stacking the individual binary logit estimates of  $\boldsymbol{\beta}$  (without the constant terms). The remaining complication in the computation is the asymptotic covariance matrix, which is computed as follows (using Brant’s results):

$$\text{Est. Asy. Cov}[\hat{\boldsymbol{\beta}}_j, \hat{\boldsymbol{\beta}}_m] = \left[ \sum_{i=1}^N \hat{\Lambda}_{ij} (1 - \hat{\Lambda}_{ij}) \mathbf{x}_i \mathbf{x}_i' \right]^{-1} \left[ \sum_{i=1}^N \hat{\Lambda}_{im} (1 - \hat{\Lambda}_{ij}) \mathbf{x}_i \mathbf{x}_i' \right] \left[ \sum_{i=1}^N \hat{\Lambda}_{im} (1 - \hat{\Lambda}_{im}) \mathbf{x}_i \mathbf{x}_i' \right]^{-1}$$

and  $\hat{\Lambda}_{ij} = \Lambda(\hat{\beta}_{0j} + \hat{\beta}'_j \mathbf{x}_i)$ . The test can be carried out for specific coefficients by removing all but the desired rows of  $\mathbf{R}$  in the computation of the statistic. By this device, for example, one can carry out the test for particular coefficients.

There are some loose ends in the computation. If the probabilities in the covariance matrix are based on the individual binary logit models, then the ordering of the probabilities is not preserved, and  $\Lambda_{ij} - \Lambda_{i,j-1} < 0$  is a possibility even though the theory rules it out. Brant suggests using the parameters of the restricted (basic ordered choice) model instead. Even with this practical fix, it remains true that the parameter estimates used in the test, each of which does have its own constant term, do not preserve the ordering of the probabilities in the model.

Table 9 displays the results of the Brant test for our ordered logit model of health satisfaction. The proportional odds restriction is clearly rejected. Loosely, it appears that the income coefficient displays the greatest variation across the cells. Both education and income appear to fail the test when it is applied individually.

**Table 9 Brant Test for Parameter Homogeneity**

```

+-----+
| Brant specification test for equal coefficient |
| vectors in the ordered logit model. The model |
| implies that logit[Prob(y>j|x)]=beta(j)*x - mj |
| for all j = 0,..., 3. The chi squared test is |
| H0:beta(0) = beta(1) = ... beta( 3)         |
| Chi squared test statistic =      71.76435    | (78.76988 based on the
| Degrees of freedom      =      15           | normal distribution)
| P value                  =      .00000      |
+-----+
=====
Specification Tests for Individual Coefficients in Ordered Logit Model
Degrees of freedom for each of these tests is 3
=====
Variable | Brant Test | Coefficients in implied model Prob(y > j). |
          | Chi-sq  P value | 0 | 1 | 2 | 3 |
=====
AGE      | 6.28   .09864 | -.0398 | -.0292 | -.0328 | -.0248 |
EDUC     | 19.89  .00018 | .1212 | .0786 | .0630 | -.0044 |
INCOME   | 13.32  .00398 | 1.9576 | .4959 | .1790 | -.0206 |
MARRIED  | 1.87   .59962 | .0674 | -.0228 | -.1486 | -.0896 |
KIDS     | 7.24   .06476 | .3218 | .2158 | .0189 | -.1231 |
=====

```

All this does naturally lead at least to some question of the model specification. For reasons we examine in more detail below, the non-proportional odds formulation is not a valid specification for the ordered logit model. Among the obvious reasons, the probabilities in the non-proportional odds model do not sum to one. If all the parameters can vary freely, as they do above, then each of the  $J$  binary choice models has been treated separately, and with no connection, there is no restriction on the sum of the probabilities. Moreover, there is no parametric restriction other than the one we seek to avoid that will preserve the ordering of the probabilities for all values of the data – that it does so for some data sets, or is a good “approximation” still leaves open the question of what specification failure makes sense to explain the finding, such as ours above.

Brant speculates at length about what model failures might lead to rejection of the hypothesis. The possibilities he lists include:

- (1) Misspecification of the latent regression,  $\beta'x$ ,
- (2) Heteroscedasticity of  $\varepsilon$  - “nonhomogeneous dispersion of the latent variable with varying  $x$ .”
- (3) Misspecification of the distributional form for the latent variable, i.e., “nonlogistic link function.”

He also considers a type of measurement error, such as the problem of “differential misclassification in the  $y$  observations.” Brant expresses little optimism that the test will likely uncover failures (1) or (2), reasoning that if the index or the variance are misspecified in the structural model, the misspecification will distort the estimators in the binary choice models similarly. For the distributional assumption, however, he shows that if some other distribution applies, such as the extreme value distribution, then the appropriate model should echo something similar to Anderson’s (1984) stereotype model, that is, with  $j$ -specific parameter vectors,  $(\theta_j, \phi_j, \beta)$ . In this case, rejection of the common  $\beta$  form in favor of the more general form would be expected. Note, though that even under this assumption, this does not suggest that one should expect to find completely separate  $\beta_j$ s. The differential multiple follows from the fact that even under the alternative distribution, the function is still parameterized in terms of a single index function. The scale factor is being induced by the different (from the logit) shape of the cdf with that same index function as its argument.

A more direct approach to testing against the distributional assumption is proposed by Johnson (1996) and Glewwe (1997). For this purpose, the null model is the ordered probit model based on the normal distribution. His Lagrange multiplier test is constructed by nesting the normal distribution within the broader Pearson family of distributions then testing against the null hypothesis of certain values of the parameters in the general form. [See Johnson, Kotz and Balakrishnan (1994).] It is noteworthy, at the end of the analysis, Glewwe (1997, p. 12) comes to the same juncture we have here. “A final question is what an applied econometrician should do when an ordered probit model does not pass the specification test.” Like all specification tests, the “alternative” is not well defined. He surmises that the test might be picking up an altogether different failure, such as an incorrect functional form. He does suggest some alternative strategies, and ultimately suggests that if the failure of the LM test persists, perhaps an ordered logit might be preferable.

The Brant test is easily transported to the ordered probit model. Using the usual approximation, each maximum likelihood binary choice estimator converges to

$$\hat{\beta}_j = \beta_j + \mathbf{H}_j^{-1} \mathbf{g}_j + o(1/N)$$

where  $\mathbf{H}_j^{-1}$  is the inverse of the information matrix and  $\mathbf{g}_j$  is the gradient of the log likelihood. Relying on the information matrix equality and the results of Berndt, Hall, Hall and Hausman (1974), we can estimate the matrix using the outer product of gradients and estimate the covariances of the derivatives with the sum of cross products. For the binary probit models,

$$\mathbf{g}_{ij} = \frac{(2q_{ij} - 1)\phi(\alpha_j + \beta_j'x_i)}{\Phi[(2q_{ij} - 1)\alpha_j + \beta_j'x_i]}(x_i)$$

where  $q_{ij} = 1(y_i > y_j)$ . The estimators of the submatrices needed for the test are

$$Est.Asy.Cov[\hat{\beta}_j, \hat{\beta}_m] = \left[ \sum_{i=1}^N \mathbf{g}_{ij} \mathbf{g}'_{ij} \right]^{-1} \left[ \sum_{i=1}^N \mathbf{g}_{ij} \mathbf{g}'_{im} \right] \left[ \sum_{i=1}^N \mathbf{g}_{im} \mathbf{g}'_{im} \right]^{-1}$$

Evidently this is not the explanation for the finding in Table 9. When we repeated the computations in Table 9 based on the ordered probit model, the chi squared statistic rose to 78.77.)

An intriguing point of the argument here is that it is not suggested that rejection of the supposed null hypothesis argues in favor of the non-proportional odds model as the alternative model. That model is not a viable alternative model, which leaves unanswered the fundamental question, what failure of the model does the Brant test reveal? Brant dwells on this question in his conclusion,

As previously mentioned, assessment of the proportionality assumption can also be based on fitting the augmented models (2.1) [the non-proportional odds model], as in Hutchison (1985) and Ekholm and Palmgren (1989). Similarly, a more directed approach can be based on fitting (3.2) [Anderson's (1984) stereotype model]. The augmented model approach is attractive in that it provides a more standard theoretical framework for developing tests. One drawback, however, is that specialized algorithms must be developed to fit the augmented models. A more serious problem is inherent in the models themselves. For example, if one wishes to extend the use of model (2.1) beyond the values of  $\mathbf{x}$ 's actually observed, the  $\beta_j$ 's must be constrained to ensure monotonicity of the extrapolated  $\gamma_j$ 's. Similar difficulties pertain to (3.2). Depending on the range of admissible values of  $\mathbf{x}$ , this can lead to technical difficulties in fitting and the need for nonstandard likelihood theory to allow for the possibility of estimates falling on the boundary of the parameter space. *It may be best then to view (2.1) and (3.2) not as scientifically meaningful models, but as directional alternatives helpful in validating the simpler proportional odds model.* [Emphasis added.]

We conclude that the Brant test is useful for supporting or for casting doubt on the basic model. It does not seem to be useful for pointing toward what might appear superficially to be an alternative specification based on freeing the parameter vectors in  $\gamma_j$ .

We note, finally, the response of some analysts to the failure of the base model (the ordered choice model), say as evidenced by the Brant test, is to switch to the unordered multinomial logit model as an alternative. Williams (2006, p. 5) dismisses this approach because the alternative proliferates parameters and is difficult to interpret. In fact, switching to the multinomial logit model as an alternative to the ordered choice model, assuming that some ordered choice model was appropriate to begin with, substitutes a manifestly misspecified model for one that was merely suspect and, probably, in need of refinement. The multinomial logit model for unordered choices is applicable to a different situation entirely. It produces coefficients, but it would be arduous at best to translate them into something meaningful to describe the behavior of an ordered random variable, such as the outcome of an attitude survey. So, following Williams, we will eschew further consideration of the multinomial logit model for unordered choices in this review.

### 5.1.3 Generalized Ordered Logit Model (1)

Quednau (1988), Clogg and Shihadeh (1994), Fahrmeir and Tutz (1994), McCullagh and Nelder (1989) have proposed versions of the ordered choice models based essentially on the “non-proportional odds” form given above. Fu (1998) and Williams (2006) have recently provided working papers and a *Stata* program (GOLogit and GOLogit2) that implement and refine the model. Williams (2006) suggests that his development is an extension of Fu’s so we focus on the latter. Motivated by the frequent rejection of the null hypothesis by Brant’s (1990) test [see Williams (2006, p. 3)], a suggested alternative model derives from the core specification

$$\text{Prob}(y_i > j) = F(\alpha_j + \beta_j' \mathbf{x}_i) = \frac{\exp(\alpha_j + \beta_j' \mathbf{x}_i)}{1 + \exp(\alpha_j + \beta_j' \mathbf{x}_i)}, j = 0, 1, \dots, J-1,$$

where, now,  $\mathbf{x}_i$  does not contain a constant term. (Note that this is the form used by Brant to motivate his analysis.) The implication is

$$\begin{aligned} \text{Prob}(y_i = 0 | \mathbf{x}_i) &= 1 - F(\alpha_0 + \beta_0' \mathbf{x}_i) \\ \text{Prob}(y_i = 1 | \mathbf{x}_i) &= F(\alpha_0 + \beta_0' \mathbf{x}_i) - F(\alpha_1 + \beta_1' \mathbf{x}_i) \\ \text{Prob}(y_i = j | \mathbf{x}_i) &= F(\alpha_{j-1} + \beta_{j-1}' \mathbf{x}_i) - F(\alpha_j + \beta_j' \mathbf{x}_i) \\ \text{Prob}(y_i = J | \mathbf{x}_i) &= F(\alpha_{J-1} + \beta_{J-1}' \mathbf{x}_i). \end{aligned}$$

We label this the “(1)” form of the generalized ordered choice model. We will examine two other forms, with (unfortunately) the same name. The “(1)” does not indicate first chronologically; that would be Terza’s (1985) formulation. It is simply the first one presented in this review. This model is related to, but is not quite the same as the implied alternative in Brant’s analysis. In fact, Brant’s alternative model, which is equivalent to  $\text{logit}(\gamma_{ij}) = \alpha_j + \beta_j' \mathbf{x}_i$ , treats each of the  $J+1$  outcomes of  $y_i$  as a separate event – the probabilities vary completely independently and need not even sum to one or a number less than one. As he notes, it should not be viewed as a valid model as it stands. In the model suggested above, the ordering aspect of the observed variable is preserved somewhat, in that the formulation implies a connection between the events  $y_i = j$  and  $y_i = j-1$ . On the other hand, with no constraints imposed on the parameters of the model, although the probabilities sum to one by construction, there is no assurance that they are positive. Brant anticipated this uncomfortable feature of the model in the conclusion related earlier. Long and Freese (2006, p. 221) observe this as well, but note that “To ensure that the  $\text{Pr}(y=j|\mathbf{x})$  is between 0 and 1, the condition  $(\tau_j - \beta_j' \mathbf{x}) \geq (\tau_{j-1} - \beta_{j-1}' \mathbf{x})$  must hold.” (The inequality must actually be strong if the probabilities are to be nonzero as well.) Rewrite the restriction as  $(\tau_j - \tau_{j-1}) > (\beta_j - \beta_{j-1})' \mathbf{x}$ . The only way to ensure that this is true for *every* possible configuration of  $\mathbf{x}$  is to have  $\tau_j > \tau_{j-1}$  and  $\beta_j = \beta_{j-1}$ , which is where we began.

The problem of negative probabilities was raised much earlier. Williams (2006) invoking McCullagh and Nelder (1989, p. 155) observes

“The usefulness of non-parallel regression models is limited to some extent by the fact that the lines must eventually intersect. Negative fitted values are then unavoidable for some values of  $\mathbf{x}$ , though perhaps not in the observed range. If such intersections occur in a sufficiently remote region of the  $\mathbf{x}$ -space, this flaw in the model need not be serious.”

This seems to be a fairly rare occurrence, and when it does occur there are often other problems with the model, e.g. the model is overly complicated and/or there are very small Ns for some categories of the dependent variable. *gologit2* will give a warning message whenever any in-sample predicted probabilities are negative. If it is just a few cases, it

may not be worth worrying about, but if there are many cases you may wish to modify your model, data, or sample, or use a different statistical technique altogether.

The prescription relates to fitting the function to the data, but not to the underlying model. I.e., the “flaw” in the model is not that it sometimes produces negative fitted probabilities; it is that it does not impose the positivity of the fitted probabilities in the structure to begin with. In practical terms, as Williams (2006) suggests, the model is usually estimable, and the problem does not arise. If one begins the iterations with starting values obtained from the “constrained” ordered logit model, then at least at the starting values, one is assured that all probabilities are positive. As the iterate moves away from the starting values, as any probability associated with an observed outcome moves toward zero, it will impose a large penalty on the log likelihood – in principle if a probability for an observation becomes negative, it exerts an infinite penalty. The practical upshot is that it seems reasonable that in spite of its potential for internal inconsistency, this model is likely to be estimable. Table 10 below shows the results for our ordered choice example. (Williams (2006) has published a *Stata* program (`GOLogit2`) for this purpose. We used the `MAXIMIZE` command in *NLOGIT*.) The estimates in Table 10 have been reordered so that coefficients associated with specific independent variables are grouped contiguously, rather than coefficients associated with specific outcomes. Inspection of the sets of estimates certainly suggests that the coefficients differ substantially across  $j$ . A likelihood ratio test would be based on

$$\chi^2[15] = 2(-5713.579 - (-5747.822)) = 68.486.$$

The 95% critical value from the table is 24.996. Thus, the hypothesis of the restricted model is decisively rejected.

A peculiarity of this “generalization” of the ordered logit model is that it does not appear to define a random variable. The specification states that “If  $y_i = j$ , then the probability that  $y_i$  equals  $j$  is as follows.” In spite of its appearance, the model does not state that the probability that a well defined random variable is equal to the given value is equal to the function. There is no underlying continuous variable that can be structured so as to produce the observed outcome. The latent regression approach is not available to motivate the outcome variable; “ $y^* = \alpha_j + \beta_j'x + \varepsilon$  then  $y^* = j$  under some condition,” since in order to generate  $y^*$ , one would need to know the appropriate  $j$  in advance. Consider, for example, that it is not possible to simulate the values of the random variable,  $y$ , defined in the probability statement. In order to assign a probability to the outcome we would first have to know what the outcome is. No data generating process produces the random variable described in the probability statement. This model, *as stated*, has the uncomfortable feature that it does not define what the “random variable “ $y$ ” is; it defines  $y$  in terms of itself. Ultimately, the problem is the ordered nature of the observed response. The ordering is incompatible with that much free parameter variation in the statement of the probabilities. If a model of an ordered random variable is to be complete and internally consistent, then ultimately, the observed response must be derived as a classification of a set of underlying events. The early writers on this model, Aitchison, McCullagh, Snell, etc., returned repeatedly to the theme of the underlying continuous variable for this reason.

**Table 10 Estimated Ordered Logit and Generalized Ordered Logit (1)**

Ordered Probability Model					
Underlying probabilities based on Logit					
Dependent variable	HEALTH				
Number of parameters	9				
Log likelihood function	-5747.822				
Restricted log likelihood	-5873.696				
Chi squared	251.7485				
Degrees of freedom	5				
Prob[ChiSq > value] =	.0000000				

Variable	Coefficient	Standard Error	b/St.Er.	P[ Z >z]	Mean of X
-----+Index function for probability					
Constant	3.51646***	.20386393	17.249	.0000	
AGE	-.03213***	.00287527	-11.175	.0000	43.445213
EDUC	.06467***	.01249042	5.178	.0000	11.416711
INCOME	.42434**	.18676718	2.272	.0231	.3488957
MARRIED	-.06451	.07455295	-.865	.3869	.7525106
KIDS	.11452*	.06686343	1.713	.0868	.3796028
-----+Threshold parameters for index					
Mu(1)	2.12143***	.03705411	57.252	.0000	
Mu(2)	4.43343***	.03901916	113.622	.0000	
Mu(3)	5.37670***	.05199838	103.401	.0000	

User Defined Optimization		Generalized Ordered Logit Model (1)			
Maximum Likelihood Estimates					
Log likelihood function	-5713.579				

Variable	Coefficient	Standard Error	b/St.Er.	P[ Z >z]	Ordered Logit Estimates
B01	2.69537***	.60687427	4.441	.0000	$\alpha = 3.51646$
B11	1.04676***	.25130943	4.165	.0000	$\alpha - \mu_1 = 1.39503$
B21	-.67133***	.25379817	-2.645	.0082	$\alpha - \mu_2 = -.91697$
B31	-1.09368***	.36891087	-2.965	.0030	$\alpha - \mu_3 = -1.86024$
B02	-.04080***	.00765120	-5.332	.0000	AGE
B12	-.02925***	.00342622	-8.538	.0000	-0.03213
B22	-.03261***	.00375790	-8.677	.0000	
B32	-.02427***	.00496850	-4.885	.0000	
B03	.12009***	.03870888	3.102	.0019	EDUC
B13	.07635***	.01552693	4.917	.0000	0.06467
B23	.06222***	.01572984	3.956	.0001	
B33	-.00252	.02338475	-.108	.9141	
B04	1.98158***	.45270808	4.377	.0000	INCOME
B14	.51201**	.21458584	2.386	.0170	0.42434
B24	.18838	.23361123	.806	.4200	
B34	-.11631	.28567597	-.407	.6839	
B05	.05870	.17101512	.343	.7314	MARRIED
B15	-.02514	.08628965	-.291	.7708	-.06451
B25	-.15166	.09659029	-1.570	.1164	
B35	-.07179	.12962378	-.554	.5797	
B06	.34731*	.18409495	1.887	.0592	KIDS
B16	.21913***	.08186585	2.677	.0074	0.11452
B26	.01939	.08827954	.220	.8261	
B36	-.11322	.12160210	-.931	.3518	

Note: \*\*\*, \*\*, \* = Significance at 1%, 5%, 10% level.

In our application, we began the computations by collapsing several categories of the dependent variable, for example, combining categories 0,1,2 into the observed “0.” Likewise, Boes and Winkelmann (2006a) combined the lowest three categories of their observed satisfaction measure. The implication of the generalized ordered probability model (1) would be either that in the collapsed model, the coefficient vector associated with the zero outcome is an ambiguous mixture of the original three coefficient vectors, or in the original model, the lowest three categories have the same coefficient vector – that would legitimize the aggregation of the three cells. It is a matter of interpretation. The implication, however, is that the population “parameters ( $\alpha_j, \beta_j$ ) exist as a function of the way that the analyst codes the dependent variable. More to the point, the model parameters, e.g., the data generating mechanism, cannot consistently exist apart from the observed data themselves. This returns to the characteristic that it is not possible to simulate a well defined random variable that obeys the probability laws defined above. This might seem to be the case in the base case model, since the “cut points” are identified with the outcomes. However, it is not the case there, since  $\mu_j$  exists (in theory) as an unknown location on the real line, independently of the random variable that drives the model,  $y^* = \beta'x + \varepsilon$ . There is no counterpart to  $y^*$  in the Generalized Ordered Logit Model (1).

All this said, it remains true that the “parameters” of the model *can* be computed, as we have done in Table 10. The least favorable view is that this is just curve fitting. However, if so, and if the ordered logit model (same  $\beta$ ) really is appropriate, then one should replicate, at least approximately the original “constrained” model. To some degrees, as evident below, that is what occurs; this could be viewed as a (numerically) inefficient estimator of the original model. But, in the same spirit as the Brant test, the same question emerges. To the extent that this procedure does not mimic the original model – the separate parameter vectors really do differ, as ours do in Table 10 – then what has it found? Since the model, such as it is, is not a valid probability model, the same loose end emerges. It must be picking up *some* failure of the original model. One might guess that Brant’s speculations about a set of explanations for rejection of the null hypothesis by his test would be helpful here as well.

We have labeled the model discussed here the “Generalized Ordered Choice Model (1).” Forms “(2)” and “(3)” are discussed below. The preceding is an orthodox interpretation of the model specification. Later, in Section 5.2, we will find that with a straightforward reinterpretation of what is ultimately the same model structure, an internally consistent model of a random variable does emerge. Since the models are only superficially different, we will label the threshold models in Section 5.2 the “(2)” forms of the Generalized Ordered Choice Model.”

### 5.1.4 The Single Crossing Feature of the Ordered Choice Model

The partial effects shown in the preceding examples vary with the data and the parameters. Since the probabilities must sum to one, the partial effects for each variable must sum to zero across the probabilities. It can also be shown that for the probit and logit models, this set of partial derivatives will change sign exactly once in the sequence from 0 to  $J$ , a property that Boes and Winkelmann (2006b) label the “single crossing” characteristic. [Crawford, Pollak and Vella (1988) explore this feature of the model at length.] For a positive coefficient,  $\beta_k$ , the signs moving from 0 to  $J$  will begin with negative and switch once to positive at some point in the sequence. The following is extracted from Table 4 in Boes and Winkelmann (2006b, page 22). (The “0-2” bracket is obtained by grouping the relatively low number of observations with the three lowest values in the original data.) Partial effects are shown with estimated standard errors in parentheses.

**Table 11 Boes and Winkelmann Estimated Partial Effects**

Response	0-2	3	4	5	6	7	8	9	10
Men									
OProbit	-0.016 (0.003)	-0.014 (0.001)	-0.016 (0.001)	-0.037 (0.003)	-0.020 (0.009)	0.003 (0.003)	0.059 (0.009)	0.027 (0.005)	0.014 (0.005)
GOProbit	-0.020 (0.007)	-0.022 (0.006)	-0.014 (0.004)	-0.027 (0.005)	-0.037 (0.006)	-0.005 (0.007)	0.088 (0.033)	0.039 (0.109)	-0.002 (0.089)
Women									
OProbit	-0.004 (0.002)	-0.005 (0.001)	-0.005 (0.001)	-0.016 (0.005)	-0.008 (0.012)	-0.003 (0.003)	0.020 (0.011)	0.012 (0.004)	0.008 (0.006)
GOProbit	-0.009 (0.008)	0.005 (0.016)	-0.011 (0.020)	-0.036 (0.015)	-0.040 (0.013)	0.038 (0.029)	0.064 (0.116)	-0.008 (0.125)	-0.003 (0.027)

The same effect can be seen in Table 3 for our application.

The “GOProbit” results – a probit version of Williams’s (2006) GOLogit approach – show the effect of relaxing the single crossing restriction. However, for men, the model seems to be preserving the restriction on its own – the second crossing at  $y = 10$ , produces a marginal effect that differs only trivially from zero, with a “z-value” of only 0.022. For women, however, one is in the uncomfortable position of now explaining four crossings which make the model seem a bit unstable. None of the estimated effects are statistically significant, in contrast to the ordered probit model, and in fact, two of the crossings rest on what looks like a maverick finite sample outcome at  $y=3$ . One the other hand, the results that remain force the analyst into a counterintuitive position of arguing that higher incomes are associated with lowered probabilities of reporting a high subjective well being – perhaps a widespread *Richard Cory* effect. The authors’ description of the results (from their pages 12 and 13) suggests the appeal of a less sharp statement about specific outcomes; the right tail result is suggested to reflect a zero effect, which of course removes the remaining extra crossing.

Table 4 summarizes the marginal probability effects of income by gender. Consider, for example, the results for men and take the ceteris paribus effect of increasing logarithmic household income by a small amount on the probability of responding a SWB level of “8”. Table 4 shows a value of 0.059 for the standard model. This means that the probability of a response of “8” increases by 0.059 percentage points if we increase logarithmic income by 0.01, which corresponds approximately to a one-percent increase in level income. A doubling of income, i.e., a change in logarithmic income by 0.693, increases the probability of response “8” by about  $0.059 \times 0.693 \times 100$ , or about 4.09 percentage points, ceteris paribus.

Comparing the MPE’s among the three different models and over all possible outcomes, we obtain the following main results. For men all models suggest that more income significantly reduces the probability of low SWB (0-5), and significantly

increases the probability of response “8”. For high SWB responses (9-10), the standard model predicts a strong positive relationship between income and SWB, whereas the generalized model and also the binary models do not find a significant effect. Since the restricted OProbit is clearly rejected, we conclude that income has no effect on positive well-being. *Our preferred specification supports the asymmetry hypothesis for men: higher income decreases the probability of negative well-being (low SWB), but it does not affect the probability of positive well-being (high SWB).* [Emphasis added.] For women the relationship between income and SWB is relatively weak. While the standard model finds small but significant effects for low and high SWB responses, *the generalized model predicts a significant negative effect only on the probability of responses “5” and “6”.* [Emphasis added.] The gender difference might be explained by social norms that assign the role of primary income earner to men and therefore make income a relatively more important determinant of male well-being (see also Lalive and Stutzer 2004).

Figure 11 shows graphically the values in Table 11. The ordered probit and generalized ordered probit models do not seem to be giving different accounts. The latter does seem to be exaggerating the outcome at choice 8, or perhaps suggesting a significant spike associated with that outcome, that then would want some explanation. The force of the model extension seems to be to produce a much more pronounced effect in the middle of the distribution. The fact that the heightened impact is negative for  $y = 6$  and positive for  $y = 8$ , followed for both genders by a sharp return to zero at  $y = 9$ , seems a bit counterintuitive.

The shortcoming of the ordered choice model that produces the single crossing result is the linearity of the single index formulation. One can achieve the same result as above without resort to the generalized model simply by building the desired curvature into the index function itself. In the figure below, we have re-estimated our original model using not “Health” coded 0 to 4, but the original Health Satisfaction variable, coded 0 to 10, the same as in Boes and Winkelmann’s study. (They are subsets of the same data base.) Income is included in linear, squared and cubed form, so that the marginal effect of income on any outcome is

$$\delta_{\text{INCOME}}(j) = [f(\mu_{j-1} - \beta'x) - f(\mu_j - \beta'x)] \times (\beta_{\text{INCOME}} + 2\beta_{\text{INCOME-SQ}} \text{INCOME} + 3\beta_{\text{INCOME-CUBE}} \text{INCOME}^2)$$

We have evaluated this at the means of all the variables in the model. The results are shown in Figure 12 along with the results from the original model. While the effects still only cross zero once, the formulation does not force this – we will accept the data’s word for it that the partial effect of income does indeed (at least seem to) start negative and become positive, conforming to intuition that greater income is broadly associated with greater health satisfaction. It is interesting as well that the linear index model produces essentially the same results.

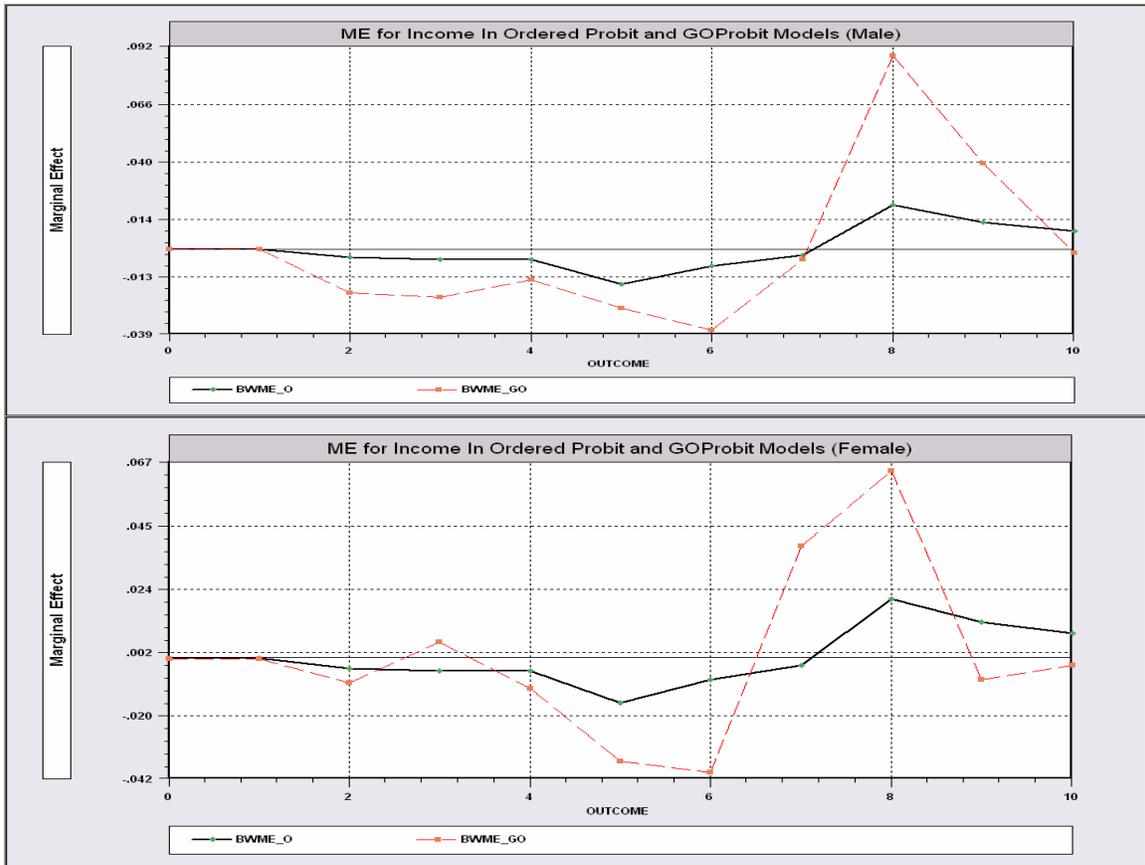


Figure 11 Estimated Partial Effects in Boes and Winkelmann (2006b) Models

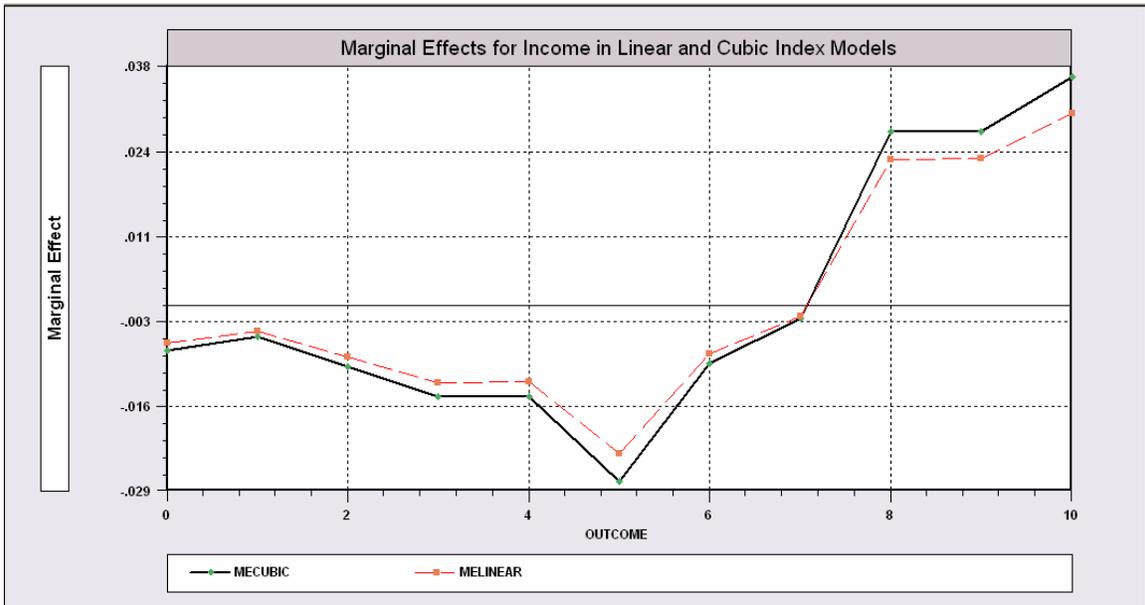


Figure 12 Estimated Partial Effects for Linear and Nonlinear Index Functions

### 5.1.5 Choice Invariant Ratios of Partial Effects

Boes and Winkelmann (2006a) note that for any two continuous covariates,  $x_{ik}$  and  $x_{il}$

$$\frac{\partial \text{Prob}[y_i = j | \mathbf{x}_i] / \partial x_{i,k}}{\partial \text{Prob}[y_i = j | \mathbf{x}_i] / \partial x_{i,l}} = \frac{\beta_k}{\beta_l},$$

which is independent of the outcomes. This is a feature of the assumed underlying utility function, the same as in any regression model. Any single index function model that is of the form

$$\text{Prob}(y_i = j | \mathbf{x}) = G(\boldsymbol{\beta}'\mathbf{x}_i)$$

will have this feature; it is a consequence of the chain rule of the calculus. It is unclear what the behavioral implications will be; that would be specific to the application. Boes and Winkelmann (2006a) develop this theme in some detail. In their application,

$$SWB = \alpha + \beta_{INCOME} INCOME + \beta_{UNEMPLOYMENT} UNEMPLOYED + \dots + \varepsilon.$$

The authors are interested in the notion of “compensating variation.” For their purpose, “what is the income increase required to offset the negative well-being effect of unemployment?” (They finesse the binary nature of the unemployment variable by considering the issue from the point of view of the population unemployment *rate*.) By equating the total differential of  $\text{Prob}(y = j | \mathbf{x})$  to zero, they find that the interesting “tradeoff ratio” is the negative of the ratio of the partial effects, as shown above. The implication of the standard model is that the tradeoff ratios are the same for all outcomes.

In the semiparametric models developed in Section 8, in which it is not possible to compute the CDF or the density – the semiparametric aspect of the model is to dispense with the assumption of a specific density – ratios of coefficients become important outputs of the estimation process. Stewart (2003, 2005) develops this idea at some length.

The common feature of this and the extensions preceding it are that the functional form is built around the outcomes. The single index models considered thus far do not provide sufficient curvature to accommodate what Anderson (1984) called the “dimensionality” of the problem. The greater fit achieved by the expanded model may have less to do with describing the underlying data generating process than with matching the fitted function to the pattern in the observed data. The modifications of the ordered choice model described in the next sections also achieve some of this increased “fit” but do so within the structure of the original behavioral model.

### 5.1.6 Methodological Issues

The various generalizations of the model suggested above do deal with the problems of parallel regressions and single crossing, but potentially create new ones. The heterogeneity in the parameter vector is an artifact of the coding of the dependent variable, not a manifestation of underlying heterogeneity in the dependent variable induced by behavioral differences. It is unclear what it means for the marginal utility parameters to be structured in this way. To put a better face on it, we might best interpret this as a semiparametric approach to modeling what is apparently underlying heterogeneity, however, again, it is not clear why this should be manifest in parameter variation across the outcomes instead of across the individuals in the sample. One

would assume that the failure of the Brant test to support the model with parameter homogeneity is, indeed, signalling some failure of the model. But, it is unclear what that failure is. The more difficult problem of this generalization of the model is that the probabilities in this model need not be positive, and there is no parametric restriction (other than the restrictive model one we started with) that could achieve this. The restrictions would have to be functions of the data. (The problem is noted by Williams (2006), but dismissed as a minor issue. Boes and Winkelmann suggest that the problem could be handled through a “nonlinear specification.”)

One might still argue that there are differences across the individuals at the “low” end vs. the “high” end of the distribution. The excerpt from Boes and Winkelmann above would suggest this. In fact, the single crossing aspect of the model accommodates this feature. Still, something more akin to a latent class structure would seem to apply under this interpretation. In such a setting, one is likely to find that the high outcomes are more likely for some classes than others. The advantage of this approach would be that the class structure can be assumed to be exogenous. One is not forced to make the model structure endogenous to the observed outcomes.

## 5.2 Accommodating Heterogeneity

The presence or absence of individual heterogeneity not contained explicitly in the model is likely the most fundamental difference between the bioassay and social science applications of ordered choice models. In the analysis of a population of fruit flies or aphids, the analyst is probably safe in assuming that the population is homogeneous enough to treat with a zero mean, homoscedastic disturbance in the latent tolerance equation and single parameter, homogeneous thresholds in the observation mechanism. The analysis of a population of congressional representatives or heads of households responding to a survey about health satisfaction or subjective well being will be far from that situation. Consider, as well, the fundamental difference in the underlying equation. For a simple insecticide experiment, the implied underlying regression will be

$$Tolerance_{ir}^* = \alpha + \beta Treatment_{ir} + \varepsilon_{ir}$$

where  $i$  indicates a group (treatment level) and  $r$  indicates a member of that group. The entire “behavioral” aspect of the model is embedded in the random term, the “tolerance” to the treatment. The ordered “choice” is

$$y_{ir} = \begin{array}{ll} 0 & \text{if } (Tolerance_{ir} - \alpha - \beta Treatment_{ir}) \leq \alpha_1 \quad (\text{dead}) \\ 1 & \text{if } (\alpha_1 < Tolerance_{ir} - \alpha - \beta Treatment_{ir} \leq \alpha_2) \quad (\text{moribund}) \\ 2 & \text{if } (\alpha_2 < Tolerance_{ir} - \alpha - \beta Treatment_{ir}) \quad (\text{alive}) \end{array}$$

It seems safe to assume that the individual observations are sufficiently homogeneous in dimensions that one could hope to measure that the simple, canonical model above is an adequate description of the outcome variable that we will ultimately observe. In contrast, for the *subjective* well being (SWB) application, the right hand side of the behavioral equation will include variables such as *Income, Education, Marital Status, Children, Working Status, Health*, and a host of other measurable and unmeasurable, and *measured* and *unmeasured* variables. In individual level behavioral models, such as

$$SWB_{it} = \beta'x_{it} + \varepsilon_{it},$$

the relevant question is whether a zero mean, homoscedastic  $\varepsilon_{it}$ , can be expected to satisfactorily accommodate the likely amount of heterogeneity in the underlying data, and whether it is reasonable to assume that the same thresholds should apply to each individual.

Beginning with Terza (1985), analysts have questioned the adequacy of the ordered choice model from this direction. As shown below, many of the proposed extensions of the model, such as heteroscedasticity, parameter heterogeneity, etc., parallel developments in other modeling contexts (such as binary choice modeling and modeling counts such as number of doctor visits or hospital visits). The regression based ordered choice model analyzed here does have a unique feature, that the thresholds are part of the behavioral specification. This aspect of the specification has been considered as well.

### 5.2.1 Threshold Models – The Generalized Ordered Probit Model (2)

The model analyzed thus far assumes that the thresholds  $\mu_j$  are the same for every individual in the sample. Terza (1985), Pudney and Shields (2000), Boes and Winkelmann (2006a), Greene, Harris, Hollingsworth and Maitra (2008) and Greene and Hensher (2008), all present cases that suggest individual variation in the set of thresholds is a degree of heterogeneity that is likely to be present in the data, but is not accommodated in the model. A precursor to this literature is Farewell (1982), who proposes an ordered Weibull model,

$$\text{Prob}(y_i > j \mid \mathbf{x}_i) = \exp(-\exp(\theta_j - \boldsymbol{\beta}'\mathbf{x}_i)).$$

To accommodate the possibility of latent heterogeneity, he suggests

$$\theta_{ij} = \theta_j^* + \eta_i$$

with  $\theta_0 = 0$ , so that the spacing between thresholds is preserved, but the location of the set of thresholds varies across individuals. The extreme value functional form is unique. However, the shift of the thresholds points toward the later generalizations of the model, beginning with Terza (1985).

Terza's (1985) generalization of the model is equivalent to

$$\mu_{ij} = \mu_j + \boldsymbol{\delta}'\mathbf{z}_i.$$

(This is the special case of the generalized model that he used in his application – his fully general case allows  $\boldsymbol{\delta}$  to differ across outcomes.) The model is reformulated later to assume that the  $\mathbf{z}_i$  in the equation for the thresholds is the same as the  $\mathbf{x}_i$  in the regression. For the moment, it is convenient to remove the constant term from  $\mathbf{x}_i$ . In Terza's application, in which there were three outcomes,

$$y_i^* = \alpha + \boldsymbol{\beta}'\mathbf{x}_i + \varepsilon_i$$

and

$$y_i = \begin{cases} 0 & \text{if } y_i^* \leq 0 \\ 1 & \text{if } 0 < y_i^* \leq \mu + \boldsymbol{\delta}'\mathbf{x}_i \\ 2 & \text{if } y_i^* > \mu + \boldsymbol{\delta}'\mathbf{x}_i. \end{cases}$$

There is an ambiguity in the model as specified. In principle, the model for three outcomes has two thresholds,  $\mu_0$  and  $\mu_1$ . It is always necessary to normalize the first,  $\mu_0 = 0$ . Therefore, the model implies the following probabilities:

$$\begin{aligned}
\text{Prob}(y = 0|\mathbf{x}) &= \Phi(-\alpha - \boldsymbol{\beta}'\mathbf{x}) &&= 1 - \Phi(\alpha_0 + \boldsymbol{\beta}_0'\mathbf{x}) \\
\text{Prob}(y = 1|\mathbf{x}) &= \Phi(\mu + \boldsymbol{\delta}'\mathbf{x}_i - \alpha - \boldsymbol{\beta}'\mathbf{x}) - \Phi(-\alpha - \boldsymbol{\beta}'\mathbf{x}) &&= \Phi(\alpha_0 + \boldsymbol{\beta}_0'\mathbf{x}) - \Phi(\alpha_1 + \boldsymbol{\beta}_1'\mathbf{x}) \\
\text{Prob}(y = 2|\mathbf{x}) &= \Phi(\alpha + \boldsymbol{\beta}'\mathbf{x} - \mu - \boldsymbol{\delta}'\mathbf{x}) &&= \Phi(\alpha_1 + \boldsymbol{\beta}_1'\mathbf{x})
\end{aligned}$$

where  $\alpha_0 = \alpha$ ,  $\boldsymbol{\beta}_0 = \boldsymbol{\beta}$ ,  $\alpha_1 = \alpha - \mu$ ,  $\boldsymbol{\beta}_1 = (\boldsymbol{\beta} - \boldsymbol{\delta})$ . This is precisely Williams's (2006) "Generalized Ordered Probit Model." That is, at this juncture, Terza's heterogeneous thresholds model and the generalized ordered probit model are indistinguishable. For direct applications of Terza's approach, see, e.g., Kerkhofs and Lindeboom (1995) and Lindeboom and van Doorslaer (2003).

The result carries over generically to the generalized ordered logit and probit models examined earlier. The motivation in these earlier instances, was to work around the parallel regressions assumption. The model specified is

$$\text{Prob}(y_i = j | \mathbf{x}_i) = F(\mu_j - \boldsymbol{\beta}_j'\mathbf{x}_i) - F(\mu_{j-1} - \boldsymbol{\beta}_{j-1}'\mathbf{x}_i).$$

Ostensibly, the generalization to allow a different parameter vector for each outcome. Boes and Winkelmann (2006a, 2006b) proposed the same model, motivated by the single crossing feature of the restricted model. But, when the regressor vector is the same in each cell, the implied "generalized threshold model"

$$\mu_{ij} = \mu_j + \boldsymbol{\gamma}_j'\mathbf{x}_i.$$

is also indistinguishable from the model with an outcome specific parameter vector;

$$\boldsymbol{\beta}_j = \boldsymbol{\gamma}_j - \boldsymbol{\beta}.$$

We can deduce a comparison of the two models from Terza's results. Terza reports results for a model with five regressors,  $\mathbf{x} = (\text{CFIE}, \text{LTIA}, \text{NIIA}, \text{TA}, \text{CVIA})$ . The numerical results in Table 12 are reported in the article (reported estimated standard errors are omitted):

**Table 12 Estimated Generalized Ordered Probit Models**

	Ordered Probit	Generalized Ordered Probit		Sample
	$\boldsymbol{\beta}$	$\boldsymbol{\beta}$	$\boldsymbol{\delta}$	Mean
Constant	-2.779	-17.862	-28.617	1.000
$x_1$	0.604	1.305	2.831	3.069
$x_2$	3.642	17.788	11.007	0.447
$x_3$	16.079	124.518	167.130	0.056
$x_4$	0.0012	0.0007	0.0009	1490.762
$x_5$	2.865	3.893	10.282	0.176
$[\mu]$	[1.955]	2.419		

The estimated value of  $\mu$  is not reported, but we should be able to approximate it. The sample consists of 222 observations in which the sample counts are 39, 100, 83, so the proportions are  $P_0 = 0.176$ ,  $P_1 = 0.450$ ,  $P_2 = 0.374$ , respectively. For the middle cell, at least approximately, at the means of the data, we should have,

$$P_1 \approx \Phi[\mu - (a + \mathbf{b}'\bar{\mathbf{x}})] - \Phi[-(a + \mathbf{b}'\bar{\mathbf{x}})]$$

The index function evaluated at the means is approximately 2.026. Using 0.45 for  $P_1$  and the inverse normal function, we obtain a value of  $\mu$  of approximately 1.955. The log likelihood values are not reported so it is not possible to compare the two models directly. In the

generalized model, the index function evaluated at the means is 2.796. Note that the coefficients have changed wildly; the second has increased by a factor of 4 and the third by a factor of 10. However, when we compute these at the sample means of the data, we find the index function is 2.796 compared to 2.026 previously and the implied threshold value is 2.419 compared to 1.955. Thus, the changes in the model are fairly moderate. The three predicted probabilities evaluated at the means are (.021382,.450315,.528302) for the first model and (.002587,.340501,.646909) for the second. (The model would not impose that these mimic the sample, even at the means, as it would in a multinomial (unordered) logit model, so these differences from the sample proportions are to be expected.) The very large swings in the parameter estimates attest to the need to use partial effects to scale them for comparisons across models.

Terza notes (p. 6) that the model formulation does not impose an ordering on the threshold coefficients. He suggests an inequality constrained maximization of the log likelihood, which is likely to be extremely difficult if there are many variables in  $\mathbf{x}$ . As a “less rigorous but apparently effective remedy,” he proposes to drop from the model variables in the threshold equations that are insignificant in the initial (unconstrained) model.

The analysis of this model continues with Pudney and Shields’s (2000) “Generalized Ordered Probit Model,” [also “(2)”] whose motivation, like Terza’s was to accommodate observable individual heterogeneity in the threshold parameters as well as in the mean of the regression. (Pudney and Shields discuss a clear example in the context of job promotion in which the steps on the promotion ladder for nurses are somewhat individual specific. In their setting, in contrast to Terza’s, the variables in the threshold equations are explicitly different from those in the regression. We (and Pudney and Shields) note an obvious problem of identification in this specification. Consider the generic probability with their extension,

$$\text{Prob}[y_i \leq j \mid \mathbf{x}_i, \mathbf{z}_i] = F(\mu_j + \boldsymbol{\delta}'\mathbf{z}_i - \boldsymbol{\beta}'\mathbf{x}_i) = F[\mu_j - (\boldsymbol{\delta}^*\mathbf{z}_i + \boldsymbol{\beta}'\mathbf{x}_i)], \boldsymbol{\delta}^* = -\boldsymbol{\delta}.$$

It is less than obvious whether the variables  $\mathbf{z}_i$  are actually in the threshold or in the mean of the regression. Either interpretation is consistent with the model. Pudney and Shields argue that the distinction is of no substantive consequence for their analysis.

## 5.2.2 Nonlinear Specifications – A Hierarchical Ordered Probit Model

The linearity of the regression specification has presented two significant obstacles to building the model. It has rendered indistinguishable the heterogeneous thresholds case and the “generalized” model that has heterogeneous parameter vectors. Second, it has produced a model that will be internally inconsistent at least for some data vectors; that is, it cannot ensure that the probabilities are always positive. One might consider modifying the thresholds directly. Greene (2007a) proposes a “Hierarchical Ordered Probit Model,”

$$y_i^* = \boldsymbol{\beta}'\mathbf{x}_i + \varepsilon_i$$

$$y_i = j \text{ if } \mu_{i,j-1} \leq y_i^* < \mu_{ij}$$

$$\mu_0 = 0,$$

$$\mu_j = \exp(\theta_j + \boldsymbol{\delta}'\mathbf{z}_i) \quad [\text{Case 1}]$$

or 
$$\mu_j = \exp(\theta_j + \boldsymbol{\delta}_j'\mathbf{z}_i) \quad [\text{Case 2}].$$

[The choice of the term “Hierarchical” model might be unfortunate, as it conflicts with a large literature on random parameter models such as the one discussed in Section 5.2.5, which is a “Hierarchical” model in the sense used in that literature. See, e.g., Raudenbush and Bryk (2002).]

Note that case 2 is the Terza(1985) and Pudney and Shields (2000) model with the exponential rather than linear function for the thresholds. It is, however, strongly distinct from Williams's model. This formulation addresses two problems; (i) the thresholds are mathematically distinct from the regression; (ii) by this construction, the threshold parameters must be positive. With a slight modification, to be pursued later, the ordering of the thresholds can also be assured; For the first case, for example, one might use

$$\mu_j = [\exp(\theta_1) + \exp(\theta_2) + \dots + \exp(\theta_j)] \times \exp(\delta'z)$$

and, in the second,

$$\mu_j = \mu_{j-1} + \exp(\theta_j + \delta_j'z_i)$$

In practical terms, the model can now be fit with the constraint that all predicted probabilities are greater than zero. This is a numerical solution to the problem of ordering the thresholds for all data vectors.

This model is a template case of *identification through functional form*. The contemporary literature views with some skepticism models that are unidentified without a change in functional form such as shown above. On the other hand, while this is true, it is also true that the underlying theory of the model does not insist on linearity of the thresholds (or the regression model, for that matter), and one might equally criticize the original model for being unidentified *because the model builder insists on a linear form*. That is, there is no obvious reason that the threshold parameters must be linear functions of the variables, or that linearity enjoys some claim to first precedence in the regression function. Of course, this is a methodological issue that cannot be resolved here.

The partial effects in this model are more involved than have been considered thus far. The Case 2 model implies

$$\text{Prob}(y = j | \mathbf{x}, \mathbf{z}) = F(\mu_j - \beta'x) - F(\mu_{j-1} - \beta'x).$$

Thus,

$$\frac{\partial \text{Prob}(y = j | \mathbf{x}, \mathbf{z})}{\partial \mathbf{x}} = [f(\mu_{j-1} - \beta'x) - f(\mu_j - \beta'x)]\beta$$

$$\frac{\partial \text{Prob}(y = j | \mathbf{x}, \mathbf{z})}{\partial \mathbf{z}} = f(\mu_j - \beta'x)\mu_j\delta_j - f(\mu_{j-1} - \beta'x)\mu_{j-1}\delta_{j-1}$$

(An obvious restriction is imposed if Case 1 applies.) If a variable appears in both  $\mathbf{x}$  and  $\mathbf{z}$ , then the two effects are added. It is clear on inspection, that this formulation has also circumvented the parallel regressions restriction, the single crossing feature and, with separate  $\delta$  vectors, the restriction that ratios of partial effects be the same for all outcomes. We conclude that at least as regards the question of functional form, the assumption of linearity has imposed a heavy cost on the construction of the model.

Numerically, the formulation shares the problem that its predecessors have. Without constraints or the modification suggested earlier, it does not impose the ordering of the threshold parameters. This, in the general form, unordered thresholds remain a possibility. We have found that the problem seems not to arise very often. As before, starting the iterations at the basic ordered probit or logit model estimates begins the process with a model in which all probabilities are positive. (At the starting values,  $\theta_j = \log \mu_j$  from the simple model.) As the iterations move the parameters away from the starting values, estimates that move the probabilities toward the proscribed regions begin to impose a heavy penalty on the log likelihood. As before, this appears generally to characterize the optimization process – it is, of course, not a prescription for how to

carry it out. We do note, it places a large value on a search method with a sensitive line search – a crude method such as Newton’s method (which uses none) is likely to fail early on.

Table 13 presents estimates of the ordered probit models using the same formulation as we used earlier. We have modeled the thresholds in terms of *INCOME*, *AGE* and *HANDDUM*, a dummy variable that indicates whether the individual reports a physical handicap. The table at the top of the listing shows that each successive generalization of the model brings a significant improvement in the log likelihood – the hypothesis of the restrictions of the preceding model is decisively rejected in all three cases (even if the significance level is adjusted for the sequential testing procedure). This seems consistent with the results found earlier for the Generalized (1) model. There is also a sizable increase (50%) in the Pseudo- $R^2$ , which we will explore in Table 15. The estimated coefficients in the index function seem to be relatively stable, save for the coefficient on *INCOME*, which increases substantially as the restrictions of the model are relaxed. This is consistent with the findings reported by Boes and Winkelmann (2006a). It is a bit less surprising when we recall that our data are drawn from the same data base, the GSOEP, as theirs. We may well be examining some of the same individuals.

Table 14 displays the partial effects for the three estimated models. Partial effects for the two binary variables that are marked with “\*” are computed by discrete changes in the probabilities with other variables held at their means. The effects are strikingly stable in spite of the changes in the coefficients from one model to the next. Table 15 suggests the payoff to the generalization. The prediction is the most probable cell computed at the individual observation. The counts of correct predictions for each model are shown in boldface/underline in the table. The effect of the generalization as one moves from left to right is to predict fewer values with  $y = 2$  correctly, but more with  $y = 1$ , and the difference is more than compensated. This would not predict the increase in the pseudo  $R^2$  seen in Table 13, but it is consistent with it.

**Table 13. Estimated Hierarchical Ordered Probit Models**

	No Model	Ordered Probit	HO-Case 1	HO-Case 2
Log likelihood function	-5875.096	-5752.985	-5690.804	-5665.088
Degrees of Freedom		5	3	6
Chi squared test of restr.	0	244.222	124.362	51.342
Info. Criterion: AIC	2.62284	2.57059	2.54419	2.53540
McFadden Pseudo R-squared	0.00000	.0207847	.0313684	.0357455

Variable	Coefficient	Standard Error	b/St.Er.	P[ Z >z]	Mean of X
-----+Index function for probability: Ordered Probit Model					
Constant	1.97882431	.11616998	17.034	.0000	
AGE	-.01807622	.00161885	-11.166	.0000	43.4401071
EDUC	.03556164	.00713213	4.986	.0000	11.4180864
INCOME	.25868689	.10387504	2.490	.0128	.34874007
MARRIED	-.03099645	.04203080	-.737	.4608	.75217488
KIDS	.06064631	.03823694	1.586	.1127	.37943342
-----+Threshold parameters for index					
Mu(1)	1.14834948	.02115847	54.274	.0000	
Mu(2)	2.54781466	.02161803	117.856	.0000	
Mu(3)	3.05638664	.02646225	115.500	.0000	
-----+Index function for probability: HOPIT, Case 1 Model					
Constant	2.02991736	.15908675	12.760	.0000	
AGE	-.02181435	.00254547	-8.570	.0000	43.4401071
EDUC	.03439864	.00738856	4.656	.0000	11.4180864
INCOME	.73490389	.15007529	4.897	.0000	.34874007
MARRIED	-.04347631	.04077144	-1.066	.2863	.75217488
KIDS	.05451516	.03887413	1.402	.1608	.37943342
-----+Estimates of t(j) in $\mu(j)=\exp[t(j)+d*j]$					
Theta(1)	.19501197	.06188197	3.151	.0016	
Theta(2)	.99052349	.05506325	17.989	.0000	
Theta(3)	1.17234814	.05418979	21.634	.0000	
-----+Threshold covariates $\mu(j)=\exp[t(j)+d*j]$					
AGE	-.00387305	.00107526	-3.602	.0003	
INCOME	.28304762	.05715368	4.952	.0000	
HANDDUM	.32483190	.02350983	13.817	.0000	
-----+Index function for probability: HOPIT Case 2 Model					
Constant	1.93645831	.18686446	10.363	.0000	
AGE	-.02120820	.00312110	-6.795	.0000	43.4401071
EDUC	.03404315	.00740298	4.599	.0000	11.4180864
INCOME	.94316189	.17336761	5.440	.0000	.34874007
MARRIED	-.04583214	.04103167	-1.117	.2640	.75217488
KIDS	.05089910	.03897042	1.306	.1915	.37943342
-----+Estimates of t(j) in $\mu(j)=\exp[t(j)+d(j)*z]$					
Theta(1)	.15119553	.12956102	1.167	.2432	
Theta(2)	.90995297	.06775050	13.431	.0000	
Theta(3)	1.18370905	.05889704	20.098	.0000	
-----+Threshold covariates $\mu(j)=\exp[t(j)+d(j)*z]$ . d(j) in sets of					
d1_AGE	-.00573712	.00250851	-2.287	.0222	
d1_INCOM	.54874525	.12216824	4.492	.0000	
d1_HANDD	.50575608	.04212225	12.007	.0000	
d2_AGE	-.00218205	.00132337	-1.649	.0992	
d2_INCOM	.31524490	.06854023	4.599	.0000	
d2_HANDD	.25780603	.03499892	7.366	.0000	
d3_AGE	-.00453013	.00118676	-3.817	.0001	
d3_INCOM	.33529293	.05730749	5.851	.0000	
d3_HANDD	.18051687	.04065755	4.440	.0000	

**Table 14. Estimated Partial Effects for Ordered Probit Models**

Variable	Y=00	Y=01	Y=02	Y=03	Y=04
Ordered Probit Model					
AGE	.0017	.0045	-.0012	-.0022	-.0028
EDUC	-.0034	-.0089	.0024	.0042	.0056
INCOME	-.0248	-.0644	.0177	.0309	.0406
*MARRIED	.0029	.0077	-.0020	-.0037	-.0049
*KIDS	-.0057	-.0151	.0040	.0072	.0096
Hierarchical Ordered Probit Model: Case 1					
AGE	.0020	.0055	-.0016	-.0026	-.0032
EDUC	-.0031	-.0087	.0026	.0042	.0051
INCOME	-.0669	-.1860	.0548	.0888	.1093
*MARRIED	.0039	.0110	-.0030	-.0053	-.0066
*KIDS	-.0049	-.0138	.0039	.0066	.0082
Hierarchical Ordered Probit Model: Case 1					
AGE	.0019	.0053	-.0015	-.0024	-.0034
EDUC	-.0031	-.0085	.0024	.0038	.0054
INCOME	-.0861	-.2363	.0666	.1065	.1493
*MARRIED	.0041	.0115	-.0030	-.0052	-.0074
*KIDS	-.0046	-.0127	.0035	.0058	.0081

**Table 15 Predicted Outcomes from Ordered Probit Models**

Cross tabulation of predictions. Row is actual, column is predicted. Predicted Outcome is the one with the largest probability.													
Model	0	1	1	0	1	2	0	1	2	0	1	2	
Actual	Row	Sum	y=0	y=1	y=1	y=2	y=2	y=2	y=3	y=3	y=4	y=4	
0	230	<u>0</u>	<u>0</u>	<u>0</u>	0	60	107	230	170	123	0	0	0
1	1113	0	0	0	<u>0</u>	<u>112</u>	<u>215</u>	1113	1001	898	0	0	0
2	2226	0	0	0	0	84	149	<u>2226</u>	<u>2142</u>	<u>2077</u>	0	0	0
3	500	0	0	0	0	2	10	500	498	490	<u>0</u>	<u>0</u>	<u>0</u>
4	414	0	0	0	0	5	10	414	409	404	0	0	<u>0</u>
Col Sum	4483	0	0	0	0	263	491	4483	4220	3992	0	0	0

### 5.2.3 Heterogeneous Scaling (Heteroscedasticity) of Random Utility

Considerably less attention has been focused on specification of the conditional variance in the regression model than on the conditional mean and the thresholds. In microeconomic data, scaling of the underlying preferences is surely as important a source of heterogeneity as displacement of the mean, perhaps even more so. But, it has received considerably less attention than heterogeneity in location. One would expect the problem of heterogeneity of the variance to be a persistent feature of individual level data. Researchers questioned its implications as early as Cox (1970). [See, also, Cox (1995).] Nonetheless, formal treatment of the issue is a relatively recent extension of the model.

A heteroscedastic ordered choice model is a minor extension of the basic model; the following form of the model based on Harvey (1976) appears in earlier versions of *LIMDEP* [Econometric Software (1997)] and *Stata* [Stata, Version 8] as a natural extension of the binary probit and logit models. The ordered choice model with heteroscedasticity would be

$$\begin{aligned}
 y_i^* &= \boldsymbol{\beta}'\mathbf{x}_i + \varepsilon_i \\
 y_i &= 0 \text{ if } \mu_{-1} < y_i^* \leq \mu_0, \\
 &= 1 \text{ if } \mu_0 < y_i^* \leq \mu_1, \\
 &= 2 \text{ if } \mu_1 < y_i^* \leq \mu_2 \\
 &= \dots \\
 &= J \text{ if } \mu_{J-1} < y_i^* \leq \mu_J.
 \end{aligned}$$

$$\text{Var}[\varepsilon_i|\mathbf{h}_i] \propto [\exp(\boldsymbol{\gamma}'\mathbf{h}_i)]^2$$

The model, itself is discussed in some detail in Williams (2006) and is also a feature of `GOLogit` and `GOProbit`. A search of the literature will turn up hundreds of recent applications of binary and ordered choice models with this form of heteroscedasticity [e.g., Hensher (2006)]. The binary probit and logit models with this form of heteroscedasticity are obvious extensions of the basic probit model, and appear, e.g., in Greene (1990) and Allison (1999).

Recall, at the outset of the discussion, it emerged that the lack of information on scaling of  $\varepsilon$  and therefore  $y^*$  is a signature feature of the ordered choice model. This same result will have major implications for building heteroscedasticity into the model. Consider the formulation of the model used in Chen and Khan (2003),

$$y_i^* = \boldsymbol{\beta}'\mathbf{x}_i + [\exp(\boldsymbol{\gamma}'\mathbf{h}_i)]\varepsilon_i$$

where  $\varepsilon_i$  is still  $N[0,1]$ . It follows that the observation mechanism is now

$$\text{Prob}(y_i = j | \mathbf{x}_i, \mathbf{h}_i) = F\left(\frac{\mu_j - \boldsymbol{\beta}'\mathbf{x}_i}{\exp(\boldsymbol{\gamma}'\mathbf{h}_i)}\right) - F\left(\frac{\mu_{j-1} - \boldsymbol{\beta}'\mathbf{x}_i}{\exp(\boldsymbol{\gamma}'\mathbf{h}_i)}\right).$$

This straightforward extension of the model should bring a substantive improvement in the correspondence of the model to the underlying data. Greene (2007a) proposes to blend this model with the hierarchical model of the previous section. The resulting functional form,

$$\text{Prob}(y_i = j | \mathbf{x}_i, \mathbf{h}_i) = F\left(\frac{\exp(\theta_j + \boldsymbol{\delta}'_j \mathbf{z}_i) - \boldsymbol{\beta}' \mathbf{x}_i}{\exp(\boldsymbol{\gamma}' \mathbf{h}_i)}\right) - F\left(\frac{\exp(\theta_{j-1} + \boldsymbol{\delta}'_{j-1} \mathbf{z}_i) - \boldsymbol{\beta}' \mathbf{x}_i}{\exp(\boldsymbol{\gamma}' \mathbf{h}_i)}\right).$$

should be intricate enough to overcome the parallel regressions, single crossing and constant ratios features of the basic model.

Unlike the linear regression case, unaccounted for heteroscedasticity is potentially disastrous for estimation of the parameters in the model. In the presence of latent heteroscedasticity that involves the variables that are in the model, or variables that are correlated with the variables in the model, the maximum likelihood estimator will be inconsistent, potentially seriously so. It is easy to see why in the formulation above. Unlike the linear regression model, in which latent heteroscedasticity will merely taint the standard errors, in the ordered (and binary) choice model, it will masquerade as a change in the functional form. Consider the model above, which can be written in equivalent form

$$y_i^{**} = \boldsymbol{\beta}' \mathbf{x}_i / [\exp(\boldsymbol{\gamma}' \mathbf{h}_i)] + \varepsilon_i$$

$$y_i = j \text{ if } \mu_{i,j} \leq y_i^{**} < \mu_{i,j+1}$$

where  $\varepsilon_i \sim N[0,1]$

but  $\mu_{i,j} = \mu_j / [\exp(\boldsymbol{\gamma}' \mathbf{h}_i)]$ ,

That is, the equivalent form of the model is one with a highly nonlinear conditional mean function and heterogeneous thresholds. Recall, the data contain no independent information on scaling of the underlying variable – any such information is determined from the conditional means and the functional form adopted for the variance. Estimating the model as if the disturbance were homoscedastic ignores both of these facts. Note that computing a “robust” covariance matrix for the estimator does nothing to redeem it. The estimator is inconsistent, potentially seriously so; the robust covariance matrix estimator is a moot point. Keele and Park (2005) have examined this model and its implications for bias in estimation. Chen and Khan (2003) have reconsidered the estimation of this model using robust methods that allow estimation of  $\boldsymbol{\beta}$  even in the presence of heteroscedasticity. (We note, estimation of  $\boldsymbol{\beta}$  solves only part of the model builder’s problem. If the measured outcome takes more than three values, then partial effects will be required to make much sense of the estimates. Without information about the underlying variance, or the underlying distribution, the scaling needed for the transformation is not computable.

As in other cases, the modification of the model alters the partial effects. For this case (omitting the hierarchical probit effects), the marginal effects are

$$\frac{\partial \text{Prob}(y_i = j | \mathbf{x}_i, \mathbf{h}_i)}{\partial \mathbf{x}_i} = \left[ f\left(\frac{\mu_{j-1} - \boldsymbol{\beta}' \mathbf{x}_i}{\exp(\boldsymbol{\gamma}' \mathbf{h}_i)}\right) - f\left(\frac{\mu_j - \boldsymbol{\beta}' \mathbf{x}_i}{\exp(\boldsymbol{\gamma}' \mathbf{h}_i)}\right) \right] \exp(-\boldsymbol{\gamma}' \mathbf{h}_i) \mathbf{x}_i$$

$$\frac{\partial \text{Prob}(y_i = j | \mathbf{x}_i, \mathbf{h}_i)}{\partial \mathbf{h}_i} = \left[ \begin{array}{c} f\left(\frac{\mu_{j-1} - \boldsymbol{\beta}' \mathbf{x}_i}{\exp(\boldsymbol{\gamma}' \mathbf{h}_i)}\right) \left(\frac{\mu_{j-1} - \boldsymbol{\beta}' \mathbf{x}_i}{\exp(\boldsymbol{\gamma}' \mathbf{h}_i)}\right) \\ - f\left(\frac{\mu_j - \boldsymbol{\beta}' \mathbf{x}_i}{\exp(\boldsymbol{\gamma}' \mathbf{h}_i)}\right) \left(\frac{\mu_j - \boldsymbol{\beta}' \mathbf{x}_i}{\exp(\boldsymbol{\gamma}' \mathbf{h}_i)}\right) \end{array} \right] \mathbf{h}_i$$

For a variable that appears in both  $\mathbf{x}_i$  and  $\mathbf{h}_i$ , the two parts are added. In such a case, the interpretation of the element of  $\beta$  associated with a particular variable becomes even more ambiguous than before.

Table 16 displays the estimates of the heteroscedastic ordered probit model using our earlier specification but adding *INCOME*, *AGE* and gender (*FEMALE*) to the variance equation. The basic slope parameters are quite similar to the earlier model (shown in the lower panel of Table 16 for convenience) But, the evidence of heteroscedasticity with respect to age and income is statistically significant, both individually and using the likelihood ratio test for the larger model. (The value of chi squared with 3 degrees of freedom is  $2(5752.985-5741.624) = 22.722$ . The tabled critical value is 7.814, so on this basis and based on the individual tests, the hypothesis of homoscedasticity would be rejected. It seems likely that this is yet another possible explanation for the finding of the Brant test carried out earlier.

**Table 16 Estimated Heteroscedastic Ordered Probit Model**

-----+-----					
Ordered Probability Model					
Dependent variable HEALTH					
Log likelihood function: Hetero. Homosk.					
-5741.624 -5752.985					
Info. Criterion: AIC: 2.56686 2.57059					
+-----+-----					
Variable	Coefficient	Standard Error	b/St.Er.	P[ Z >z]	Mean of X
+-----+-----					
-----+Index function for probability					
Constant	2.19351352	.17779847	12.337	.0000	
AGE	-.01992862	.00212062	-9.398	.0000	43.4401071
EDUC	.03904511	.00801855	4.869	.0000	11.4180864
INCOME	.24987395	.08630652	2.895	.0038	.34874007
MARRIED	-.03056402	.04443990	-.688	.4916	.75217488
KIDS	.06977840	.04168753	1.674	.0942	.37943342
-----+Variance function					
INCOME	-.23586591	.06074855	-3.883	.0001	.34874007
FEMALE	.01676897	.02491178	.673	.5009	.48405086
AGE	.00371284	.00111258	3.337	.0008	43.4401071
-----+Threshold parameters for index					
Mu(1)	1.28169216	.08114704	15.795	.0000	
Mu(2)	2.80192157	.15915064	17.605	.0000	
Mu(3)	3.35086805	.18739919	17.881	.0000	
+-----+-----					
Ordered Probit with Homoscedastic Disturbances					
-----+Index function for probability: Ordered Probit Model					
Constant	1.97882431	.11616998	17.034	.0000	
AGE	-.01807622	.00161885	-11.166	.0000	43.4401071
EDUC	.03556164	.00713213	4.986	.0000	11.4180864
INCOME	.25868689	.10387504	2.490	.0128	.34874007
MARRIED	-.03099645	.04203080	-.737	.4608	.75217488
KIDS	.06064631	.03823694	1.586	.1127	.37943342
-----+Threshold parameters for index					
Mu(1)	1.14834948	.02115847	54.274	.0000	
Mu(2)	2.54781466	.02161803	117.856	.0000	
Mu(3)	3.05638664	.02646225	115.500	.0000	

Table 17 displays the partial effects from both the restricted model and the heteroscedastic model. The latter are decomposed into the mean effects ( $\partial P(\cdot)/\partial \mathbf{x}$ ), the variance effects, ( $\partial P(\cdot)/\partial \mathbf{h}$ ) and the total equal to the sum of the two. The parts are marked; the total effects are shown in boldface. The partial effects from the restricted model are shown in parentheses for comparison. In contrast to the raw coefficients, the partial effects have shown some fairly substantial changes. The effects of *AGE* and *INCOME* are quite different (and changes sign twice), while the partial effects for *EDUC*, *MARRIED* and *KIDS* are quite similar to their earlier values.

**TABLE 17 Partial Effects in Heteroscedastic Ordered Probit Model**

Marginal Effects for Ordered Probit						
Variable	HEALTH=0	HEALTH=1	HEALTH=2	HEALTH=3	HEALTH=4	
AGE	.00169	.00463	-.00128	-.00216	-.00288	Mean
AGE	.00618	.00103	-.01647	.00086	.00839	Variance
<b>AGE</b>	<b>.00787</b>	<b>.00566</b>	<b>-.01775</b>	<b>-.00130</b>	<b>.00551</b>	<b>Total</b>
(AGE)	(.0017)	(.0045)	(-.0012)	(-.0022)	(-.0028)	Restricted
EDUC	-.00332	-.00906	.00251	.00423	.00564	Total
(EDUC)	(-.0034)	(-.0089)	(.0024)	(.0042)	(.0056)	restricted
INCOME	-.02122	-.05800	.01607	.02704	.03611	Mean
INCOME	.34732	.05785	-.92501	.04858	.47126	Variance
<b>INCOME</b>	<b>.32610</b>	<b>-.00015</b>	<b>-.90894</b>	<b>.07562</b>	<b>.50737</b>	<b>Total</b>
(INCOME)	(-.0248)	(-.0644)	(.0177)	(.0309)	(.0406)	Restricted
MARRIED	.00260	.00709	-.00197	-.00331	-.00442	Total
(MARRIED)	(.0029)	(.0077)	(-.0020)	(-.0037)	(-.0049)	Restricted
KIDS	-.00593	-.01620	.00449	.00755	.01008	Total
(KIDS)	(-.0057)	(-.0151)	(.0040)	(.0072)	(.0096)	Restricted
Pure Variance Effect						
FEMALE	-.00316	-.00053	.00840	-.00044	-.00428	Total

## 5.2.4 Individually Heterogeneous Marginal Utilities

Greene (2002, 2008a) argues that the fixed parameter version of the ordered choice model (and more generally, many microeconomic specifications) do not adequately account for the underlying heterogeneity likely to be present in observed data. Further extensions of the ordered choice model presented there include full random parameters treatments and discrete approximations under the form of latent class, or finite mixture models. These two specific extensions are also listed by Boes and Winkelmann (2006a).

The preceding lists the received “generalizations” of the ordered choice model. (The many other modified ordered choice models, such as bivariate ordered choice models, models with sample selection, and zero inflation models, that appear elsewhere have not been mentioned, as they are proposed to deal with features of the data other than heterogeneity. We will describe some of them in the sections to follow.) In what follows, we will propose a formulation of the ordered choice model that relaxes the restrictions listed above but treats heterogeneity in a unified, internally consistent fashion. The model contains three points at which individual heterogeneity can substantively appear, in the random utility model (the marginal utilities), in the threshold parameters, and in the scaling (variance) of the random components. As argued above, this form of treatment seems more likely to capture the salient features of the data generating mechanism than the received “generalized ordered logit model.”

## 5.2.5 Random Parameters Models

Formal modeling of heterogeneity in the parameters as representing a feature of the underlying data, appears in Greene (2002) (version 8.0) and Boes and Winkelmann (2006), both of whom suggest a full random parameters (RP) approach to the model. In Boes and Winkelmann, however, it is noted that the nature of an RP specification induces heteroscedasticity, and could be modeled as such. The model would appear as follows:

$$\beta_i = \beta + \mathbf{u}_i$$

where  $\mathbf{u}_i \sim N[\mathbf{0}, \mathbf{\Omega}]$ .

### Implied Heteroscedasticity

Boes and Winkelmann's treatment of a zero constant term and a full set of threshold parameters will prove less convenient than including a constant in  $\mathbf{x}_i$  and setting  $\mu_0 = 0$ , instead. We will maintain the latter formulation used heretofore. Inserting the expression above in the latent regression model, we obtain

$$\begin{aligned} y_i^* &= \beta_i' \mathbf{x}_i + \varepsilon_i \\ &= \beta' \mathbf{x}_i + \varepsilon_i + \mathbf{x}_i' \mathbf{u}_i. \end{aligned}$$

The observation mechanism is the same as earlier. The result is an ordered probit model in which the disturbance has variance  $\text{Var}[\varepsilon_i + \mathbf{x}_i' \mathbf{u}_i] = 1 + \mathbf{x}_i' \mathbf{\Omega} \mathbf{x}_i$ ; that is, a heteroscedastic ordered probit model. The resulting model has

$$\text{Prob}[y_i \leq j \mid \mathbf{x}_i] = \text{Prob}[\varepsilon_i + \mathbf{x}_i' \mathbf{u}_i \leq \mu_j - \beta' \mathbf{x}_i] = F \left( \frac{\mu_j - \beta' \mathbf{x}_i}{\sqrt{1 + \mathbf{x}_i' \mathbf{\Omega} \mathbf{x}_i}} \right),$$

which, it is suggested, can be estimated by ordinary means, albeit with a new source of nonlinearity – the elements of  $\mathbf{\Omega}$  must now be estimated as well. (The authors' suggestion that this could be handled semiparametrically without specifying a distribution for  $\mathbf{u}_i$  is incorrect, because the resulting heteroscedastic ordered choice model as written above only preserves the standard normal form assumed if  $\mathbf{u}_i$  is normally distributed as well as  $\varepsilon_i$ .) They did not pursue this approach. This computation will present a series of difficulties owing to the need to force  $\mathbf{\Omega}$  to be a positive definite matrix. One cannot simply insert the function above into the log likelihood and be optimistic that the estimated unconstrained matrix will, indeed, stay positive definite. At worst, it will become indefinite and it will become impossible to compute the log likelihood. A standard remedy is to use a Cholesky decomposition of  $\mathbf{\Omega}$ ; write  $\mathbf{\Omega} = \mathbf{L} \mathbf{D}^2 \mathbf{L}'$  where  $\mathbf{D}$  is a diagonal matrix with strictly positive elements and  $\mathbf{L}$  is a lower triangular matrix with ones on the diagonal. The log likelihood is then maximized with respect to the elements of  $\mathbf{L}$  and  $\mathbf{D}$  in addition to  $\beta$  and  $\mu_1, \dots, \mu_{j-1}$ . This will preserve the positive definiteness of the implied covariance matrix. Elements of  $\mathbf{\Omega}$  can be deduced after estimation.

Partial effects in this model can be obtained by differentiating the probabilities as if the parts in the numerators and denominators are functions of different variables, then adding them. An expression for this result is given in Boes and Winkelmann (2006a). An application in the study was done under the assumption that  $\mathbf{\Omega}$  is diagonal, which then requires only that the variances of the random parameters be positive.

## Maximum Simulated Likelihood Estimation

Greene (2002, 2007a, 2008a,b) analyzes the same model, but estimates the parameters by maximum simulated likelihood. First, write the random parameters as

$$\boldsymbol{\beta}_i = \boldsymbol{\beta} + \mathbf{LD}\mathbf{w}_i$$

where  $\mathbf{w}_i$  has a multivariate standard normal distribution, and  $\mathbf{LD}^2\mathbf{L}' = \boldsymbol{\Omega}$ . The probability for an observation is

$$\begin{aligned} \text{Prob}(y_i = j | \mathbf{x}_i, \mathbf{w}_i) &= \left[ \Phi(\mu_j - \boldsymbol{\beta}'_i \mathbf{x}_i) - \Phi(\mu_{j-1} - \boldsymbol{\beta}'_i \mathbf{x}_i) \right] \\ &= \left[ \Phi(\mu_j - \boldsymbol{\beta}' \mathbf{x}_i - (\mathbf{LD}\mathbf{w}_i)' \mathbf{x}_i) - \Phi(\mu_{j-1} - \boldsymbol{\beta}' \mathbf{x}_i - (\mathbf{LD}\mathbf{w}_i)' \mathbf{x}_i) \right]. \end{aligned}$$

In order to maximize the log likelihood, we must first integrate out the elements of the unobserved  $\mathbf{w}_i$ . Thus, the contribution to the unconditional log likelihood for observation  $i$  is

$$\log L_i = \log \int_{\mathbf{w}_i} \left[ \Phi(\mu_j - \boldsymbol{\beta}' \mathbf{x}_i - (\mathbf{LD}\mathbf{w}_i)' \mathbf{x}_i) - \Phi(\mu_{j-1} - \boldsymbol{\beta}' \mathbf{x}_i - (\mathbf{LD}\mathbf{w}_i)' \mathbf{x}_i) \right] F(\mathbf{w}_i) d\mathbf{w}_i.$$

The log likelihood for the sample is then the sum over the observations. Computing the integrals is an obstacle that must now be overcome. It has been simplified considerably already by decomposing  $\boldsymbol{\Omega}$  explicitly in the log likelihood, so that  $F(\mathbf{w}_i)$  is the multivariate standard normal density. The *Stata* routine, GLAMM [Rabe-Hesketh, Skrondal and Pickles (2005)] that is used for some discrete choice models does the computation using a form of Hermite quadrature. An alternative, generally substantially faster method of maximizing the log likelihood is maximum simulated likelihood. The integration is replaced with a simulation over  $R$  draws from the multivariate standard normal population. The simulated log likelihood is, then

$$\log L_S = \sum_{i=1}^N \log \frac{1}{R} \sum_{r=1}^R \left[ \Phi(\mu_j - \boldsymbol{\beta}' \mathbf{x}_i - (\mathbf{LD}\mathbf{w}_{ir})' \mathbf{x}_i) - \Phi(\mu_{j-1} - \boldsymbol{\beta}' \mathbf{x}_i - (\mathbf{LD}\mathbf{w}_{ir})' \mathbf{x}_i) \right].$$

The simulations are speeded up considerably by using Halton draws [see Train (2003)] rather than random draws. Further details on this method of estimation are also given in Greene (2007b, 2008a). Partial effects and predicted probabilities must be simulated as well. For the partial effects,

$$\begin{aligned} \frac{\partial \text{Prob}(y_i = j | \mathbf{x}_i)}{\partial \mathbf{x}_i} &= \\ &= \int_{\mathbf{w}_i} \left[ \phi(\mu_{j-1} - \boldsymbol{\beta}' \mathbf{x}_i - (\mathbf{LD}\mathbf{w}_i)' \mathbf{x}_i) - \phi(\mu_j - \boldsymbol{\beta}' \mathbf{x}_i - (\mathbf{LD}\mathbf{w}_i)' \mathbf{x}_i) \right] (\boldsymbol{\beta} + \mathbf{LD}\mathbf{w}_i) F(\mathbf{w}_i) d\mathbf{w}_i. \end{aligned}$$

As in the earlier formulations, this is a scalar multiple of the main parameter vector. We use simulation to compute

$$Est. \frac{\partial \text{Prob}(y_i = j | \mathbf{x}_i)}{\partial \mathbf{x}_i} = \left\{ \frac{1}{R} \sum_{r=1}^R \left[ \phi(\hat{\mu}_{j-1} - \hat{\boldsymbol{\beta}}' \mathbf{x}_i - (\hat{\mathbf{L}} \hat{\mathbf{D}} \mathbf{w}_{ir})' \mathbf{x}_i) - \phi(\hat{\mu}_j - \hat{\boldsymbol{\beta}}' \mathbf{x}_i - (\hat{\mathbf{L}} \hat{\mathbf{D}} \mathbf{w}_{ir})' \mathbf{x}_i) \right] \right\} (\hat{\boldsymbol{\beta}} + \hat{\mathbf{L}} \hat{\mathbf{D}} \mathbf{w}_{ir}).$$

Table 18 gives the estimates of the random parameters model for our familiar specification. The estimator produces estimates of  $\mathbf{L}$  and  $\mathbf{D}$ . The implied estimate of  $\boldsymbol{\Omega}$  is given in Table 19 with the estimates of the square roots of the diagonal elements of  $\boldsymbol{\Omega}$  and the implied correlation matrix obtained by  $\boldsymbol{\Sigma}^{-1} \boldsymbol{\Omega} \boldsymbol{\Sigma}^{-1}$  where  $\boldsymbol{\Sigma}$  is a diagonal matrix containing the estimated standard deviations from  $\boldsymbol{\Omega}$ . The estimates of the partial effects are shown in Table 20 with their counterparts from the basic model. A likelihood ratio test of the null hypothesis that the basic model applies against the alternative of this generalization is based on a chi squared statistic of  $2(5752.985 - 5705.592) = 94.786$  with 20 degrees of freedom. The null hypothesis would be rejected.

### Conditional Mean Estimation in the Random Parameters Model

The random parameters model is couched in terms of  $(\boldsymbol{\beta}_i, \mu_1, \dots, \mu_{J-1})$ , specific to the individual. Recall in the structure,

$$\boldsymbol{\beta}_i = \boldsymbol{\beta} + \mathbf{u}_i.$$

It would be useful to estimate  $\boldsymbol{\beta}_i$  rather than the population parameters,  $\boldsymbol{\beta}$ , if that were possible. It is not, of course, as that would require estimation of  $\mathbf{u}_i$  which is “noise.” However, in the same spirit as its Bayesian counterpart, one can compute an estimate of  $E[\boldsymbol{\beta}_i | y_i, \mathbf{x}_i]$ , which will contain more information than the natural, unconditional estimator,  $\boldsymbol{\beta}$ . The approach proceeds as follows: The density of  $y_i | \mathbf{x}_i, \boldsymbol{\beta}_i$  is

$$P(y_i | \mathbf{x}_i, \boldsymbol{\beta}_i) = \text{Prob}(y_i = j | \mathbf{x}_i, \boldsymbol{\beta}_i) = \left[ \Phi(\mu_j - \boldsymbol{\beta}_i' \mathbf{x}_i) - \Phi(\mu_{j-1} - \boldsymbol{\beta}_i' \mathbf{x}_i) \right].$$

The marginal density of  $\boldsymbol{\beta}_i$  assuming  $\mathbf{u}_i \sim N[\mathbf{0}, \boldsymbol{\Omega}]$  is  $N[\boldsymbol{\beta}, \boldsymbol{\Omega}]$ . The joint density of  $y_i$  and  $\boldsymbol{\beta}_i$  is

$$P(y_i, \boldsymbol{\beta}_i | \mathbf{x}_i) = \text{Prob}(y_i = j | \mathbf{x}_i, \boldsymbol{\beta}_i) P(\boldsymbol{\beta}_i).$$

Using Bayes Theorem, then,

$$\begin{aligned} P(\boldsymbol{\beta}_i | y_i, \mathbf{x}_i) &= \frac{P(y_i | \mathbf{x}_i, \boldsymbol{\beta}_i) P(\boldsymbol{\beta}_i)}{P(y_i)} \\ &= \frac{P(y_i | \mathbf{x}_i, \boldsymbol{\beta}_i) P(\boldsymbol{\beta}_i)}{\int_{\boldsymbol{\beta}_i} P(y_i | \mathbf{x}_i, \boldsymbol{\beta}_i) P(\boldsymbol{\beta}_i) d\boldsymbol{\beta}_i} \end{aligned}$$

The conditional mean is then,

$$E(\boldsymbol{\beta}_i | y_i, \mathbf{x}_i) = \frac{\int_{\boldsymbol{\beta}_i} \boldsymbol{\beta}_i P(y_i | \mathbf{x}_i, \boldsymbol{\beta}_i) P(\boldsymbol{\beta}_i) d\boldsymbol{\beta}_i}{\int_{\boldsymbol{\beta}_i} P(y_i | \mathbf{x}_i, \boldsymbol{\beta}_i) P(\boldsymbol{\beta}_i) d\boldsymbol{\beta}_i}$$

The integrals must be computed by simulation. The result is easily obtained as a byproduct of the estimation process. To see how, first insert the components of the probabilities, and replace the integration with simulation, as we did in computing the log likelihood. Then,

$$Est.E(\beta_i | y_i, \mathbf{x}_i) = \frac{\frac{1}{R} \sum_{r=1}^R \beta_{ir} \sum_{j=0}^J m_{ij} [\Phi(\mu_j - \beta'_{ir} \mathbf{x}_i) - \Phi(\mu_{j-1} - \beta'_{ir} \mathbf{x}_i)]}{\frac{1}{R} \sum_{r=1}^R \sum_{j=0}^J m_{ij} [\Phi(\mu_j - \beta'_{ir} \mathbf{x}_i) - \Phi(\mu_{j-1} - \beta'_{ir} \mathbf{x}_i)]}$$

where  $m_{ij} = 1(y_i = j)$ . The terms in square brackets are the simulated probabilities that enter the log likelihood. The draws on  $\beta_{ir}$  are obtained during the simulation; they are

$$\beta_{ir} = \beta + \mathbf{L} \mathbf{D} \mathbf{w}_{ir}.$$

That is, the same simulation that was done to maximize the log likelihood. It is illuminating to write this in a different form. Write this, using our final estimates of the model parameters,

$$0 < \hat{a}_{ir} = \frac{[\Phi(\hat{\mu}_j - \hat{\beta}'_{ir} \mathbf{x}_i) - \Phi(\hat{\mu}_{j-1} - \hat{\beta}'_{ir} \mathbf{x}_i)]}{\frac{1}{R} \sum_{r=1}^R [\Phi(\hat{\mu}_j - \hat{\beta}'_{ir} \mathbf{x}_i) - \Phi(\hat{\mu}_{j-1} - \hat{\beta}'_{ir} \mathbf{x}_i)]} < 1$$

where

$$\hat{\beta}_{ir} = \hat{\beta} + \hat{\mathbf{L}} \hat{\mathbf{D}} \mathbf{w}_{ir}.$$

Then, our estimator is

$$Est.E(\beta_i | y_i, \mathbf{x}_i) = \frac{1}{R} \sum_{r=1}^R \hat{a}_{ir} \hat{\beta}_{ir}$$

Other functions of the parameters, such as partial effects or probabilities for individual observations, could be simulated in the same way, just by replacing  $\beta_i$  with the desired function of  $\beta_i$  in the simulation.

Before illustrating the method, we emphasize two aspects of the computation. First, it must be borne in mind, this is not a direct estimator of  $\beta_i$ ; it is an estimator of the mean of the conditional distribution from which  $\beta_i$  is drawn. In the classical framework we are using here, this is as well as we can do, in terms of using the sample information, to estimate  $\beta_i$ . Second, this estimator is a counterpart to the Bayesian posterior mean, which would estimate the same parameters in the same way. A difference would be that the Bayesian posterior variance would be smaller than the variance of the conditional distribution if we computed it above. The reason is that our classical estimator uses the asymptotic distribution of the estimator while the Bayesian posterior mean is conditioned only on the observed sample. There is a degree of imprecision in the classical estimator above that is absent from the posterior mean because the simulations plug in the estimates of the parameters as if they were known, while the Bayesian counterpart is based on the exact, finite sample distribution of the estimators conditioned on the data in hand. This latter difference is likely to be extremely small in a sample as large as the one in use here. Figure 13 shows a kernel density estimator for the distribution of estimates of  $E[\beta_{INCOME} | y_i, \mathbf{x}_i]$  across the sample.

**Table 18 Estimated Random Parameters Ordered Probit Model**

Random Coefficients OrdProbs Model					
Dependent variable		HEALTH			
Number of observations		4483			
Log likelihood function		-5705.592			
Number of parameters		30			
Info. Criterion: AIC =		2.55882			
Ordered probit (normal) model					
LHS variable = values 0,1,..., 4					
Simulation based on 25 Halton draws					
Variable	Coefficient	Standard Error	b/St.Er.	P[ Z >z]	Mean of X
+-----+Means for random parameters					
Constant	3.21422***	.13661724	23.527	.0000	
AGE	-.02975***	.00181656	-16.379	.0000	43.440107
EDUC	.05994***	.00772363	7.760	.0000	11.418086
INCOME	.55843***	.11634679	4.800	.0000	.3487401
MARRIED	-.11403**	.04606014	-2.476	.0133	.7521749
KIDS	.11667***	.04188596	2.785	.0053	.3794334
+-----+Diagonal elements of Cholesky matrix					
Constant	1.59018***	.13023214	12.210	.0000	
AGE	.00021	.00131986	.155	.8765	
EDUC	.00311	.00401774	.773	.4394	
INCOME	.52568***	.08223135	6.393	.0000	
MARRIED	.18289***	.02664746	6.863	.0000	
KIDS	.25075***	.02857019	8.777	.0000	
+-----+Below diagonal elements of Cholesky matrix					
LAGE_ONE	.00809***	.00182840	4.427	.0000	
LEDU_ONE	-.02642***	.00777706	-3.397	.0007	
LEDU_AGE	-.02035***	.00542323	-3.753	.0002	
LINC_ONE	-1.56213***	.12129203	-12.879	.0000	
LINC_AGE	-1.37532***	.11797323	-11.658	.0000	
LINC_EDU	.99637***	.11897454	8.375	.0000	
LMAR_ONE	-.21299***	.04776366	-4.459	.0000	
LMAR_AGE	1.25275***	.04903494	25.548	.0000	
LMAR_EDU	.43215***	.04204805	10.277	.0000	
LMAR_INC	.57377***	.04023452	14.261	.0000	
LKID_ONE	-.72797***	.04382901	-16.609	.0000	
LKID_AGE	.33841***	.03991671	8.478	.0000	
LKID_EDU	-1.06815***	.04111858	-25.977	.0000	
LKID_INC	-.63225***	.03823865	-16.534	.0000	
LKID_MAR	-.14459***	.03709603	-3.898	.0001	
+-----+Threshold parameters for probabilities					
MU(1)	1.92835***	.05189777	37.157	.0000	
MU(2)	4.18431***	.06756770	61.928	.0000	
MU(3)	5.00177***	.07426448	67.351	.0000	
+-----+Note: ***, **, * = Significance at 1%, 5%, 10% level.					

**Table 19 Implied Estimates of Parameter Matrices**

$LD^2L = W$  = Implied covariance matrix of random parameters

```

2.5287
 0.0128715    6.55609e-005
-0.0420126   -0.00021803    0.00112193
-2.48407     -0.0129266    0.0723595    5.60087
-0.338695   -0.00146703   -0.0185281   -0.65802    2.16416
-1.1576     -0.00582305    0.00902707  -0.724866   -0.27181    2.26893
  
```

Square roots of diagonal elements of

$W$  = Implied standard deviations of random parameters

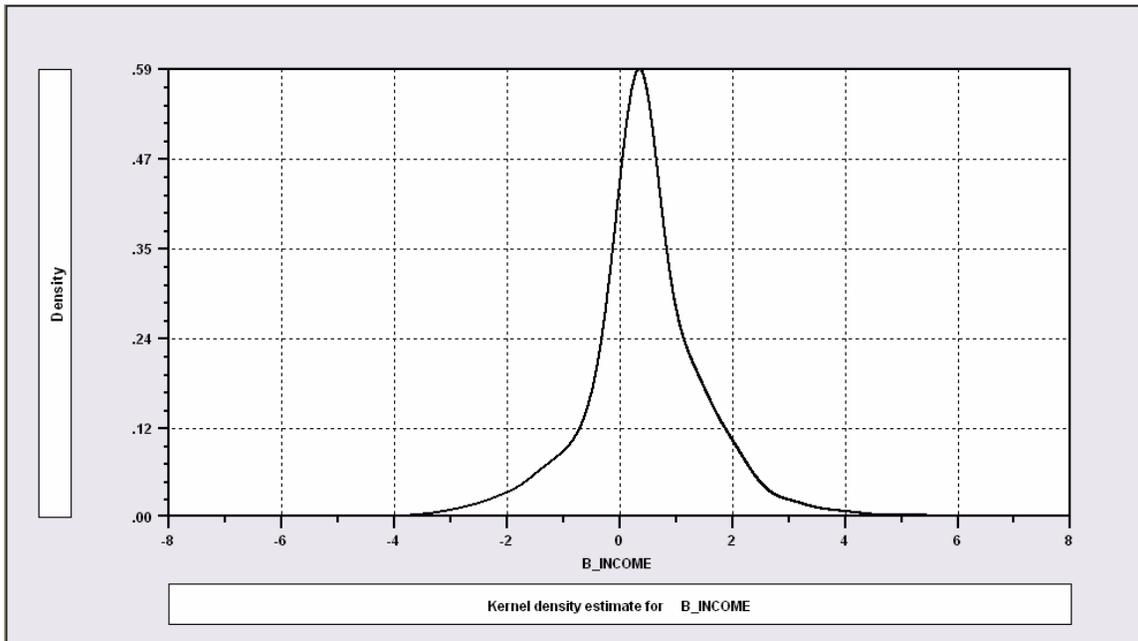
```

      1
+-----+
1 | 1.59018
2 | .00810
3 | .03350
4 | 2.36662
5 | 1.47111
6 | 1.50630
  
```

Implied correlation matrix of random parameters

```

1.000000
 0.999679    1.000000
-0.788771   -0.803913    1.000000
-0.660071   -0.674583    0.912818    1.000000
-0.144783   -0.123161   -0.376012   -0.189002    1.000000
-0.483286   -0.477438    0.178917   -0.203338   -0.122662    1.000000
  
```



**Figure 13 Kernel Density for Estimate of the Distribution of Means of Income Coefficient**

**Table 20 Estimated Partial Effects from Random Parameters Model**

=====  
	Summary of Marginal Effects for Ordered Probability Model (probit)	
	Effects computed at means. Effects for binary variables are	
	computed as differences of probabilities, other variables at means.	
 =====

Continuous Variable AGE			
Outcome	Effect	dPy<=nn/dX	dPy>=nn/dX
Y = 00	.00026	.00026	.00000
Y = 01	.00814	.00840	-.00026
Y = 02	-.00410	.00430	-.00840
Y = 03	-.00334	.00096	-.00430
Y = 04	-.00096	.00000	-.00096

Continuous Variable EDUC			Continuous Variable INCOME			
Outcome	Effect	dPy<=nn/dX	dPy>=nn/dX	Effect	dPy<=nn/dX	dPy>=nn/dX
Y = 00	-.00053	-.00053	.00000	-.00495	-.00495	.00000
Y = 01	-.01640	-.01693	.00053	-.15279	-.15774	.00495
Y = 02	.00827	-.00866	.01693	.07703	-.08071	.15774
Y = 03	.00673	-.00193	.00866	.06269	-.01803	.08071
Y = 04	.00193	.00000	.00193	.01803	.00000	.01803

Binary(0/1) Variable MARRIED			Binary(0/1) Variable KIDS			
Outcome	Effect	dPy<=nn/dX	dPy>=nn/dX	Effect	dPy<=nn/dX	dPy>=nn/dX
Y = 00	.00094	.00094	.00000	-.00100	-.00100	.00000
Y = 01	.03049	.03143	-.00094	-.03157	-.03256	.00100
Y = 02	-.01426	.01717	-.03143	.01535	-.01721	.03256
Y = 03	-.01324	.00393	-.01717	.01332	-.00390	.01721
Y = 04	-.00393	.00000	-.00393	.00390	.00000	.00390

**Estimated Partial Effects from Ordered Probit Model**

Continuous Variable AGE			
Outcome	Effect	dPy<=nn/dX	dPy>=nn/dX
Y = 00	.00173	.00173	.00000
Y = 01	.00450	.00623	-.00173
Y = 02	-.00124	.00499	-.00623
Y = 03	-.00216	.00283	-.00499
Y = 04	-.00283	.00000	-.00283

Continuous Variable EDUC			Continuous Variable INCOME			
Outcome	Effect	dPy<=nn/dX	dPy>=nn/dX	Effect	dPy<=nn/dX	dPy>=nn/dX
Y = 00	-.00340	-.00340	.00000	-.02476	-.02476	.00000
Y = 01	-.00885	-.01225	.00340	-.06438	-.08914	.02476
Y = 02	.00244	-.00982	.01225	.01774	-.07141	.08914
Y = 03	.00424	-.00557	.00982	.03085	-.04055	.07141
Y = 04	.00557	.00000	.00557	.04055	.00000	.04055

Binary(0/1) Variable MARRIED			Binary(0/1) Variable KIDS			
Outcome	Effect	dPy<=nn/dX	dPy>=nn/dX	Effect	dPy<=nn/dX	dPy>=nn/dX
Y = 00	.00293	.00293	.00000	-.00574	-.00574	.00000
Y = 01	.00771	.01064	-.00293	-.01508	-.02081	.00574
Y = 02	-.00202	.00861	-.01064	.00397	-.01684	.02081
Y = 03	-.00370	.00491	-.00861	.00724	-.00960	.01684
Y = 04	-.00491	.00000	-.00491	.00960	.00000	.00960

## 5.2.6 Latent Class and Finite Mixture Modeling

Latent class modeling [see McLachlan and Peel (2000)] provides an alternative approach to accommodating heterogeneity. [Applications include Everitt (1988) and Uebersax (1999).] The natural approach assumes that parameter vectors,  $\beta_i$  are distributed among individuals with a discrete distribution, rather than the continuous distribution of the previous section. Thus, it is assumed that the population consists of a finite number,  $Q$ , of groups of individuals. The groups are heterogeneous, with common parameters,  $\gamma_q = (\beta_q, \mu_q)$  for the members of the group, but the groups themselves are different from one another. The analyst does not know from the data which observation is in which class. (Hence the term *latent* classes.)

The model assumes that individuals are distributed heterogeneously with a discrete distribution in a population. Two other interpretations of the model are useful. The latent class model can also be viewed as a discrete approximation to the continuous distribution. This follows the development of Heckman and Singer (1984) who used this approach to modeling heterogeneity in a study of duration. Alternatively, the *finite mixture* model may be used as a technique to model the distribution in its own right. This technique is often used to mix normal distributions to obtain a non-normal mixture distribution, but it can be used more generally, as we show below.

### The Latent Class Ordered Choice Model

For modeling purposes, class membership is distributed with discrete distribution,

$$\text{Prob}(\text{individual } i \text{ is a member of class } = q) = \pi_{iq} = \pi_q$$

This statement needs its own interpretation. It can be given a long run frequency interpretation in that the probability that an individual drawn at random from the full population is a member of the particular class. Alternatively, it reflects the priors of the analyst over the same random outcome. Under either interpretation, then

$$\text{Prob}(y_i = j | \mathbf{x}_i) = \sum_q \text{Prob}(y_i = j | \mathbf{x}_i, \text{class} = q) \text{Prob}(\text{class} = q).$$

Combining terms from earlier, then, a latent class ordered probit model would be

$$\text{Prob}(y_i = j | \mathbf{x}_i) = \sum_{q=1}^Q \pi_q \left[ \Phi(\mu_{j,q} - \beta'_q \mathbf{x}_i) - \Phi(\mu_{j-1,q} - \beta'_q \mathbf{x}_i) \right].$$

(We will use the probit formulation for this discussion. A logit model is obtained trivially by changing the assumed cdf and density – it will be a simple change of notation.) By this construction, the implied estimator of the cell probabilities would be a mixture of the class specific probabilities, using the estimated class probabilities,  $\pi_q$  for the mixture. Likewise, the partial effects would be

$$\frac{\partial \text{Prob}(y_i = j | \mathbf{x}_i)}{\partial \mathbf{x}_i} = \sum_{q=1}^Q \pi_q \left[ \phi(\mu_{j-1,q} - \beta'_q \mathbf{x}_i) - \phi(\mu_{j,q} - \beta'_q \mathbf{x}_i) \right] \beta_q,$$

that is, the same weighted mixture of the class specific partial effects.

## Estimation by Maximum Likelihood

The estimation problem now includes estimation of  $(\boldsymbol{\beta}_q, \boldsymbol{\mu}_q, \pi_q), q = 1, \dots, Q$ . The class probabilities are estimated with the other parameters. It is necessary to force the class probabilities to be between zero and one and to sum to one. A convenient way to do so is to use a multinomial logit parameterization of the class probabilities.

$$\pi_q = \frac{\exp(\theta_q)}{\sum_{q=1}^Q \exp(\theta_q)}, \quad q = 1, \dots, Q, \quad \theta_Q = 0.$$

Assembling the parts, then, the full log likelihood for the parameters, given the observed data is

$$\log L = \sum_{i=1}^N \log \left\{ \sum_{q=1}^Q \frac{\exp(\theta_q)}{\sum_{q=1}^Q \exp(\theta_q)} \sum_{j=0}^J m_{ij} \left[ \Phi(\boldsymbol{\mu}_{j,q} - \boldsymbol{\beta}'_q \mathbf{x}_i) - \Phi(\boldsymbol{\mu}_{j-1,q} - \boldsymbol{\beta}'_q \mathbf{x}_i) \right] \right\}$$

where

$$m_{ij} = 1 \text{ if } y_i = j \text{ and } 0 \text{ otherwise, } j = 0, \dots, J; \quad i = 1, \dots, N,$$

and the full vector of parameters to be estimated is

$$\boldsymbol{\Theta} = (\boldsymbol{\beta}_1, \boldsymbol{\mu}_1, \dots, \boldsymbol{\beta}_Q, \boldsymbol{\mu}_Q, \theta_1, \dots, \theta_Q)$$

with several constraints,  $\mu_{-1,q} = -\infty$ ,  $\mu_{0,q} = 0$ ,  $\mu_{J,q} = +\infty$ ,  $q = 1, \dots, Q$  and  $\theta_Q = 0$ .

We have assumed to this point that the number of classes,  $Q$ , is known. This will rarely be the case, so a question naturally arises, how can the analyst determine  $Q$ ? Since  $Q$  is not a free parameter, a likelihood ratio test is not appropriate, though, in fact,  $\log L$  will increase when  $Q$  increases. Researchers typically use an information criterion, such as AIC, to guide them toward the appropriate value. (Heckman and Singer (1984) note a practical guidepost. If the model is fit with too many classes, then estimates will become imprecise, even varying wildly. Signature features of a model that has been overfit will be exceedingly small estimates of the class probabilities (see below), wild values of the structural parameters and huge estimated standard errors.)

Statistical inference about the parameters can be done in the familiar fashion. The Wald test or likelihood ratio tests will probably be more convenient. There are a couple cautions that should be borne in mind. Hypothesis tests across classes are unlikely to be meaningful. For example, suppose we fit a three class model. Tests about the equality of some of the coefficients in one class to those in another would probably be ambiguous, because the classes, themselves are indeterminate. It is rare that one can even put a name on the classes, other than, "1," "2," etc. Likewise, testing about the number of classes is an uncertain exercise. Consider our two class example below. If the parameters of the two classes are identical, it would seem that there is a single class. The number of restrictions would seem to be the number of model parameters. However, there remain two class probabilities,  $\pi_1$  and  $\pi_2$ . If the parameter vectors are the same, then regardless of the values of  $\pi_1$  and  $\pi_2$ , there is only one class. Thus, the degrees of freedom for this test are ambiguous. The same log likelihood will emerge for any pair of probabilities that sum to one.

## The EM Algorithm

The log likelihood can be maximized using conventional gradient methods. [See Econometric Software (2007).] An alternative method, the *EM* algorithm [Dempster, Laird and Rubin (1977)], is particularly well suited to latent class modeling. Though generally slower than gradient methods such as Broyden, Fletcher, Goldfarb and Shanno [see Greene (2008a)], the *EM* method does have the advantage of great stability.

The *EM* algorithm is most effective in estimating the parameters of “missing data models.” In the model we are examining, the missing data are

$$d_{iq} = 1 \text{ if individual } i \text{ is a member of class } q \text{ and } 0 \text{ if not.}$$

If  $d_{iq}$  were observed, then the “complete data” log likelihood could be written

$$\begin{aligned} \log L | \mathbf{d} &= \sum_{i=1}^N \sum_{q=1}^Q d_{iq} \log \sum_{j=0}^J m_{ij} \left[ \Phi(\mu_{j,q} - \beta'_q \mathbf{x}_i) - \Phi(\mu_{j-1,q} - \beta'_q \mathbf{x}_i) \right] \\ &= \sum_{q=1}^Q \sum_{i=1}^N d_{iq} \log \sum_{j=0}^J m_{ij} \left[ \Phi(\mu_{j,q} - \beta'_q \mathbf{x}_i) - \Phi(\mu_{j-1,q} - \beta'_q \mathbf{x}_i) \right] \\ &= \sum_{q=1}^Q \left\{ \sum_{i=1}^{N_q} \log \sum_{j=0}^J m_{ij} \left[ \Phi(\mu_{j,q} - \beta'_q \mathbf{x}_i) - \Phi(\mu_{j-1,q} - \beta'_q \mathbf{x}_i) \right] \right\} \end{aligned}$$

(In the second line, we have only reversed the order of the summations.) That is, if  $d_{iq}$  were known, then we could partition the log likelihood into separate log likelihoods for the  $Q$  classes and maximize each one separately. Maximization of this log likelihood would be done by separating the observations into the  $Q$  now known groups and estimating a separate ordered choice model for each group of  $N_q$  observations.

Since  $d_{iq}$  is not observed, we must maximize the earlier log likelihood instead. The *E* (expectation) step of the *EM* algorithm requires derivation of the expectation of  $\log L | \mathbf{d}$  given the observed data,  $y_i, \mathbf{x}_i, i=1, \dots, N$  and the parameters of the class specific models,  $\beta_q$  and  $\mu_q$ . This, in turn requires deriving  $E[d_{iq} | y_i, \mathbf{x}_i, \beta_q, \mu_q]$ . Unconditionally,  $E[d_{iq}] = \pi_q$ . However, there is more information in the sample. The conditional mean function,  $E[d_{iq} | y_i, \mathbf{x}_i, \beta_q, \mu_q]$  (which is the expectation conditioned on  $m_{ij}, \mathbf{x}_i$  and the parameters  $\beta_q, \mu_q$ ) is found as follows: the joint density of  $y_i$  and  $d_{iq}$  is

$$\begin{aligned} P(y_i, d_{iq} | \mathbf{x}_i, \beta_q, \mu_q) &= P(y_i | d_{iq}, \mathbf{x}_i, \beta_q, \mu_q) P(d_{iq}) \\ &= \text{Prob}(y_i = j | \text{class} = q, \mathbf{x}_i, \beta_q, \mu_q) \times \pi_q. \end{aligned}$$

Using Bayes Theorem,

$$\begin{aligned} P(d_{iq} | y_i, \mathbf{x}_i, \beta_q, \mu_q) &= \frac{P(y_i | d_{iq}, \mathbf{x}_i, \beta_q, \mu_q) P(d_{iq})}{P(y_i | \mathbf{x}_i, \beta_q, \mu_q)} \\ &= \frac{P(y_i | d_{iq}, \mathbf{x}_i, \beta_q, \mu_q) P(d_{iq})}{\sum_{q=1}^Q P(y_i | d_{iq}, \mathbf{x}_i, \beta_q, \mu_q) P(d_{iq})} \\ &= \frac{\left\{ \sum_{j=0}^J m_{ij} \left[ \Phi(\mu_{j,q} - \beta'_q \mathbf{x}_i) - \Phi(\mu_{j-1,q} - \beta'_q \mathbf{x}_i) \right] \right\} \pi_q}{\sum_{q=1}^Q \left\{ \sum_{j=0}^J m_{ij} \left[ \Phi(\mu_{j,q} - \beta'_q \mathbf{x}_i) - \Phi(\mu_{j-1,q} - \beta'_q \mathbf{x}_i) \right] \right\} \pi_q}. \end{aligned}$$

The conditional mean is, then,

$$\begin{aligned}
E(d_{iq} | y_i, \mathbf{x}_i, \boldsymbol{\beta}_q, \boldsymbol{\mu}_q) &= \sum_{q=1}^Q d_{iq} \frac{\left\{ \sum_{j=0}^J m_{ij} \left[ \Phi(\mu_{j,q} - \boldsymbol{\beta}'_q \mathbf{x}_i) - \Phi(\mu_{j-1,q} - \boldsymbol{\beta}'_q \mathbf{x}_i) \right] \right\} \pi_q}{\sum_{q=1}^Q \left\{ \sum_{j=0}^J m_{ij} \left[ \Phi(\mu_{j,q} - \boldsymbol{\beta}'_q \mathbf{x}_i) - \Phi(\mu_{j-1,q} - \boldsymbol{\beta}'_q \mathbf{x}_i) \right] \right\} \pi_q} \\
&= \frac{\left\{ \sum_{j=0}^J m_{ij} \left[ \Phi(\mu_{j,q} - \boldsymbol{\beta}'_q \mathbf{x}_i) - \Phi(\mu_{j-1,q} - \boldsymbol{\beta}'_q \mathbf{x}_i) \right] \right\} \pi_q}{\sum_{q=1}^Q \left\{ \sum_{j=0}^J m_{ij} \left[ \Phi(\mu_{j,q} - \boldsymbol{\beta}'_q \mathbf{x}_i) - \Phi(\mu_{j-1,q} - \boldsymbol{\beta}'_q \mathbf{x}_i) \right] \right\} \pi_q} \\
&= \hat{w}_{iq}
\end{aligned}$$

(There is only one nonzero term in the summation in the first line.) The  $M$  (maximization) step of the EM algorithm consists of maximizing  $E[\log L | \mathbf{d}]$  by replacing  $d_{iq}$  with the expectations derived above. (Note, we are conditioning on an existing (previous) value of  $(\boldsymbol{\beta}_q, \boldsymbol{\mu}_q)$ , so  $\hat{w}_{iq}$  is not a function of the parameters in the expected log likelihood.) Thus, the  $M$  step consists of maximizing

$$\log L | E = \sum_{q=1}^Q \sum_{i=1}^N \hat{w}_{iq} \log \sum_{j=0}^J m_{ij} \left[ \Phi(\mu_{j,q} - \boldsymbol{\beta}'_q \mathbf{x}_i) - \Phi(\mu_{j-1,q} - \boldsymbol{\beta}'_q \mathbf{x}_i) \right].$$

Since the weights,  $\hat{w}_{iq}$ , are now known, this maximand can be partitioned into  $Q$  separate weighted log likelihoods that can be maximized separately;

$$\log L | E, q = \sum_{i=1}^N \hat{w}_{iq} \log \sum_{j=0}^J m_{ij} \left[ \Phi(\mu_{j,q} - \boldsymbol{\beta}'_q \mathbf{x}_i) - \Phi(\mu_{j-1,q} - \boldsymbol{\beta}'_q \mathbf{x}_i) \right], q = 1, \dots, Q.$$

To assemble the parts, then, the  $EM$  algorithm for latent class modeling – this is a general template that we can use for our ordered choice model – is

- (1) Obtain starting values for  $\boldsymbol{\beta}_q, \boldsymbol{\mu}_q, q = 1, \dots, Q$ .
- (2) Compute weights  $\hat{w}_{iq}, i = 1, \dots, N$  based each of the  $q$  parameter vectors.
- (3) Using the weights obtained in step (2), compute  $Q$  new sets of parameters by maximizing  $Q$  separate weighted log likelihoods.
- (4) Return to step (2) if the new estimates are not sufficiently close to the previous ones. Otherwise, exit the iterations

As noted earlier, this algorithm can take many iterations. However, each iteration is simple. Adding weights to the log likelihood we have been manipulating all along is a trivial modification. Moreover, as shown by Dempster et al. (1977), the log likelihood increases with every iteration – that is the stability aspect that is not necessarily achieved by other gradient methods.

Before leaving this discussion of the EM algorithm, we note a few practical points: (1) It would be tempting to obtain the starting values by using for each class the single class estimates obtained by maximizing the log likelihood for the sample without the latent class structure. Unfortunately, this leads to a frustrating result. If the parameters in the classes are the same, then the sets of weights for the classes will also be the same, which means that the next set of parameter estimates will again be the same. The end result is that these starting values will prevent the iterations from ever reaching the solution. A practical expedient is a small, different perturbation of the original estimates for each class. (2) The EM algorithm finds the maximizer

of the log likelihood function, but unlike other gradient methods, it does not automatically produce an estimate of the asymptotic covariance matrix of the estimator. That must be obtained separately after the estimation is done. Note that the second derivatives matrix (or an approximation to it) computed from the weighted log likelihood function is not an appropriate estimator of the asymptotic covariance matrix of the class specific parameter vector. (3) To this point, we have not obtained an estimator of  $\pi_q$ . The appropriate estimator, perhaps not surprisingly, is

$$\hat{\pi}_q = \frac{\sum_{i=1}^N \hat{w}_{iq}}{N} = \bar{\hat{w}}_q$$

### Estimating the Class Assignments

There is a secondary estimation problem in the latent class setting, known as the “classification problem.” Ex post, it would be useful to be able to assign observations to classes. Of course, if we could do this, then the classes would not be latent, and the model would be superfluous. However, one’s best guess of the class from which observation  $i$  is drawn would be based on the posterior,

$$\text{Prob}(\text{individual } i \text{ is in class } q | y_i, \mathbf{x}_i, \boldsymbol{\beta}_q, \boldsymbol{\mu}_q) = \hat{w}_{iq}$$

as computed earlier. Thus, the EM algorithm provides the sample estimator for the classification problem automatically. If the EM algorithm has not been used, it is still possible to compute  $\hat{w}_{iq}$  using the estimated parameters, simply using the definition given earlier. The end result would be to estimate the class membership for individual  $i$  as that  $q$  associated with the maximum value of  $\hat{w}_{iq}$  for  $q = 1, \dots, Q$ .

### A Latent Class Model Extension

The latent class interpretation of the model suggests a useful extension of the class probabilities model. Thus far, the specification provides no prior information about the class membership. That is, the prior class probabilities are constants,

$$\text{Prob}(\text{class} = q) = \pi_q.$$

If there were sample information that were useful, though not definitive (in which case, the classes would not be latent) for determining class membership, then we might write

$$\text{Prob}(\text{class} = q | \mathbf{z}_i) = \pi_q(\mathbf{z}_i).$$

where presumably,  $\mathbf{z}_i$  does not appear in the main model. For example, in our ordered probit model, it might be suspected that gender or working status has an influence on the class probabilities for health satisfaction. This is straightforward to build into the multinomial logit model, in the form

$$\pi_{iq} = \frac{\exp(\theta_q + \boldsymbol{\delta}'_q \mathbf{z}_i)}{\sum_{q=1}^Q \exp(\theta_q + \boldsymbol{\delta}'_q \mathbf{z}_i)}, \quad q = 1, \dots, Q, \quad \theta_Q = 0, \quad \boldsymbol{\delta}_Q = \mathbf{0}.$$

Estimation is also only slightly more complicated. The log likelihood for the full model would now be

$$\log L = \sum_{i=1}^N \log \left\{ \sum_{q=1}^Q \frac{\exp(\theta_q + \boldsymbol{\delta}'_q \mathbf{z}_i)}{\sum_{q=1}^Q \exp(\theta_q + \boldsymbol{\delta}'_q \mathbf{z}_i)} \sum_{j=0}^J m_{ij} \left[ \Phi(\mu_{j,q} - \boldsymbol{\beta}'_q \mathbf{x}_i) - \Phi(\mu_{j-1,q} - \boldsymbol{\beta}'_q \mathbf{x}_i) \right] \right\}$$

The EM algorithm would require a slight modification. We would add a step (2a) that would be estimation of the logit parameters,  $(\theta_q, \boldsymbol{\delta}_q), q=1, \dots, Q-1$  (with  $\theta_Q = 0$  and  $\boldsymbol{\delta}_Q = \mathbf{0}$ ). This (sub)step is done by fitting a multinomial logit model to the weights,  $\hat{w}_{iq}$  based on proportions, rather than individual data. The implied log likelihood function is

$$\begin{aligned} \log L(\boldsymbol{\theta}, \boldsymbol{\Delta}) &= \sum_{i=1}^N \sum_{q=1}^Q \hat{w}_{iq} \log \left[ \frac{\exp(\theta_q + \boldsymbol{\delta}'_q \mathbf{z}_i)}{\sum_{q=1}^Q \exp(\theta_q + \boldsymbol{\delta}'_q \mathbf{z}_i)} \right] \\ &= \sum_{i=1}^N \sum_{q=1}^Q \hat{w}_{iq} \log \Lambda_{iq}. \end{aligned}$$

The solution to this estimation is also straightforward using Newton's method. The first order conditions are revealing of the structure of the problem;

$$\frac{\partial \log L(\boldsymbol{\theta}, \boldsymbol{\Delta})}{\partial \begin{pmatrix} \theta_q \\ \boldsymbol{\delta}_q \end{pmatrix}} = \sum_{i=1}^N (\hat{w}_{iq} - \Lambda_{iq}) \begin{pmatrix} 1 \\ \mathbf{z}_i \end{pmatrix}.$$

The first of the equations would imply  $\sum_i (\hat{w}_{iq} - \Lambda_{iq}) = 0$ . If there were no covariates,  $\mathbf{z}_i$  in the equation, this would return the original solution for  $\hat{\pi}_q$  that was shown earlier. Thus, we find (as expected) that the ordinary methods and the EM method find the same maximizer of the log likelihood.

## Application

Table 21 presents estimates of a two class latent class model using our base specification. The single class estimates are presented for comparison. The estimates for the two class model, as expected, bracket the one class estimates. Although the log likelihood has increased substantially (from -5752.985 to -5716.627), the class definition does not appear to have greatly changed the results. The estimated prior class probabilities are near 50%. In Table 22, we have listed estimates of an extended model in which gender (*FEMALE*), handicapped (*HANDDUM*) and work status (*WORKING*) enter the class probabilities. This modification does appear to add significantly to the class segregation. Evidently *HANDDUM* and *WORKING*, though not *FEMALE*, are significant determinants. The log likelihood for the extended model jumps to -5683.202. The chi squared for the extension is  $2(5716.627 - 5683.202) = 66.85$  with 3 degrees of freedom, which is also highly significant. Partial effects for the two models are shown in Table 23.

**Table 21 Estimated Two Class Latent Class Ordered Probit Model**

Latent Class / Panel OrdProbs Model					
Dependent variable	HEALTH				
Number of observations	4483				
Log likelihood function	-5716.627				
Info. Criterion: AIC =	2.55883				
Ordered probit (normal) model					
LHS variable = values 0,1,..., 4					
Model fit with 2 latent classes.					
Variable	Coefficient	Standard Error	b/St.Er.	P[ Z >z]	Mean of X
-----+Model parameters for latent class 1					
Constant	2.95021***	.41309506	7.142	.0000	
AGE	-.01199***	.00363216	-3.302	.0010	43.440107
EDUC	.00658	.01995384	.330	.7415	11.418086
INCOME	-.89315***	.32110301	-2.782	.0054	.3487401
MARRIED	-.00384	.08409903	-.046	.9635	.7521749
KIDS	-.06006	.08588745	-.699	.4844	.3794334
MU(1)	1.05941***	.20888397	5.072	.0000	
MU(2)	2.99143***	.24112250	12.406	.0000	
MU(3)	3.26386***	.19738957	16.535	.0000	
-----+Model parameters for latent class 2					
Constant	1.33844***	.31513662	4.247	.0000	
AGE	-.03136***	.00489766	-6.403	.0000	43.440107
EDUC	.07599***	.02139757	3.551	.0004	11.418086
INCOME	1.87668***	.48437067	3.874	.0001	.3487401
MARRIED	-.11059	.09619434	-1.150	.2503	.7521749
KIDS	.21815**	.10142270	2.151	.0315	.3794334
MU(1)	1.53999***	.16607984	9.273	.0000	
MU(2)	2.47627***	.16415343	15.085	.0000	
MU(3)	3.71906***	.34228520	10.865	.0000	
-----+Estimated prior probabilities for class membership					
Class1Pr	.57532***	.08598171	6.691	.0000	
Class2Pr	.42468***	.08598171	4.939	.0000	
-----+Ordered Probit with					
-----+Index function for probability: Ordered Probit Model					
Constant	1.97882431	.11616998	17.034	.0000	
AGE	-.01807622	.00161885	-11.166	.0000	43.4401071
EDUC	.03556164	.00713213	4.986	.0000	11.4180864
INCOME	.25868689	.10387504	2.490	.0128	.34874007
MARRIED	-.03099645	.04203080	-.737	.4608	.75217488
KIDS	.06064631	.03823694	1.586	.1127	.37943342
-----+Threshold parameters for index					
Mu(1)	1.14834948	.02115847	54.274	.0000	
Mu(2)	2.54781466	.02161803	117.856	.0000	
Mu(3)	3.05638664	.02646225	115.500	.0000	
-----+Note: ***, **, * = Significance at 1%, 5%, 10% level.					

**Table 22 Estimated Extended Latent Class Ordered Probit Model**

Latent Class / Panel OrdProbs Model					
Dependent variable		HEALTH			
Number of observations		4483			
Log likelihood function		-5683.202			
Info. Criterion: AIC =		2.54526			
Ordered probit (normal) model					
LHS variable = values 0,1,..., 4					
Model fit with 2 latent classes.					
Variable	Coefficient	Standard Error	b/St.Er.	P[ Z >z]	Mean of X
-----+Model parameters for latent class 1					
Constant	2.67403***	.98769977	2.707	.0068	
AGE	-.01683***	.00306469	-5.491	.0000	43.440107
EDUC	.05650***	.01406183	4.018	.0001	11.418086
INCOME	-.07221	.20533188	-.352	.7251	.3487401
MARRIED	-.12503*	.07141561	-1.751	.0800	.7521749
KIDS	.05849	.06950909	.841	.4001	.3794334
MU(1)	1.24272*	.72874720	1.705	.0881	
MU(2)	3.10037***	.97525584	3.179	.0015	
MU(3)	3.81235***	1.04497621	3.648	.0003	
-----+Model parameters for latent class 2					
Constant	1.78822***	.35857003	4.987	.0000	
AGE	-.02218***	.00500768	-4.428	.0000	43.440107
EDUC	.00627	.02981904	.210	.8335	11.418086
INCOME	.44730	.34811347	1.285	.1988	.3487401
MARRIED	-.06115	.11423990	-.535	.5924	.7521749
KIDS	.12432	.11100482	1.120	.2627	.3794334
MU(1)	1.45292***	.30105346	4.826	.0000	
MU(2)	2.39382***	.39210119	6.105	.0000	
MU(3)	2.39382***	.30766315	7.781	.0000	
-----+Estimated prior probabilities for class membership					
ONE_1	.73833	.76214396	.969	.3327	
FEMALE_1	-.04306	.12780401	-.337	.7362	
HANDDU_1	-1.22319***	.23890142	-5.120	.0000	
WORKIN_1	.40969***	.15117501	2.710	.0067	
ONE_2	.000***	.....(Fixed Parameter).....			
FEMALE_2	.000***	.....(Fixed Parameter).....			
HANDDU_2	.000***	.....(Fixed Parameter).....			
WORKIN_2	.000***	.....(Fixed Parameter).....			
Note: ***, **, * = Significance at 1%, 5%, 10% level.					
-----+Prior class probabilities at data means for LCM variables					
Class 1	Class 2	Class 3	Class 4	Class 5	
.70182	.29818	.00000	.00000	.00000	

**Table 23. Estimated Partial Effects from Latent Class Models**

```

=====
|| Summary of Marginal Effects for Ordered Probability Model (probit) ||
|| Effects computed at means. Effects for binary variables are ||
|| computed as differences of probabilities, other variables at means. ||
=====
||
|| Continuous Variable AGE
Outcome Effect dPy<=nn/dX dPy>=nn/dX
Y = 00 .00137 .00137 .00000
Y = 01 .00529 .00666 -.00137
Y = 02 -.00123 .00543 -.00666
Y = 03 -.00309 .00234 -.00543
Y = 04 -.00234 .00000 -.00234
=====
||
|| Continuous Variable EDUC Continuous Variable INCOME
Outcome Effect dPy<=nn/dX dPy>=nn/dX Effect dPy<=nn/dX dPy>=nn/dX
Y = 00 -.00244 -.00244 .00000 -.01919 -.01919 .00000
Y = 01 -.00943 -.01187 .00244 -.07405 -.09323 .01919
Y = 02 .00219 -.00968 .01187 .01720 -.07603 .09323
Y = 03 .00552 -.00417 .00968 .04332 -.03272 .07603
Y = 04 .00417 .00000 .00417 .03272 .00000 .03272
=====
||
|| Binary(0/1) Variable MARRIED Binary(0/1) Variable KIDS
Outcome Effect dPy<=nn/dX dPy>=nn/dX Effect dPy<=nn/dX dPy>=nn/dX
Y = 00 .00326 .00326 .00000 -.00389 -.00389 .00000
Y = 01 .01281 .01607 -.00326 -.01516 -.01904 .00389
Y = 02 -.00272 .01335 -.01607 .00335 -.01570 .01904
Y = 03 -.00756 .00579 -.01335 .00891 -.00679 .01570
Y = 04 -.00579 .00000 -.00579 .00679 .00000 .00679
=====
Partial Effects From Expanded Latent Class Model
=====
||
|| Continuous Variable AGE
Outcome Effect dPy<=nn/dX dPy>=nn/dX
Y = 00 .00092 .00092 .00000
Y = 01 .00203 .00294 -.00092
Y = 02 .00129 .00424 -.00294
Y = 03 -.00123 .00300 -.00424
Y = 04 -.00300 .00000 -.00300
=====
||
|| Continuous Variable EDUC Continuous Variable INCOME
Outcome Effect dPy<=nn/dX dPy>=nn/dX Effect dPy<=nn/dX dPy>=nn/dX
Y = 00 -.00331 -.00331 .00000 .00579 .00579 .00000
Y = 01 -.00733 -.01064 .00331 .01284 .01863 -.00579
Y = 02 -.00467 -.01531 .01064 .00818 .02681 -.01863
Y = 03 .00446 -.01085 .01531 -.00781 .01900 -.02681
Y = 04 .01085 .00000 .01085 -.01900 .00000 -.01900
=====
||
|| Binary(0/1) Variable MARRIED Binary(0/1) Variable KIDS
Outcome Effect dPy<=nn/dX dPy>=nn/dX Effect dPy<=nn/dX dPy>=nn/dX
Y = 00 .00688 .00688 .00000 -.00300 -.00300 .00000
Y = 01 .01566 .02254 -.00688 -.00668 -.00967 .00300
Y = 02 .01087 .03341 -.02254 -.00433 -.01400 .00967
Y = 03 -.00944 .02396 -.03341 .00405 -.00995 .01400
Y = 04 -.02396 .00000 -.02396 .00995 .00000 .00995
=====

```

## Endogenous Class Assignment and A Generalized Ordered Choice Model

Greene, Harris, Hollingsworth and Maitra (2008) analyzed obesity in a sample of 12,601 men and 15,259 women in the U.S. National Health Interview Survey from 2005. The central feature of their model is a three outcome ordered choice model for weight class defined as normal, overweight and obese. Obesity is measured by the World Health Organization's standard body mass index, or *BMI*. *BMI* is computed as the weight in Kg divided the square of the height in meters. Values under 18.5 are classified by WHO as underweight. The 2% of their sample in this class was deleted. The remaining three classes are normal (18.5,25], overweight (25,30] and obese, (30,∞). There are great differences across individuals in body fat and conditioning, and the *BMI* classification is at best only a loose categorization of the desired health level indicated. The authors reasoned that the latent regression model *with known thresholds* that might seem superficially to apply,

$$\begin{aligned}
 BMI^* &= \beta'x + \varepsilon, \varepsilon \sim N[0, \sigma^2] \\
 BMI &= 0 \text{ if } BMI^* \leq 25 \\
 &1 \text{ if } 25 < BMI^* \leq 30 \\
 &2 \text{ if } BMI^* > 30,
 \end{aligned}$$

would be too narrow, and would neglect several sources of heterogeneity. They opted instead for an ordered "choice" model, defined as

$$\begin{aligned}
 BMI^* &= \beta'x + \varepsilon, \varepsilon \sim N[0, 1] \\
 WT &= 0 \text{ if } BMI^* \leq 0 \\
 &1 \text{ if } 0 < BMI^* \leq \mu \\
 &2 \text{ if } BMI^* > \mu.
 \end{aligned}$$

A recent study in *Science* [Herbert, Gerry and McQueen (2006)] suggests that an obesity predisposing geno-type is present in 10% of individuals. In the sample, roughly 25% of the sample is categorized as obese. This suggests that a latent class model might be appropriate and that the class division depends on more than just this (unobserved) geno-type. The study used a two class model, with

$$\begin{aligned}
 class^* &= \alpha'w + u, u \sim N[0, 1]. \\
 class &= 0 \text{ if } class^* \leq 0 \\
 &1 \text{ if } class^* > 0.
 \end{aligned}$$

The rigidity of the BMI classification, itself, might have produced erroneous classifications. For examples, athletes with high *BMI* levels due to high percentages of muscle mass, rather than fat, could be misclassified. To accommodate this sort of heterogeneity, the authors specified a heterogeneous threshold model,

$$\mu = \exp(\theta + \delta'r).$$

Finally, reasoning that the *BMI* outcome and the latent class assignment would likely depend on common features, both observed (in  $x$  and  $w$ ) and unobserved (in  $\varepsilon$  and  $u$ ), they specified a joint

normal distribution for  $\varepsilon$  and  $u$  with correlation  $\rho$ . (This is the first application of this model extension that we have seen. Since this is a cross section analysis, the natural extension is straightforward to build into the specification. If the sample were a panel, it would make sense to build a time invariant random effect into the main equation and allow that to be correlated with  $u$  in the class assignment. Some more elaborate specification would be necessary of the model specified more than two classes.

Combining all of the components, we have

*Outcome Model*

$$\begin{aligned} (BMI^*|class=c) &= \boldsymbol{\beta}_c' \mathbf{x} + \varepsilon_c, \varepsilon_c \sim N[0,1] \\ WT|class=c &= 0 \text{ if } BMI^*|class=c \leq 0 \\ &= 1 \text{ if } 0 < BMI^*|class=c \leq \mu_c \\ &= 2 \text{ if } BMI^*|class=c > \mu_c. \end{aligned}$$

$$\text{Threshold}|class=c: \mu_c = \exp(\theta_c + \boldsymbol{\delta}_c' \mathbf{r})$$

*Class Assignment*

$$\begin{aligned} c^* &= \boldsymbol{\alpha}' \mathbf{w} + u, u \sim N[0,1]. \\ c &= 0 \text{ if } c^* \leq 0 \\ &= 1 \text{ if } c^* > 0. \end{aligned}$$

*Endogenous Class Assignment*

$$(\varepsilon_c, u) \sim N_2[(0,0), (1, \rho_c, 1)]$$

Formation of the probabilities for the observed outcomes is a bit more complicated than previously due to the correlation between the class assignment and the *BMI* outcome. Generically,

$$\text{Prob}[WT=j | class=c] = \text{Prob}[WT=j, class=c] / \text{Prob}(class=c).$$

To form the likelihood, we require the joint probabilities, not the conditional;

$$\log L = \sum_{i=1}^N \log \left[ \sum_{class=0}^1 \text{Prob}(class=c) \text{Prob}(WT=j | class=c) \right].$$

The joint probability is a bivariate normal probability. (See Section 7.3.1 below.) To reach the components of the log likelihood and the probabilities to analyze for the partial effects, we begin with

$$\text{Prob}(WT=j | class=c) = \frac{\begin{cases} \Phi_2[(\mu_{i,j,c} - \boldsymbol{\beta}'_c \mathbf{x}), ((2c-1)\boldsymbol{\alpha}'_c \mathbf{z}), ((2c-1)\rho_c)] - \\ \Phi_2[(\mu_{i,j-1,c} - \boldsymbol{\beta}'_c \mathbf{x}), ((2c-1)\boldsymbol{\alpha}'_c \mathbf{z}), ((2c-1)\rho_c)] \end{cases}}{\Phi[(2c-1)\boldsymbol{\alpha}'_c \mathbf{z}]}$$

The unconditional probability is

$$\begin{aligned}
\text{Prob}(WT = j) &= \sum_{c=0}^1 \text{Prob}(\text{class} = c) \frac{\left\{ \begin{array}{l} \Phi_2[(\mu_{i,j,c} - \beta'_c \mathbf{x}), ((2c-1)\alpha' \mathbf{z}), ((2c-1)\rho_c)] - \\ \Phi_2[(\mu_{i,j-1,c} - \beta'_c \mathbf{x}), ((2c-1)\alpha' \mathbf{z}), ((2c-1)\rho_c)] \end{array} \right\}}{\Phi[(2c-1)\alpha' \mathbf{z}]} \\
&= \sum_{c=0}^1 \Phi[(2c-1)\alpha' \mathbf{z}] \frac{\left\{ \begin{array}{l} \Phi_2[(\mu_{i,j,c} - \beta'_c \mathbf{x}), ((2c-1)\alpha' \mathbf{z}), ((2c-1)\rho_c)] - \\ \Phi_2[(\mu_{i,j-1,c} - \beta'_c \mathbf{x}), ((2c-1)\alpha' \mathbf{z}), ((2c-1)\rho_c)] \end{array} \right\}}{\Phi[(2c-1)\alpha' \mathbf{z}]} \\
&= \sum_{c=0}^1 \left\{ \begin{array}{l} \Phi_2[(\mu_{i,j,c} - \beta'_c \mathbf{x}), ((2c-1)\alpha' \mathbf{z}), ((2c-1)\rho_c)] - \\ \Phi_2[(\mu_{i,j-1,c} - \beta'_c \mathbf{x}), ((2c-1)\alpha' \mathbf{z}), ((2c-1)\rho_c)] \end{array} \right\}
\end{aligned}$$

Combining all terms, then,

$$\log L = \sum_{i=1}^N \log \sum_{c=0}^1 \sum_{j=0}^2 m_{ij} \left\{ \begin{array}{l} \Phi_2[(\mu_{i,j,c} - \beta'_c \mathbf{x}), ((2c-1)\alpha' \mathbf{z}), ((2c-1)\rho_c)] - \\ \Phi_2[(\mu_{i,j-1,c} - \beta'_c \mathbf{x}), ((2c-1)\alpha' \mathbf{z}), ((2c-1)\rho_c)] \end{array} \right\}$$

where  $m_{ij} = j$  if  $WT_i = j$ ,  $j = 0, 1, 2$ ,  $\mu_{i,-1,c} = -\infty$ ,  $\mu_{i,0,c} = 0$ ,  $\mu_{i,1,c} = \exp(\theta_c + \delta_c' \mathbf{r}_i)$ ,  $\mu_{i,2,c} = +\infty$ .

In order to simplify the derivation of the partial effects, assume for the present that  $\mathbf{x}$ ,  $\mathbf{z}$  and  $\mathbf{r}$  all contain the same variables, labeled  $\mathbf{w}$ . The partial effects will contain three terms, for the latent regression, the class assignment and the threshold model. For variables that appear in more than one part, the partial effect will be obtained by adding the terms. For convenience, we will drop the observation subscript. Partial effects will typically be computed at the means of the variables, or by averaging the partial effects over all observations. Define the quantities

$$\begin{aligned}
A_{j,c} &= \mu_{j,c} - \beta'_c \mathbf{w}, \\
B_{j,c} &= (2c-1)\alpha' \mathbf{w}, \\
\tau_c &= (2c-1)\rho_c
\end{aligned}$$

Then,

$$\text{Prob}(WT = j | \mathbf{w}) = \sum_{c=0}^1 \Phi_2[A_{j,c}, B_{j,c}, \tau_c] - \Phi_2[A_{j-1,c}, B_{j-1,c}, \tau_c]$$

The partial effects are

$$\frac{\partial \text{Prob}(WT = j | \mathbf{x})}{\partial \mathbf{w}} = \sum_{c=0}^1 \left[ \begin{array}{l} \Phi \left( \frac{B_c - \tau_c A_{j,c}}{\sqrt{1 - \tau_c^2}} \right) [\phi(A_{j,c}) - \phi(A_{j-1,c})] (-\beta_c) + \\ \Phi \left( \frac{B_c - \tau_c A_{j,c}}{\sqrt{1 - \tau_c^2}} \right) [\phi(A_{j,c}) \mu_{j,c} - \phi(A_{j-1,c}) \mu_{j-1,c}] (\delta_c) + \\ \phi(B_c) \left[ \Phi \left( \frac{A_{j,c} - \tau_c B_c}{\sqrt{1 - \tau_c^2}} \right) - \Phi \left( \frac{A_{j-1,c} - \tau_c B_c}{\sqrt{1 - \tau_c^2}} \right) \right] (2c-1)(\alpha) \end{array} \right].$$

### 5.2.7 Generalized Ordered Choice Model (3)

In this section, we combine the features of the preceding generalized models in a single internally consistent model framework. The model contains random parameters, heterogeneous thresholds and heteroscedasticity. We depart from the base case,

$$\text{Prob}[y_i = j | \mathbf{x}_i] = F(\mu_j - \boldsymbol{\beta}'\mathbf{x}_i) - F(\mu_{j-1} - \boldsymbol{\beta}'\mathbf{x}_i) > 0, j = 0, 1, \dots, J.$$

The intrinsic heterogeneity across individuals is captured by writing

$$\boldsymbol{\beta}_i = \boldsymbol{\beta} + \Delta\mathbf{z}_i + \Gamma\mathbf{v}_i$$

where  $\Gamma$  is a lower triangular matrix and  $\mathbf{v}_i \sim N[\mathbf{0}, \mathbf{I}]$ . Thus,  $\boldsymbol{\beta}_i$  is normally distributed across individuals with conditional mean

$$E[\boldsymbol{\beta}_i | \mathbf{x}_i, \mathbf{z}_i] = \boldsymbol{\beta} + \Delta\mathbf{z}_i$$

and conditional variance

$$\text{Var}[\boldsymbol{\beta}_i | \mathbf{x}_i, \mathbf{z}_i] = \Gamma\Gamma' = \boldsymbol{\Omega}.$$

This is a random parameters formulation that appears elsewhere, e.g., Greene (2002, 2005) and Jones and Hensher (2004). It is the same as the random parameters model developed in the previous section, with the addition of the nonzero mean,  $\Delta\mathbf{z}_i$ , to the distribution of the heterogeneity in  $\boldsymbol{\beta}_i$ .

The thresholds are modeled as

$$\mu_{ij} = \mu_{i,j-1} + \exp(\alpha_j + \boldsymbol{\delta}'\mathbf{r}_i + \sigma_j w_{ij}), \mu_0 = 0, \mu_{-1} = -\infty, \mu_J = +\infty, w_{ij} \sim N[0, 1].$$

Integrating the difference equation, we obtain

$$\begin{aligned} \mu_1 &= \exp(\alpha_1 + \boldsymbol{\delta}'\mathbf{r}_i + \sigma_1 w_{j1}) \\ &= \exp(\boldsymbol{\delta}'\mathbf{r}_i) \exp(\alpha_1 + \sigma_1 w_{j1}) \\ \mu_2 &= \exp(\boldsymbol{\delta}'\mathbf{r}_i) [\exp(\alpha_1 + \sigma_1 w_{j1}) + \exp(\alpha_2 + \sigma_2 w_{j2})], \\ \mu_j &= \exp(\boldsymbol{\delta}'\mathbf{r}_i) \left( \sum_{m=1}^j \exp(\alpha_m + \sigma_m w_{jm}) \right) \\ \mu_J &= +\infty \text{ is imposed by } \alpha_J = +\infty \text{ and } \sigma_J = 0. \end{aligned}$$

This preserves the ordering of the thresholds and incorporates the necessary normalizations. Note that the thresholds, like the regression itself, are shifted by both observable ( $\mathbf{r}_i$ ) and unobservable ( $w_{ij}$ ) heterogeneity. [Theoretical models that produce ordered choice situations with stochastic thresholds are explored by Carneiro et al. (2001, 2003) and by Cunha et al. (2007). The models described by these authors have elements in common with the one described here, but do not appear to be implemented – the focus of these papers is on underlying structural models and on identification. The present model is transparently identified; our interest is in implementation.] The model is fully consistent in that probabilities are all positive and sum to one by construction. Finally, the disturbance variance is allowed to be heteroscedastic, as before, randomly as well as deterministically; thus,

$$\text{Var}[\varepsilon_i|\mathbf{h}_i] = \exp(\boldsymbol{\gamma}'\mathbf{h}_i + \tau e_i)$$

where  $e_i \sim N[0,1]$ .

Let  $\mathbf{v}_i = (v_{i1}, \dots, v_{iK})'$  and  $\mathbf{w}_i = (w_{i1}, \dots, w_{i,J-1})'$ . Combining terms, the conditional probability of outcome  $j$  is

$$\text{Prob}[y_i = j | \mathbf{x}_i, \mathbf{z}_i, \mathbf{h}_i, \mathbf{r}_i, \mathbf{v}_i, \mathbf{w}_i, e_i] = F \left[ \frac{\mu_{ij} - \boldsymbol{\beta}'_i \mathbf{x}_i}{\sqrt{\exp(\boldsymbol{\gamma}'\mathbf{h}_i + \tau e_i)}} \right] - F \left[ \frac{\mu_{i,j-1} - \boldsymbol{\beta}'_i \mathbf{x}_i}{\sqrt{\exp(\boldsymbol{\gamma}'\mathbf{h}_i + \tau e_i)}} \right].$$

The term that enters the log likelihood function is unconditioned on the unobservables. Thus, after integrating out the unobservable heterogeneity, we have

$$\begin{aligned} \text{Prob}[y_i = j | \mathbf{x}_i, \mathbf{z}_i, \mathbf{h}_i, \mathbf{r}_i] = \\ \int_{\mathbf{v}_i, \mathbf{w}_i, e_i} \left( F \left[ \frac{\mu_{ij} - \boldsymbol{\beta}'_i \mathbf{x}_i}{\sqrt{\exp(\boldsymbol{\gamma}'\mathbf{h}_i + \tau e_i)}} \right] - F \left[ \frac{\mu_{i,j-1} - \boldsymbol{\beta}'_i \mathbf{x}_i}{\sqrt{\exp(\boldsymbol{\gamma}'\mathbf{h}_i + \tau e_i)}} \right] \right) f(\mathbf{v}_i, \mathbf{w}_i, e_i) d\mathbf{v}_i d\mathbf{w}_i de_i. \end{aligned}$$

where

$$\boldsymbol{\beta}_i = \boldsymbol{\beta} + \Delta \mathbf{z}_i + \mathbf{L} \mathbf{D} \mathbf{v}_i$$

and

$$\mu_{ij} = \exp(\boldsymbol{\delta}' \mathbf{r}_i) \left( \sum_{m=1}^j \exp(\alpha_m + \sigma_m w_{im}) \right), j = 1, \dots, J-1.$$

The model is estimated by maximum simulated likelihood. The simulated log likelihood function is

$$\begin{aligned} \log L_S(\boldsymbol{\beta}, \Delta, \boldsymbol{\alpha}, \boldsymbol{\delta}, \boldsymbol{\gamma}, \mathbf{L}, \mathbf{D}, \boldsymbol{\sigma}, \tau) = \\ \sum_{i=1}^n \log \frac{1}{R} \sum_{r=1}^R \left( F \left[ \frac{\mu_{ij,r} - \boldsymbol{\beta}'_{i,r} \mathbf{x}_i}{\sqrt{\exp(\boldsymbol{\gamma}'\mathbf{h}_i + \tau e_{i,r})}} \right] - F \left[ \frac{\mu_{i,j-1,r} - \boldsymbol{\beta}'_{i,r} \mathbf{x}_i}{\sqrt{\exp(\boldsymbol{\gamma}'\mathbf{h}_i + \tau e_{i,r})}} \right] \right). \end{aligned}$$

(This is the model in its full generality. Whether a particular data set is rich enough to support this much parameterization, particularly the elements of the covariances of the unobservables in  $\boldsymbol{\Gamma}$ , is an empirical question that will depend on the application.)

The model contains three points at which changes in the observed variables can induce changes in the probabilities of the outcomes, in the thresholds, in the utility function, and in the variance. For convenience in the derivations, let a vector  $\mathbf{a}_i$  denote the union of  $(\mathbf{x}_i, \mathbf{r}_i, \mathbf{z}_i, \mathbf{h}_i)$ . This allows for cases in which variables appear at more than one place in the model. The partial effect of an element of  $\mathbf{a}_i$  on the probability will depend on where it appears in the specification. For cases in which a variable appears in more than one location, the partial effect will be the sum of the two, three or four terms. To avoid a cumbersome reparameterization of the model to place zeros in the appropriate places in the various parameter vectors and matrix, we simply assume at this point that  $\mathbf{a}_i$  appears in full throughout the model. Thus, we write the probability of interest as

$$\begin{aligned} \text{Prob}(y_i = j | \mathbf{a}_i) = \\ \int_{\mathbf{v}_i, \mathbf{w}_i, e_i} \left( F \left[ \frac{\mu_{ij} - (\boldsymbol{\beta} + \Delta \mathbf{a}_i + \mathbf{L} \mathbf{D} \mathbf{v}_i)' \mathbf{a}_i}{\sqrt{\exp(\boldsymbol{\gamma}' \mathbf{a}_i + \tau e_i)}} \right] - F \left[ \frac{\mu_{i,j-1} - (\boldsymbol{\beta} + \Delta \mathbf{a}_i + \mathbf{L} \mathbf{D} \mathbf{v}_i)' \mathbf{a}_i}{\sqrt{\exp(\boldsymbol{\gamma}' \mathbf{a}_i + \tau e_i)}} \right] \right) f(\mathbf{v}_i, \mathbf{w}_i, e_i) d\mathbf{v}_i d\mathbf{w}_i de_i. \end{aligned}$$

$$\mu_{ij} = \exp(\boldsymbol{\delta}'\mathbf{a}_i) \left( \sum_{m=1}^j \exp(\alpha_m + \sigma_m w_{im}) \right), j = 1, \dots, J-1.$$

The set of partial effects is

$$\begin{aligned} \frac{\partial \text{Prob}(y_i = j | \mathbf{a}_i)}{\partial \mathbf{a}_i} = & \\ & \int_{\mathbf{v}_i, \mathbf{w}_i, e_i} \left( f \left[ \frac{\mu_{ij} - \boldsymbol{\beta}'_i \mathbf{a}_i}{\sqrt{\exp(\boldsymbol{\gamma}'_i \mathbf{a}_i + \tau e_i)}} \right] \frac{1}{\sqrt{\exp(\boldsymbol{\gamma}'_i \mathbf{a}_i + \tau e_i)}} (\boldsymbol{\beta}_i + 2\Delta \mathbf{a}_i - \frac{1}{2} \boldsymbol{\gamma} + \mu_{ij} \boldsymbol{\delta}) \right) f(\mathbf{v}_i, \mathbf{w}_i, e_i) d\mathbf{v}_i d\mathbf{w}_i de_i \\ & - \int_{\mathbf{v}_i, \mathbf{w}_i, e_i} \left( f \left[ \frac{\mu_{i,j-1} - \boldsymbol{\beta}'_i \mathbf{a}_i}{\sqrt{\exp(\boldsymbol{\gamma}'_i \mathbf{a}_i + \tau e_i)}} \right] \frac{1}{\sqrt{\exp(\boldsymbol{\gamma}'_i \mathbf{a}_i + \tau e_i)}} (\boldsymbol{\beta}_i + 2\Delta \mathbf{a}_i - \frac{1}{2} \boldsymbol{\gamma} + \mu_{i,j-1} \boldsymbol{\delta}) \right) f(\mathbf{v}_i, \mathbf{w}_i, e_i) d\mathbf{v}_i d\mathbf{w}_i de_i \end{aligned}$$

The four parts of the effect appear in the parentheses in the middle of each expression. Effects for particular variables are added from the corresponding parts.

Applications of this model appear in Eluru, Bhat and Hensher (2007) and Greene and Hensher (2008). The former is a study of extent of injuries in traffic accidents. In the latter, the authors examine the information processing strategies in commuter choices of travel routes. Table 24 below shows an application of the model (with some of its features) to our health satisfaction example. The estimated partial effects are presented in Table 25.

**Table 24 Estimated Generalized Random Thresholds Ordered Logit Model**

Variable	Coefficient	Standard Error	b/St.Er.	P[ Z >z]	Mean of X
Random Thresholds Ordered Choice Model					
Maximum Likelihood Estimates					
Dependent variable	HEALTH				
Weighting variable	None				
Number of observations	4483				
Iterations completed	101				
Log likelihood function	-5725.181				
Number of parameters	20				
Info. Criterion: AIC =	2.56310				
Info. Criterion: BIC =	2.59168				
Info. Criterion:HQIC =	2.57317				
Restricted log likelihood	-5748.656				
McFadden Pseudo R-squared	.0040836				
Chi squared	46.95039				
Degrees of freedom	10				
Prob[ChiSq > value] =	.0000000				
Underlying probabilities based on Logistic					
+-----+Latent Regression Equation					
Constant	11.7009***	1.49050482	7.850	.0000	
AGE	-.13297***	.02046869	-6.496	.0000	43.440107
EDUC	.32358***	.06666995	4.853	.0000	11.418086
INCOME	2.28773***	.77817018	2.940	.0033	.3487401
MARRIED	-.33971	.30954183	-1.097	.2724	.7521749
KIDS	.60536**	.30606283	1.978	.0479	.3794334
+-----+Intercept Terms in Random Thresholds					
Alpha-01	1.70601***	.14293233	11.936	.0000	
Alpha-02	2.27770***	.15706772	14.501	.0000	
Alpha-03	1.89260	4.81997970	.393	.6946	
+-----+Standard Deviations of Random Thresholds					
Alpha-01	.51951***	.17206133	3.019	.0025	
Alpha-02	.19948***	.06158478	3.239	.0012	
Alpha-03	4.23250	16.2463226	.261	.7945	
+-----+Standard Deviations of Random Regression Parameters					
Constant	2.50042**	1.04989065	2.382	.0172	
AGE	.04075***	.01346402	3.027	.0025	
EDUC	.00501	.06255616	.080	.9362	
INCOME	.63914	1.53472806	.416	.6771	
MARRIED	.55559*	.31455076	1.766	.0773	
KIDS	.12332	.57564957	.214	.8304	
+-----+Heteroscedasticity in Latent Regression Equation					
FEMALE	.00201	.05316741	.038	.9698	
+-----+Latent Heterogeneity in Variance of Epsilon					
Tau(v)	.30733**	.15029065	2.045	.0409	
Note: ***, **, * = Significance at 1%, 5%, 10% level.					

**Table 25 Estimated Partial Effects for Ordered Thresholds Model**

```

=====
|| Summary of Marginal Effects for Ordered Probability Model (probit) ||
|| Effects are computed by averaging over observs. during simulations. ||
=====
||
||           Regression Variable AGE
||           =====
Outcome   Effect   dPy<=nn/dX   dPy>=nn/dX
=====
Y = 00    .00574    .00574    .03235
Y = 01    .01125    .01699    .02661
Y = 02    .01125    .02824    .01537
Y = 03    .01125    .03948    .00412
Y = 04    -.00713    .03235    -.00713
=====
||
||           Regression Variable EDUC           Regression Variable INCOME
||           =====
Outcome   Effect   dPy<=nn/dX   dPy>=nn/dX   Effect   dPy<=nn/dX   dPy>=nn/dX
=====
Y = 00    -.01397   -.01397   -.07873   -.09879   -.09879   -.55665
Y = 01    -.02737   -.04134   -.06476   -.19350   -.29229   -.45786
Y = 02    -.02737   -.06871   -.03739   -.19350   -.48579   -.26436
Y = 03    -.02737   -.09608   -.01002   -.19350   -.67929   -.07086
Y = 04    .01735    -.07873    .01735    .12263    -.55665    .12263
=====
||
||           Regression Variable MARRIED       Regression Variable KIDS
||           =====
Outcome   Effect   dPy<=nn/dX   dPy>=nn/dX   Effect   dPy<=nn/dX   dPy>=nn/dX
=====
Y = 00    .01467    .01467    .08266   -.02614   -.02614   -.14730
Y = 01    .02873    .04340    .06799   -.05120   -.07734   -.12116
Y = 02    .02873    .07213    .03926   -.05120   -.12855   -.06995
Y = 03    .02873    .10087    .01052   -.05120   -.17975   -.01875
Y = 04    -.01821    .08266    -.01821    .03245   -.14730    .03245
=====

```

### 5.3 Specification Tests for Ordered Choice Models

The ordered probit model is a conventional model by the standards of maximum likelihood estimation. Under the assumptions that the model is correctly specified and the data on  $y_i$  and  $\mathbf{x}_i$  are “well behaved,” [see Greene (2008a, chapter 4)], the familiar asymptotics and testing procedures used in Section 4.5 apply. That is, we can use the familiar apparatus, Wald, Lagrange multiplier and likelihood ratio procedures to test against null hypotheses that are nested within the essential parametric model,

$$\text{Prob}[y_i = j \mid \mathbf{x}_i] = F(\mu_j - \boldsymbol{\beta}'\mathbf{x}_i) - F(\mu_{j-1} - \boldsymbol{\beta}'\mathbf{x}_i) > 0, j = 0, 1, \dots, J.$$

Since the asymptotic theory relies on central limit theorems, and not on the specific distribution of  $\varepsilon_i$ , the same devices will apply for logit and probit models. Procedures for “exact” inference based specifically on the distribution assumed [see, e.g., Mehta and Patel (1995)] have not been developed for ordered choice models.

In this section, we consider “specification tests.” That is, tests against the null specification of the model, for which often there is no clearly defined alternative. For example, a test of the appropriateness of the assumption that  $\varepsilon_i$  is normally distributed is considered against the alternative that it is not. Specification tests for the ordered choice model have been obtained essentially for two issues, functional form and distribution. The functional form question relates to the assumption about the basic model specification,

$$\text{Prob}(y_i > j \mid \mathbf{x}_i) = F(\boldsymbol{\beta}'\mathbf{x}_i - \mu_j), j = 0, \dots, J-1.$$

The linearity of the index function is the main issue, though it will be clear shortly that, because the alternative hypothesis is not clearly stated, a test against this null might pick up a variety of other failures of the model assumption. The distributional tests are specifically directed to the question of whether normality (or logisticality) is appropriate. Once again, the alternative hypothesis is unclear. For example, it seems reasonable to suggest that a test against normality might be picking up the influence of an omitted variable – perhaps one with a skewed distribution. Recognizing the essential ambiguity of the nature of these tests, we can nonetheless usefully divide them into these two broad groupings.

We note in passing, a third type of specification test that has been considered. Section 6.3 discusses a counterpart to the Hausman (1978) test for random vs. fixed effects in a panel data model. Since the test is considered in detail there, we will not reconsider it in this section.

#### 5.3.1 Model Specifications – Missing Variables and Heteroscedasticity

A number of studies have considered the null specification of the ordered choice models against specific alternatives. These tests involve three particular features of the model, missing variables, heteroscedasticity and the distribution of  $\varepsilon_i$ . Murphy (1994, 1996), for example, examines the ordered logit model from earlier as a special case of the more general model

$$y_i^* = \boldsymbol{\beta}'\mathbf{x}_i + \boldsymbol{\gamma}'\mathbf{z}_i + \sigma_i \varepsilon_i$$

$$y_i = j \text{ if } \mu_{j-1} < y_i^* \leq \mu_j$$

where

- (1)  $\mathbf{z}_i$  is a set of omitted variables that are believed to be appropriate to be in the model;
- (2)  $\sigma_i^2 = (\pi^2/3)[\exp(\boldsymbol{\alpha}'\mathbf{h}_i)]^2$ ;
- (3)  $F(\varepsilon_i) = [1 + \exp(-\varepsilon_i)]^{-\delta}$ . (This is an asymmetric distribution.)

Murphy's extended ordered logit model encompasses the familiar ordered logit model; the null hypothesis of the restricted model would be  $\boldsymbol{\gamma} = \mathbf{0}$ ,  $\boldsymbol{\alpha} = \mathbf{0}$ ,  $\delta = 1$ . In principle, the alternative model can be fit by full maximum likelihood. If so, then the tests of the three specifications can be done one at a time or jointly, using Wald or Likelihood ratio tests. Murphy proposes Lagrange multiplier tests for the three hypotheses that involve only estimating the restricted, basic model. We will consider the missing variables and heteroscedasticity tests here, and return to the distribution in the next section.

For the moment, we revert to the simpler distribution with  $\delta = 1$ , and examine the *LM* test for missing variables and heteroscedasticity. Without the special consideration of the shape of the distribution ( $\delta$ ), the testing procedures are the same for the probit and logit models, so they are given generically below. In this context, it is worth noting, since  $\mathbf{z}_i$  is observed, not much is gained by using an *LM* test for missing variables; one can just as easily fit the full model and use the *LM* or Wald test of the null hypothesis that  $\boldsymbol{\gamma} = \mathbf{0}$ . The test for heteroscedasticity is likewise straightforward if one is able to fit the full model with this form of heteroscedasticity. [The *LM* tests proposed by Murphy (1994, 1996) and Weiss (1997) actually apply to any form of heteroscedasticity such that  $\sigma_i^2 = \sigma_0^2 w(\boldsymbol{\gamma}, \mathbf{h}_i)$  such that  $w(\mathbf{0}, \mathbf{h}_i) = 1$ . [See Breusch and Pagan (1979).] Harvey's (1976) model has been the form usually used in the received applications.

Consider, first, an *LM* test for missing variables. The log likelihood function is

$$\log L = \sum_{i=1}^N \sum_{j=0}^J m_{ij} \log [F(\mu_j - \boldsymbol{\beta}'\mathbf{x}_i - \boldsymbol{\gamma}'\mathbf{z}_i) - F(\mu_{j-1} - \boldsymbol{\beta}'\mathbf{x}_i - \boldsymbol{\gamma}'\mathbf{z}_i)].$$

The *LM* test is carried out by estimating the model under the null hypothesis that  $\boldsymbol{\gamma} = \mathbf{0}$ , then obtaining the statistic,

$$LM = \left( \text{Est.} \frac{\partial \log L}{\partial \begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\mu} \\ \boldsymbol{\gamma} \end{pmatrix}} \right)'_{|\boldsymbol{\gamma}=\mathbf{0}} \left\{ \text{Est.} \text{Var} \left[ \frac{\partial \log L}{\partial \begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\mu} \\ \boldsymbol{\gamma} \end{pmatrix}} \right] \right\}^{-1} \left( \text{Est.} \frac{\partial \log L}{\partial \begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\mu} \\ \boldsymbol{\gamma} \end{pmatrix}} \right)_{|\boldsymbol{\gamma}=\mathbf{0}}$$

(We have reversed the usual order of  $\boldsymbol{\gamma}$  and  $\boldsymbol{\mu}$  for convenience.) The test statistic is used to test the hypothesis that the gradient is zero at the restricted parameter vector. When the restricted model is fit by maximum likelihood, the derivatives with respect to  $\boldsymbol{\beta}$  and  $\boldsymbol{\mu}$  evaluated at the MLEs are numerically zero, so the sample estimator of the statistic is

$$LM = \left( \begin{array}{c} \mathbf{0} \\ \mathbf{0} \\ \frac{\partial \log L}{\partial \boldsymbol{\gamma}} \end{array} \right)'_{|\boldsymbol{\gamma}=\mathbf{0}} \left\{ \text{Est.} \text{Var} \left[ \frac{\partial \log L}{\partial \begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\mu} \\ \boldsymbol{\gamma} \end{pmatrix}} \right]_{|\boldsymbol{\gamma}=\mathbf{0}} \right\}^{-1} \left( \begin{array}{c} \mathbf{0} \\ \mathbf{0} \\ \frac{\partial \log L}{\partial \boldsymbol{\gamma}} \end{array} \right)_{|\boldsymbol{\gamma}=\mathbf{0}}$$

The practical application of the test requires computation of the derivatives of the log likelihood with respect to  $\boldsymbol{\gamma}$ , evaluated at  $\boldsymbol{\gamma} = 0$ , and an estimator of the asymptotic covariance matrix, which we consider below. For the derivatives,

$$\begin{aligned} \left( \frac{\partial \log L}{\partial \boldsymbol{\gamma}} \right)_{\boldsymbol{\gamma}=0} &= \left\{ \sum_{i=1}^N \sum_{j=0}^J m_{ij} \left[ \frac{f(\boldsymbol{\mu}_j - \boldsymbol{\beta}'\mathbf{x}_i - \boldsymbol{\gamma}'\mathbf{z}_i) - f(\boldsymbol{\mu}_{j-1} - \boldsymbol{\beta}'\mathbf{x}_i - \boldsymbol{\gamma}'\mathbf{z}_i)}{F(\boldsymbol{\mu}_j - \boldsymbol{\beta}'\mathbf{x}_i - \boldsymbol{\gamma}'\mathbf{z}_i) - F(\boldsymbol{\mu}_{j-1} - \boldsymbol{\beta}'\mathbf{x}_i - \boldsymbol{\gamma}'\mathbf{z}_i)} \right] (-\mathbf{z}_i) \right\}_{\boldsymbol{\gamma}=0} \\ &= \sum_{i=1}^N \sum_{j=0}^J m_{ij} \left[ \frac{f(\boldsymbol{\mu}_j - \boldsymbol{\beta}'\mathbf{x}_i) - f(\boldsymbol{\mu}_{j-1} - \boldsymbol{\beta}'\mathbf{x}_i)}{F(\boldsymbol{\mu}_j - \boldsymbol{\beta}'\mathbf{x}_i) - F(\boldsymbol{\mu}_{j-1} - \boldsymbol{\beta}'\mathbf{x}_i)} \right] (-\mathbf{z}_i) \end{aligned}$$

It remains to obtain the appropriate asymptotic covariance matrix. A convenient estimator is the sum of the outer products of the individual gradients,  $\mathbf{H}^{-1} = [\boldsymbol{\Sigma}_i \mathbf{g}_i \mathbf{g}_i']^{-1}$ . Davidson and MacKinnon (1983, 1984), MacKinnon (1992) and Godfrey (1988) [see, also, Weiss (1997)] present persuasive evidence, however, that the finite sample properties of the *LM* statistic are notably inferior to those when it is based on the second derivatives matrix or, when possible, the expected second derivatives matrix. Precise expressions for the second derivatives matrix appear in various places, including McElvey and Zavoina (1975) and Maddala (1983). Write the second derivatives matrix in the partitioned form

$$\mathbf{H} = \begin{bmatrix} \mathbf{H}_{\beta\beta'} & \mathbf{H}_{\beta\mu'} & \mathbf{H}_{\beta\gamma'} \\ \mathbf{H}_{\mu\beta'} & \mathbf{H}_{\mu\mu'} & \mathbf{H}_{\mu\gamma'} \\ \mathbf{H}_{\gamma\beta'} & \mathbf{H}_{\gamma\mu'} & \mathbf{H}_{\gamma\gamma'} \end{bmatrix}.$$

Then, from the form of the first derivatives, it can be seen that the *LM* statistic equals the first derivatives vector times the lower right submatrix of  $\mathbf{H}^{-1}$ . Collecting terms and using the partitioned inverse form [Greene (2008a, result A-74)], this will be

$$LM = \left\{ \left( \frac{\partial \log L}{\partial \boldsymbol{\gamma}} \right)_{\boldsymbol{\gamma}=0} \right\}' \left[ \mathbf{H}_{\gamma\gamma'} - (\mathbf{H}_{\gamma\beta'} \quad \mathbf{H}_{\gamma\mu'}) \begin{bmatrix} \mathbf{H}_{\beta\beta'} & \mathbf{H}_{\beta\mu'} \\ \mathbf{H}_{\mu\beta'} & \mathbf{H}_{\mu\mu'} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{H}_{\beta\gamma'} \\ \mathbf{H}_{\mu\gamma'} \end{bmatrix} \right]^{-1} \left\{ \left( \frac{\partial \log L}{\partial \boldsymbol{\gamma}} \right)_{\boldsymbol{\gamma}=0} \right\}$$

Weiss (1997) notes an interesting interpretation of the *LM* test for omitted variables. The gradient,  $(\partial \log L / \partial \boldsymbol{\gamma})_{\boldsymbol{\gamma}=0}$  given earlier can be written

$$\begin{aligned} \left( \frac{\partial \log L}{\partial \boldsymbol{\gamma}} \right)_{\boldsymbol{\gamma}=0} &= \sum_{i=1}^N \sum_{j=1}^J m_{ij} \left[ \frac{f(\boldsymbol{\mu}_{j-1} - \boldsymbol{\beta}'\mathbf{x}_i) - f(\boldsymbol{\mu}_j - \boldsymbol{\beta}'\mathbf{x}_i)}{F(\boldsymbol{\mu}_j - \boldsymbol{\beta}'\mathbf{x}_i) - F(\boldsymbol{\mu}_{j-1} - \boldsymbol{\beta}'\mathbf{x}_i)} \right] (\mathbf{z}_i) \\ &= \sum_{i=1}^N \sum_{j=1}^J m_{ij} E[\varepsilon_i | \mathbf{x}_i, y_i = j] (\mathbf{z}_i) \\ &= \sum_{i=1}^N \sum_{j=1}^J m_{ij} GR(1)_{ij} \mathbf{z}_i. \end{aligned}$$

That is, the test is based on the covariance between the (unobserved) disturbance and the omitted variables. This is precisely the approach used in the linear regression model, where  $\varepsilon_i$  is estimated directly with the residual,  $e_i$ . In this case, the estimator is a “*generalized residual*.” [See Chesher and Irish (1987) and Gourieroux et al. (1987).]

An *LM* test for heteroscedasticity is essentially the same, save for the considerably more complicated first and second derivatives. The model with heteroscedasticity (and no missing variables) has

$$\log L = \sum_{i=1}^N \sum_{j=1}^J m_{ij} \log \left[ F \left( \frac{\mu_j - \beta' \mathbf{x}_i}{\sigma_i} \right) - F \left( \frac{\mu_{j-1} - \beta' \mathbf{x}_i}{\sigma_i} \right) \right]$$

The first derivative vector is

$$\frac{\partial \log L}{\partial \delta} = \sum_{i=1}^N \sum_{j=1}^J m_{ij} \frac{\left[ f \left( \frac{\mu_j - \beta' \mathbf{x}_i}{\sigma_i} \right) \left( \frac{\mu_j - \beta' \mathbf{x}_i}{\sigma_i} \right) - f \left( \frac{\mu_{j-1} - \beta' \mathbf{x}_i}{\sigma_i} \right) \left( \frac{\mu_{j-1} - \beta' \mathbf{x}_i}{\sigma_i} \right) \right]}{F \left( \frac{\mu_j - \beta' \mathbf{x}_i}{\sigma_i} \right) - F \left( \frac{\mu_{j-1} - \beta' \mathbf{x}_i}{\sigma_i} \right)} (-\mathbf{h}_i)$$

The remaining computations are analogous to those done for the missing variables test. Note that under the null hypothesis,  $\sigma_i = 1$ , which considerably simplifies computing (albeit not deriving) the first and second derivatives.

In many cases, test statistics such as the *LM* statistic are computable using “artificial regressions.” For many of the common applications, we may write the *LM* statistic in the form

$$LM = \left( \sum_{i=1}^N w_i(\boldsymbol{\theta}) \mathbf{g}_i(\boldsymbol{\theta}) \right)' \left( \sum_{i=1}^N \mathbf{g}_i(\boldsymbol{\theta}) \mathbf{g}_i'(\boldsymbol{\theta}) \right)^{-1} \left( \sum_{i=1}^N w_i(\boldsymbol{\theta}) \mathbf{g}_i(\boldsymbol{\theta}) \right)$$

where  $\boldsymbol{\theta}$  is the full parameter vector being analyzed. In this case, the *LM* statistic is equal to the explained sum of squares in the regression of the variable  $w_i(\boldsymbol{\theta})$  on the “regressors,”  $\mathbf{g}_i(\boldsymbol{\theta})$ . [See MacKinnon (1992) and Orme (1990).] Consider the narrower case in which  $\mathbf{g}_i(\boldsymbol{\theta})$  is the full gradient,  $w_i(\boldsymbol{\theta}) = 1$ , and the outer product of gradients (OPG) estimator of the covariance matrix is used to complete the statistic. Then, the “dependent variable” in this regression is 1 for all  $i$ , the total uncentered sum of squares is  $N$ , and  $LM = NR^2$  in the artificial regression. [See Davidson and MacKinnon (1984).] With the extensive matrix manipulation routines in contemporary software such as *Stata* (Mata), *SAS* (Proc MATRIX), *NLOGIT* (Matrix), *Gauss* and *Matlab*, the appeal of the artificial regression interpretation is now largely confined to the analytics that precede computation.

### 5.3.2 Testing Against the Logistic and Normal Distributions

Murphy (1994, 1996) proposes an alternative distribution for  $\varepsilon$  in the ordered logit model,

$$F(\varepsilon) = \frac{1}{[1 + \exp(-\varepsilon)]^\delta}.$$

This distribution of  $\varepsilon_i$  is asymmetric; called a Burr type II distribution. This has been labeled the “scobit model” (skewed logit) elsewhere and has been suggested as an alternative to the normal and logistic distributions for binary choice models. [See Murphy (1994), Smith (1989), Lechner

(1991), Nagler (1994) and *Stata* (2008) or Econometric Software (2007).] The density is

$$f(\varepsilon) = \left( \frac{\exp(-\varepsilon)}{1 + \exp(-\varepsilon)} \right) \frac{\delta}{[1 + \exp(-\varepsilon)]^\delta}$$

For  $\delta = 1$ , the model reverts to the familiar logit form. Since this is fully parameterized, the alternative model can be fit directly and a Wald or likelihood ratio test can be used to test the null hypothesis that  $\delta = 1$ . Murphy proposes a Lagrange multiplier test that is based entirely on computations from the ordered logit model ( $\delta = 1$ ).

The scobit model has not been widely used in the ordered choice literature; tests about the distribution generally revolve around alternatives to the normal. Tests of the normality assumption build on the approach developed by Bera, Jarque and Lee (1984) for limited dependent variable models. A parametric alternative to the normal distribution is the Pearson family of distributions,

$$f(\varepsilon) = \frac{\exp[q(\varepsilon)]}{\int_{-\infty}^{\infty} \exp[q(t)] dt}, \text{ where } q(t) = \int \frac{c_1 - t}{c_0 - c_1 t + c_2 t^2} dt.$$

The relationship between the moments of the random variable and the three constants is

$$\begin{aligned} c_0 &= (4\tau_4 - 3\tau_3^2) / (10\tau_4 - 12\tau_3^2 - 18) \\ c_1 &= \tau_3 (\tau_4 + 3) / (10\tau_4 - 12\tau_3^2 - 18) \\ c_2 &= (2\tau_4 - 3\tau_3^2 - 6) / (10\tau_4 - 12\tau_3^2 - 18) \end{aligned}$$

[See Weiss (1997).] (We are avoiding a potentially confusing conflict in notation by using  $\tau$  rather than the conventional  $\mu$  to denote the moments of the distribution.) For the standard normal distribution,  $\tau_3 = 0$  and  $\tau_4 = 3$ . It follows that  $c_0 = 1$ ,  $c_1 = 0$  and  $c_2 = 0$ . (It also follows that the functional form is that of the standard normal.) Bera et al. (1984) developed an *LM* test for this restriction for the censored regression model. The corresponding result for the ordered probit model is given in Johnson (1996), Glewwe (1997) and Weiss (1997).

The test is based on the generalized residuals. For the normal distribution, we are testing against the hypothesis that the third and fourth moments of  $\varepsilon$  are  $\tau_3 = 0$  and  $\tau_4 = 3$ . As before, we cannot observe  $\varepsilon$ , so the test is based on the generalized residuals

$$\begin{aligned} E[\varepsilon_i^3 | y_i = j, \mathbf{x}_i] &= GR(3)_{ij} \\ &= \frac{[2 + (\mu_{j-1} - \boldsymbol{\beta}'\mathbf{x}_i)^2]\phi(\mu_{j-1} - \boldsymbol{\beta}'\mathbf{x}_i) - [2 + (\mu_j - \boldsymbol{\beta}'\mathbf{x}_i)^2]\phi(\mu_j - \boldsymbol{\beta}'\mathbf{x}_i)}{\Phi(\mu_j - \boldsymbol{\beta}'\mathbf{x}_i) - \Phi(\mu_{j-1} - \boldsymbol{\beta}'\mathbf{x}_i)} \\ E[\varepsilon_i^4 - 3 | y_i = j, \mathbf{x}_i] &= GR(4)_{ij} \\ &= \frac{[3 + (\mu_{j-1} - \boldsymbol{\beta}'\mathbf{x}_i)^2](\mu_{j-1} - \boldsymbol{\beta}'\mathbf{x}_i)\phi(\mu_{j-1} - \boldsymbol{\beta}'\mathbf{x}_i) - [3 + (\mu_j - \boldsymbol{\beta}'\mathbf{x}_i)^2](\mu_j - \boldsymbol{\beta}'\mathbf{x}_i)\phi(\mu_j - \boldsymbol{\beta}'\mathbf{x}_i)}{\Phi(\mu_j - \boldsymbol{\beta}'\mathbf{x}_i) - \Phi(\mu_{j-1} - \boldsymbol{\beta}'\mathbf{x}_i)} \end{aligned}$$

The full derivative vector including  $c_1$  and  $c_2$  evaluated at  $c_1 = c_2 = 0$  is

$$\frac{\partial \log L}{\partial \begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\mu} \\ c_1 \\ c_2 \end{pmatrix}} = \sum_{i=1}^N \left[ \sum_{j=0}^J m_{ij} \begin{pmatrix} GR(1)_{ij} \mathbf{x}_i \\ GR(1)_{ij} \mathbf{a}_j \\ GR(3)_{ij} \\ GR(4)_{ij} \end{pmatrix} \right] = \sum_{i=1}^N \mathbf{g}_i.$$

where  $\mathbf{a}_j$  is a  $(J-1) \times 1$  vector that has a 1 in position  $j$  and a -1 in position  $j-1$  save for  $j=1$ , when the  $j-1$  position is absent. To complete the computation of the test statistic, an estimator of the covariance matrix of the gradient is needed. Notwithstanding its less than ideal finite sample properties, the usual choice is the outer products matrix,

$$\mathbf{V} = \sum_{i=1}^N \mathbf{g}_i \mathbf{g}_i'.$$

Using the maximum likelihood estimates from the ordered probit model, the first two parts of the derivative vector will be numerically zero. This, the final result for the LM statistic is

$$LM = \left[ \sum_{i=1}^N \sum_{j=0}^J m_{ij} \begin{pmatrix} GR(3)_{ij} \\ GR(4)_{ij} \end{pmatrix} \right]' \left[ \mathbf{V}_{c_0, c_1}^{-1} \right] \left[ \sum_{i=1}^N \sum_{j=0}^J m_{ij} \begin{pmatrix} GR(3)_{ij} \\ GR(4)_{ij} \end{pmatrix} \right]$$

where  $\mathbf{V}_{c_1, c_2}^{-1}$  denotes the southeast  $2 \times 2$  submatrix of  $\mathbf{V}^{-1}$ .

Glewwe (1997) discusses other methods of testing for normality (actually symmetry,  $\tau_3=0$ , and mesokurtosis,  $\tau_4=3$ ) without a full parameterization of the alternative hypothesis, by using conditional moment tests. Newey (1985), Tauchen (1985) and Pagan and Vella (1989) provide details. As Glewwe shows, the *LM* test is essentially the same test. The use of the generalized residuals above suggests why this should be expected. Even though the *LM* test is structured around the Pearson alternative, in the end, it is a test of the values of the third and moments. The use of conditional moment tests is pursued by Mora and Moro-Egido (2008). For  $J$  of the  $J+1$  outcomes (because one is redundant), the model implies a set of moment conditions,

$$E[m_{ij} - P_{ij}(\boldsymbol{\theta})] = 0,$$

based on the additional assumptions of the model that produce the precise form of the probabilities. The authors examine the effect of different choices of the estimator of the covariance matrix for the Wald tests, and different formulations of the density of  $\varepsilon$ .

### 5.3.3 Unspecified Alternatives

The Brant (1990) test developed in Section 5.1.2 is ostensibly a test against the null hypothesis

$$H_0: \boldsymbol{\beta}_0 = \boldsymbol{\beta}_1 = \dots = \boldsymbol{\beta}_{J-1}$$

in the model

$$\text{Prob}(y_i > j \mid \mathbf{x}_i) = F(\boldsymbol{\beta}_j' \mathbf{x}_i - \mu_j), j = 0, \dots, J-1.$$

The apparently natural alternative hypothesis is the generalized ordered choice model (1). However, that model is not an internally consistent model for the probabilities associated with the outcomes. The presumed alternative does not prevent negative probabilities. One might conclude that the alternative is the “generalized” model when all  $\mathbf{x}$ ’s are such that the probabilities *are* positive – that is, in a certain range of  $\mathbf{x}$ . However, that range also depends on  $\boldsymbol{\beta}$  and  $\boldsymbol{\mu}$ . So, the suggestion amounts to concluding that the model is internally consistent when it is internally consistent. There is no other way to delineate when the model is internally consistent, other than it is when it is. On the other hand, it is persuasive that that Brant test, when it rejects the “null” hypothesis, is picking up some failure of the assumptions of the model. We have examined a variety of generalizations of the ordered choice model; it seems reasonable to conclude that the Brant test might well be finding any of them as an alternative to the base case. Thus, the Brant test might reasonably be considered in the same light as other conditional moment tests. That is, under the null hypothesis, certain features should be observed (within sampling variability). The alternative is, essentially, “not the null.”

Butler and Chatterjee (1995, 1997) have reconsidered estimation of the ordered probit model using the generalized method of moments. The null model implies a set of orthogonality conditions based on the definition of the model,

$$E[m_{ij} - (F(\mu_j - \boldsymbol{\beta}'\mathbf{x}_i) - F(\mu_{j-1} - \boldsymbol{\beta}'\mathbf{x}_i))] = 0,$$

where  $m_{ij} = 1$  if  $y_i = j$  and 0 otherwise. This provides a set of orthogonality conditions,

$$E\{\mathbf{x}_i[m_{ij} - (F(\mu_j - \boldsymbol{\beta}'\mathbf{x}_i) - F(\mu_{j-1} - \boldsymbol{\beta}'\mathbf{x}_i))]\}, j = 0, \dots, J.$$

In principle, this implies  $(J+1)K$  moment conditions, but one, the last, is redundant. The implied sample moments are, then,

$$\bar{\mathbf{g}}(\boldsymbol{\beta}, \boldsymbol{\mu}) = \frac{1}{N} \sum_{i=1}^N \begin{pmatrix} \mathbf{x}_i[m_{i0} - F(-\boldsymbol{\beta}'\mathbf{x}_i)] \\ \mathbf{x}_i[m_{i1} - [F(\mu_1 - \boldsymbol{\beta}'\mathbf{x}_i) - F(-\boldsymbol{\beta}'\mathbf{x}_i)]] \\ \vdots \\ \mathbf{x}_i[m_{i,J-1} - [F(\mu_{J-1} - \boldsymbol{\beta}'\mathbf{x}_i) - F(\mu_{J-2} - \boldsymbol{\beta}'\mathbf{x}_i)]] \end{pmatrix} = \frac{1}{N} \sum_{i=1}^N \mathbf{g}_i(\boldsymbol{\beta}, \boldsymbol{\mu}).$$

The GMM estimator is then obtained by two steps: (1) Obtain a consistent estimator of  $\boldsymbol{\beta}$  and  $\boldsymbol{\mu}$ , say the MLE. then compute an estimator of  $\text{Asy. Var}[\bar{\mathbf{g}}(\boldsymbol{\beta}, \boldsymbol{\mu})]$ , such as

$$\mathbf{V} = \frac{1}{N} \left[ \frac{1}{N} \sum_{i=1}^N \mathbf{g}_i(\boldsymbol{\beta}, \boldsymbol{\mu}) \mathbf{g}_i(\boldsymbol{\beta}, \boldsymbol{\mu})' \right]$$

(2) minimize the GMM criterion

$$Nq = N \bar{\mathbf{g}}(\boldsymbol{\beta}, \boldsymbol{\mu})' \mathbf{V}^{-1} \bar{\mathbf{g}}(\boldsymbol{\beta}, \boldsymbol{\mu}).$$

The minimized value has a limiting chi squared distribution with degrees of freedom equal to the number of overidentifying restrictions. In this case, the number of moment conditions is  $J \times K$  and the number of parameters is  $K+J-1$ . The number of overidentifying restrictions is  $(J-1)(K-1)$ . The authors go on to explore the corresponding computations for a bivariate ordered probit model. This proliferates moment conditions, as there is a  $K$ -order condition for each pairing of  $y_{i1}$  and  $y_{i2}$  – though the paucity of observations in some cells might suggest dropping some of the moments.

In all cases, it is uncertain what the alternative hypothesis should be if  $Nq$  is significant. [It is not in the application studied in their paper – see Butler and Chatterjee (1995).] Two suggestions are exogeneity of the independent variabilities and, of course, the distributional assumption.

## 6. Ordered Choice Modeling with Panel Data

Development of models for panel data parallel those in other modeling settings. The departure point is the familiar fixed and random effects approaches. We then consider other types of applications including extensions of the random parameters and latent classes formulations, dynamic models and some special treatments that accommodate features peculiar to the ordered choice models.

### 6.1 Ordered Choice Models with Fixed Effects

An ordered choice model with fixed effects formulated in the most familiar fashion would be

$$\text{Prob}[y_{it} = j \mid \mathbf{x}_i] = F(\mu_j - \alpha_i - \boldsymbol{\beta}'\mathbf{x}_{it}) - F(\mu_{j-1} - \alpha_i - \boldsymbol{\beta}'\mathbf{x}_{it}) > 0, j = 0, 1, \dots, J.$$

At the outset, there are two problems that this model shares with other nonlinear fixed effects models. First, regardless of how estimation and analysis are approached, time invariant variables are precluded. Since social science applications typically include demographic variables such as gender and, for some at least, education level, that are time invariant, this is likely to be a significant obstacle. (Several of the variables in the GSOEP analyzed by Boes and Winkelmann (2006b) and others are time invariant.) Second, there is no sufficient statistic available to condition the fixed effects out of the model. That would imply that in order to estimate the model as stated, one must manipulate the full log likelihood,

$$\log L = \sum_{i=1}^N \log \left\{ \prod_{t=1}^{T_i} \left( \sum_{j=0}^J m_{ijt} \left[ \Phi(\mu_j - \alpha_i - \boldsymbol{\beta}'\mathbf{x}_{it}) - \Phi(\mu_{j-1} - \alpha_i - \boldsymbol{\beta}'\mathbf{x}_{it}) \right] \right) \right\}$$

If the sample is small enough, of course, one may simply insert the individual group dummy variables and treat the entire pooled sample as a cross section. See, e.g., Mora (2006) for a cross-country application in banking that includes separate country dummy variables. We are interested, instead, in the longitudinal data case in which this would not be feasible. The data set from which our sample used in the preceding examples is extracted from an unbalanced panel of 7,293 households, observed from 1 to 7 times each.

The full ordered probit model with fixed effects, including the individual specific constants, can be estimated by unconditional maximum likelihood using the results in Greene (2004a,b and 2008a, Section 16.9.6.c). The likelihood function is globally concave [see Pratt (1981)], so despite its superficial complexity, the estimation is straightforward. In another application, based on the full panel data set [see Greene (2008a, pp. 838-840), estimation of the full model required roughly five seconds of computation on an ordinary desktop computer.

The larger methodological problem with this approach would be at least the potential for the incidental parameters problem that has been widely documented for the binary choice case. [See, e.g., Lancaster (2000).] That is the small  $T$  bias in the estimated parameters when the full MLE is applied in panel data. For  $T = 2$  in the binary logit model, it has been shown analytically [Abrevaya (1997)] that the full MLE converges to  $2\boldsymbol{\beta}$ . [See, as well, Hsiao (1986, 2003).] No corresponding results have been obtained for larger  $T$  or for other models. However, Monte Carlo results have strongly suggested that the small sample bias persists for larger  $T$  as well, though as might be expected, it diminishes with increasing  $T$ .

No theoretical counterpart to the Hsiao (1986, 2003) and Abrevaya (1997) result on the small  $T$  bias (incidental parameters problem) of the MLE in the presence of fixed effects has been derived for the ordered probit model. The Monte Carlo results in Greene (2004b) reproduced

below in Figure 14 suggest that biases comparable to those in the binary choice models persist in the ordered probit model as well. (In the first, third and fifth rows that correspond to estimation of coefficients, the true coefficients being estimated both equal one.)

**Table 2.** Means of empirical sampling distributions,  $N = 1,000$  individuals based on 200 replications.

	$T = 2$		$T = 3$		$T = 5$		$T = 8$		$T = 10$		$T = 20$	
	$\beta$	$\delta$	$\beta$	$\delta$	$\beta$	$\delta$	$\beta$	$\delta$	$\beta$	$\delta$	$\beta$	$\delta$
Logit Coeff	2.020	2.027	1.698	1.668	1.379	1.323	1.217	1.156	1.161	1.135	1.069	1.062
Logit M.E. <sup>a</sup>	1.676	1.660	1.523	1.477	1.319	1.254	1.191	1.128	1.140	1.111	1.034	1.052
Probit Coeff	2.083	1.938	1.821	1.777	1.589	1.407	1.328	1.243	1.247	1.169	1.108	1.068
Probit M.E. <sup>a</sup>	1.474	1.388	1.392	1.354	1.406	1.231	1.241	1.152	1.190	1.110	1.088	1.047
Ord. Probit	2.328	2.605	1.592	1.806	1.305	1.415	1.166	1.220	1.131	1.158	1.058	1.068

<sup>a</sup>Average ratio of estimated marginal effect to true marginal effect.

**Figure 14 Monte Carlo Analysis of Biases in Fixed Effects MLE in Discrete Choice Models**

The preceding bode ill for unconditional fixed effects models for ordered choice. So far, the approach has little to recommend it other than the theoretical robustness of fixed effects as an alternative to random effects. Recent proposals for “bias reduction” estimators for binary choice models, including Fernandez-Val and Vella (2007), Fernandez-Val (2008), Carro (2007), Hahn and Newey (2004) and Hahn and Kuersteiner (2003) suggest some directions for further research. However, no counterparts for the ordered choice models have yet been developed. We would note, for this model, the estimation of  $\beta$  which is the focus of these estimators, is only a means to the end. As seen earlier, in order to make meaningful statements about the implications of the model for behavior, it will be necessary to compute probabilities and derivatives. These, in turn, will require estimation of the constants, or some surrogates. The problem remains to be solved.

In their application to the GSOEP panel data set, Boes and Winkelmann (2006b) further modify the heterogeneous thresholds model. Their model is a fixed effects model,

$$\text{Prob}[y_{it} = j \mid \mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{iT}] = F(\mu_{ij} - \beta_j' \mathbf{x}_{it}) - F(\mu_{i,j-1} - \beta_{j-1}' \mathbf{x}_i)$$

where

$$\mu_{ij} = \mu_j + \alpha_i.$$

Seeking to avoid the incidental parameters problem, they use Mundlak’s (1978) and Chamberlain’s (1980) device to model the fixed effect. Projecting the fixed effects on the group means of the regressors,

$$\alpha_i = \gamma_j' \bar{\mathbf{x}}_i + \sigma v_i$$

they obtain an equivalent random effects model,

$$\text{Prob}[y_{it} = j \mid \mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{iT}] = F(\mu_{ij} - \beta_j' \mathbf{x}_{it}) - F(\mu_{i,j-1} - \beta_{j-1}' \mathbf{x}_i)$$

where

$$\mu_{ij} = \mu_j + \gamma_j' \bar{\mathbf{x}}_i + \sigma v_i, v_i \sim N[0, 1]$$

and  $\sigma$  is a new parameter to be estimated. This model is estimated by using quadrature to integrate  $v_i$  out of the log likelihood. [See the next section and Butler and Moffitt (1982) for the methodology.] As observed at several earlier points, the placement of the heterogeneity in the thresholds is not substantive; it can be moved to the mean of the regression with no change in the

interpretation of the model. As usual, the placement of the fixed effects in this linear specification is not consequential. Thus, their model is functionally equivalent to a more conventional random effects model with the group means added as covariates;

$$\text{Prob}[y_{it} = j \mid \mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{iT}] = F[\mu_j - (\boldsymbol{\beta}'_j \mathbf{x}_{it} + \boldsymbol{\gamma}'_{*j} \bar{\mathbf{x}}_i + \sigma v_i)] - F[\mu_{j-1} - (\boldsymbol{\beta}'_{j-1} \mathbf{x}_{it} + \boldsymbol{\gamma}'_{*j-1} \bar{\mathbf{x}}_i + \sigma v_i)].$$

The underlying logic of the Brant test suggests an alternative approach to estimation proposed by Das and van Soest (2000). Consider the base case ordered logit model with fixed effects. The model assumptions imply that

$$\begin{aligned} \text{Prob}[y_{it} > j \mid \mathbf{x}_{it}] &= \Lambda(\alpha_i + \boldsymbol{\beta}' \mathbf{x}_{it} - \mu_j) \\ &= \Lambda[(\alpha_i - \mu_j) + \boldsymbol{\beta}' \mathbf{x}_{it}] \end{aligned}$$

Now, define a binary variable  $w_{it,j} = 1[y_{it} > j], j = 0, 1, \dots, J-1$ . It follows that

$$\begin{aligned} \text{Prob}[y_{it} > j \mid \mathbf{x}_{it}] &= \Lambda[(\alpha_i - \mu_j) + \boldsymbol{\beta}' \mathbf{x}_{it}] \\ &= \Lambda[\lambda_i + \boldsymbol{\beta}' \mathbf{x}_{it}] \\ &= \text{Prob}(w_{itj} = 1 \mid \mathbf{x}_{it}). \end{aligned}$$

The “ $j$ ” specific part of the constant is the same for all individuals so it is absorbed in  $\lambda_i$ . Thus, a fixed effects binary logit model applies to each of the  $J - 1$  binary random variables,  $w_{it,j}$ . The method of Rasch (1960), Andersen (1970) and Chamberlain (1980) can be applied to each of these binary choice models to obtain an estimator of  $\boldsymbol{\beta}$  without having to estimate the constant terms. [See also Greene (2008a, pp. 800-806).] This provides  $J - 1$  estimators of the parameter vector  $\boldsymbol{\beta}$  (but no estimator of the threshold parameters). The authors propose to reconcile these different estimators by using a minimum distance estimator of the common true  $\boldsymbol{\beta}$ . The minimum distance estimator at the second step is chosen to minimize

$$q = \sum_{l=1}^{J-1} \sum_{m=1}^{J-1} (\hat{\boldsymbol{\beta}}_l - \boldsymbol{\beta})' [\mathbf{V}^{-1}]_{lm} (\hat{\boldsymbol{\beta}}_m - \boldsymbol{\beta})$$

where  $[\mathbf{V}^{-1}]_{lm}$  is the  $l, m$  block of the inverse of the  $(J - 1)K \times (J - 1)K$  partitioned matrix  $\mathbf{V}$  that contains  $\text{Asy.Cov}(\hat{\boldsymbol{\beta}}_l, \hat{\boldsymbol{\beta}}_m)$ . The appropriate form of this matrix for a set of cross-section estimators is given in Brant (1990). Since Das and van Soest (2000) used the counterpart for Chamberlain’s fixed effects estimator, this would be inappropriate. They used, instead, a counterpart to the BHHH estimator. The  $l, m$  block of  $\mathbf{V}$  (before inversion) is computed using

$$\mathbf{V}_{lm} = \sum_{i=1}^N \left( \frac{\partial \log L_{i,l}}{\partial \boldsymbol{\beta}_l} \right) \left( \frac{\partial \log L_{i,m}}{\partial \boldsymbol{\beta}_m} \right)$$

where  $\log L_{i,m}$  is the contribution of individual  $i$  to the log likelihood for  $\boldsymbol{\beta}_l$ . The diagonal blocks of the matrix are the BHHH estimators for the asymptotic covariance matrices for the  $j$  specific estimators.

As in the binary choice case, the complication of the fixed effects model is the small  $T$  bias, not the computation. The Das and van Soest approach finesses this problem—their estimator is consistent—but at the cost of losing the information needed to compute partial effects or predicted probabilities.

Winkelmann and Winkelmann (1998) analyzed data on well being from the German Socioeconomic Panel (GSOEP). The central question under the analysis is “How satisfied are you at present with your life as a whole?” which was answered on a discrete scale from 0 to 10. (See Section 2.1 for discussion of the methodological aspects of this analysis.) The natural approach to the analysis would be an ordered choice – the authors were interested in the effect of unemployment on the response. A fixed effects ordered choice (logit) model is the starting point for the specification.. Since there is no sufficient statistic available to use to condition the fixed effects out of the log likelihood, and fitting the fixed effects model by brute force by including the dummy variables in the model (assuming it could be done) would induce the biases of the incidental parameters problem, the authors opted for a simpler strategy. They divided the responses (0 to 10) into “dissatisfied” and “satisfied” and recoded the former 0 and the latter 1, producing a binary choice model. The structure, then, is equivalent to

$$\begin{aligned}
 y_{it}^* &= \boldsymbol{\beta}'\mathbf{x}_{it} + \alpha_i + \varepsilon_{it}, \\
 y_{it} &= j \text{ if } \mu_{j-1} \leq y_{it}^* < \mu_j, j = 0, 1, \dots, 10, i = 1, \dots, N, t = 1, \dots, T_i, \\
 z_{it} &= 1 \text{ if } y_{it} > 7.
 \end{aligned}$$

(The average response on the observed  $y_{it}$  in the sample was between 7 and 8.) The transformation is equivalent to the 8<sup>th</sup> of the 10 possible binary choice models in the Das and van Soest (2000) formulation;

$$\text{Prob}(y_{it} > 7 \mid \mathbf{x}_{it}) = \Lambda(\alpha_i + \boldsymbol{\beta}'\mathbf{x}_{it} - \mu_7)$$

Once again, the constant  $\mu_7$  is absorbed in the individual specific constant term, to produce, as before,

$$\text{Prob}[z_{it} = 1 \mid \mathbf{x}_{it}] = \Lambda[\lambda_i + \boldsymbol{\beta}'\mathbf{x}_{it}].$$

The model was then fit using the same Rasch/Andersen/Chamberlain method noted earlier.

Ferrer-i-Carbonell and Frijters (2004) built on this approach in developing an alternative estimator. In their study, the response variable of interest, from the same GSOEP data set was “General Satisfaction.” One of the shortcomings of the fixed effect binary choice model (whether it is estimated conditionally as suggested above) or unconditionally by computing the full set of coefficients including  $\alpha_i$ ) is that groups that do not change outcomes in the  $T_i$  periods fall out of the sample. For the conditional model,

$$\text{Prob}(z_{i1}=1, z_{i2}=1, \dots, z_{iT}=1 \mid \sum_i z_{it} = T) = 1,$$

so the contribution of this observation group  $i$  to the log likelihood is zero if  $z_{it}$  is always equal to 1. (The same occurs if  $z_{it}$  equals zero in every period.) For the brute force approach, the likelihood equation for estimation if  $\alpha_i$  for a group in which  $z_{it}$  is the same in every period is

$$\partial \log L / \partial \alpha_i = \sum_t f(\alpha_i + \boldsymbol{\beta}'\mathbf{x}_{it}) = 0 \text{ if } z_{it} = 1 \text{ in every period,}$$

$$\partial \log L / \partial \alpha_i = \sum_t -f[-(\alpha_i + \boldsymbol{\beta}'\mathbf{x}_{it})] = 0 \text{ if } z_{it} = 0 \text{ in every period.}$$

The first order condition for estimation of  $\alpha_i$  cannot be met with a finite  $\alpha_i$  if  $z_{it}$  is always one or always zero in every period. For ordered choice data, this is likely to be a frequent occurrence, particularly at the two ends of the distribution. The implication is that the samples used for

possibly many of the of the binary choice equations in the Das and van Soest (2000) or the Winkelmann and Winkelmann (1998) estimator will lose many observations.

Ferrer-i-Carbonell and Frijters (2004) [and Frijters, Haisken-DeNew and Shields (2004)] modified the Winkelmann and Winkelmann (1998) approach. Initially, the approach is essentially the same, though it begins with a fixed effect *and* individual specific thresholds;

$$y_{it}^* = \alpha_i + \beta' \mathbf{x}_{it} + \varepsilon_{it}$$

$$y_{it} = j \text{ if } \mu_{j-1,i} \leq y_{it}^* < \mu_{j,i}, j = 0, \dots, J; i = 1, \dots, N; t = 1, \dots, T_i.$$

The ordered logit form is assumed. For each individual,  $i$ , in the sample, once again,

$$\text{Prob}[z_{it} = 1 \mid \mathbf{x}_{it}] = \Lambda[\lambda_i + \beta' \mathbf{x}_{it}].$$

The difference here is that  $z_{it}$  is defined with respect to an individual specific  $j_i^*$ , so

$$z_{it} = 1 \text{ if } y_{it} > j_i^* \text{ and } 0 \text{ otherwise.}$$

(In Winkelmann and Winkelmann's method,  $j_i^* = 7$  for all  $i$ .) The algorithm for choosing  $j_i^*$  efficiently for each individual is given in the paper. (The technical Appendix that describes their method can be downloaded from the website for the Royal Economic Society at [http://www.res.org.uk/economic/ta/pdfs/eco\\_j\\_235\\_app.pdf](http://www.res.org.uk/economic/ta/pdfs/eco_j_235_app.pdf). It is not contained in the paper, itself.) The resulting contribution to the likelihood for individual  $i$  is

$$\text{Prob}(y_{i1} > j_i^*, y_{i2} > j_i^*, \dots, y_{iT_i} > j_i^* \mid \mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{iT_i}, \sum_{t=1}^{T_i} z_{it} = c_i)$$

$$= \frac{\exp\left[\sum_{t=1}^{T_i} z_{it} \beta' \mathbf{x}_{it}\right]}{\sum_{(z_1, z_2, \dots, z_{T_i}) \in S(j_i^*, c_i)} \exp\left[\sum_{t=1}^{T_i} z_t \beta' \mathbf{x}_{it}\right]}$$

Where  $c_i = \sum_t z_{it}$  = the number of times  $y_{it}$  is greater than the chosen threshold. The threshold  $j_i^*$  is chosen so that  $c_i$  is not equal to 0 or  $T_i$ .  $S(j_i^*, c_i)$  is the set of all possible vectors,  $(z_1, z_2, \dots, z_{T_i})$ , whose elements are all zero or one and sum to  $c_i$ ; that is, the set of vectors corresponding to sets of outcomes  $y_{it}$  such that  $c_i$  of them are greater than  $j_i^*$ . The denominator of the probability is the sum over all possible arrangements of  $T_i$   $z$ 's such that the sum is  $c_i$ . [See Krailo and Pike (1984) for the computations involved.]

## 6.2 Ordered Choice Models with Random Effects

Save for an ambiguity about the mixture of distributions in an ordered logit model, a random effects version of the ordered choice model is a straightforward extension of the binary choice case developed by Butler and Moffitt (1982). An interesting application which appears to replicate, but not connect to Butler and Moffitt is Jansen (1990). Jansen estimates the equivalent of the Butler and Moffitt model with an ordered probit model, using an iterated MLE with quadrature used between iterations. Following Jansen's lead, Crouchley (1995) also designed the equivalent of the common random effects model, but embeds it in a complementary log-log form that allows, at least for his two period model, a closed form expression for the probabilities after the random effect is integrated out. Characteristically, this strand of the literature emerged completely apart from the social science counterpart, which had, by then, integrated the random effects, panel data model into a variety of single index specifications such as this one.

Crouchley's formulation of the "random-effects ordered response model" is

$$y_{ij} = \beta_0 + \boldsymbol{\beta}'\mathbf{x}_{ij} + \mathbf{b}_i'\mathbf{z}_{ij} + e_{ij}$$

where  $\mathbf{b}_i$  is a vector of individual specific random effects,  $\mathbf{x}_{ij}$  is a known design matrix, and  $e_{ij}$  is the stochastic disturbance. The model is immediately simplified to a single random effect,  $\mathbf{b}_i'\mathbf{z}_{ij} = e_i$ , which leaves

$$y_{ij} = \beta_0 + \boldsymbol{\beta}'\mathbf{x}_{ij} + e_i + e_{ij}, i = 1, \dots, N, j = 1, \dots, T_i.$$

The remainder of the treatment is an ordered complementary log-log model with random effects, which is very similar to the model we have considered so far. The difference from this point forward is in the functional form of the distributions of both  $e_i$  and  $e_{ij}$ , neither of which is assumed to be normal. Crouchley notes, the simplified dimensions can be relaxed

The structure of the random effects ordered choice model is

$$y_{it}^* = \boldsymbol{\beta}'\mathbf{x}_{it} + u_i + \varepsilon_{it}$$

$$y_{it} = j \text{ if } \mu_{j-1} \leq y_{it}^* < \mu_j$$

$$\varepsilon_{it} \sim f(\cdot) \text{ with mean zero and constant variance } 1 \text{ or } \pi^2/3 \text{ (probit or logit),}$$

$$u_i \sim g(\cdot) \text{ with mean zero and constant variance, } \sigma^2, \text{ independent of } \varepsilon_{it} \text{ for all } t.$$

If we maintain the ordered probit form and assume as well that  $u_i$  is normally distributed, then, at least superficially, we can see the implications for the estimator of ignoring the heterogeneity. Using the usual approach,

$$\begin{aligned} \text{Prob}(y_{it} = j | \mathbf{x}_{it}) &= \text{Prob}(\boldsymbol{\beta}'\mathbf{x}_{it} + u_i + \varepsilon_{it} < \mu_j) - \text{Prob}(\boldsymbol{\beta}'\mathbf{x}_{it} + u_i + \varepsilon_{it} < \mu_{j-1}) \\ &= \Phi\left(\frac{\mu_j}{\sqrt{1+\sigma^2}} - \frac{\boldsymbol{\beta}'\mathbf{x}_{it}}{\sqrt{1+\sigma^2}}\right) - \Phi\left(\frac{\mu_{j-1}}{\sqrt{1+\sigma^2}} - \frac{\boldsymbol{\beta}'\mathbf{x}_{it}}{\sqrt{1+\sigma^2}}\right) \\ &= \Phi(\tau_j - \boldsymbol{\gamma}'\mathbf{x}_{it}) - \Phi(\tau_{j-1} - \boldsymbol{\gamma}'\mathbf{x}_{it}). \end{aligned}$$

Unconditionally, then, the result is an ordered probit in the scaled threshold values and scaled coefficients. Evidently, this is what is estimated if the data are pooled and the heterogeneity is ignored. (Note that a "robust" covariance matrix estimator does not redeem the estimator.)

The likelihood function for a sample can be estimated using the method of Butler and Moffitt. It is convenient to write  $u_i = \sigma v_i$  where  $v_i$  is the standardized variable – for the moment,  $N(0,1)$ . Then, conditioned on  $v_i$ , the observations on  $y_{it}$ ,  $t = 1, \dots, T_i$  are independent, so the contribution to the conditional likelihood for individual  $i$  would be the joint probability,

$$\text{Prob}(y_{i1} = j_1, y_{i2} = j_2, \dots, y_{iT} = j_T | \mathbf{X}_i, v_i) = \prod_{t=1}^{T_i} \left[ \Phi(\mu_{j_t} - \boldsymbol{\beta}'\mathbf{x}_{it} - \sigma v_i) - \Phi(\mu_{j_t-1} - \boldsymbol{\beta}'\mathbf{x}_{it} - \sigma v_i) \right]$$

The unconditional probability would be, then,

$$P(\mathbf{y}_i = \mathbf{j}_i | \mathbf{X}_i) = \int_{-\infty}^{\infty} \prod_{t=1}^{T_i} \left[ \Phi(\mu_{j_t} - \boldsymbol{\beta}'\mathbf{x}_{it} - \sigma v_i) - \Phi(\mu_{j_t-1} - \boldsymbol{\beta}'\mathbf{x}_{it} - \sigma v_i) \right] \phi(v_i) dv_i$$

(where we have defined a shorthand for the joint probability). The unconditional log likelihood is

$$\log L = \sum_{i=1}^N \log \int_{-\infty}^{\infty} \prod_{t=1}^{T_i} \left[ \Phi(\mu_j - \beta' \mathbf{x}_{it} - \sigma v_i) - \Phi(\mu_{j-1} - \beta' \mathbf{x}_{it} - \sigma v_i) \right] \phi(v_i) dv_i$$

The remaining complication is how to compute the integral. Two methods are available. The method of Gauss-Hermite quadrature developed by Butler and Moffitt uses an approximation to the integrals;

$$\log L_H = \sum_{i=1}^N \log \sum_{m=1}^M WT_m \prod_{t=1}^{T_i} \left[ \Phi(\mu_j - \beta' \mathbf{x}_{it} - \sigma N_m) - \Phi(\mu_{j-1} - \beta' \mathbf{x}_{it} - \sigma N_m) \right]$$

where  $WT_m$  and  $N_m$  are the weights and nodes, respectively, for the quadrature. [See, e.g., Abramovitz and Stegun (1971).] The accuracy of the approximation is a function of  $M$ , the number of quadrature points. Greater accuracy is achieved with increased  $M$ , but at the cost of greater computation time. [See, e.g., Rabe-Hesketh, Skrondal and Pickles (2005).] An alternative approach to the estimation would be maximum simulated likelihood. The integral in the log likelihood is

$$\int_{v_i} (L_i | v_i) \phi(v_i) dv_i = E_{v_i} \left[ \prod_{t=1}^{T_i} \left[ \Phi(\mu_j - \beta' \mathbf{x}_{it} - \sigma v_i) - \Phi(\mu_{j-1} - \beta' \mathbf{x}_{it} - \sigma v_i) \right] \right],$$

which can be approximated using simulation. The simulated log likelihood to be maximized is

$$\log L_S = \sum_{i=1}^N \log \frac{1}{R} \sum_{r=1}^R \prod_{t=1}^{T_i} \left[ \Phi(\mu_j - \beta' \mathbf{x}_{it} - \sigma v_{ir}) - \Phi(\mu_{j-1} - \beta' \mathbf{x}_{it} - \sigma v_{ir}) \right]$$

where  $v_{ir}$ ,  $r = 1, \dots, R$  is a set of random draws from the standard normal population (the same set, reused every time the function is calculated for individual  $i$ ). [See Train (2003) and Greene (2008a, Chapter 17) for details on simulation based estimation.] Neither method of computation has an obvious advantage in this one dimensional integration problem. (In terms of computational time, the advantage shifts significantly in favor of simulation when the number of dimensions (the order of the integration) increases past two.)

The random effects model extends naturally to the ordered probit model if the heterogeneity is viewed as the sum of small influences – a central limit theorem could be invoked to justify the layering of the normally distributed heterogeneity,  $u_i$ , on the normally distributed disturbance,  $\varepsilon_{it}$ . That does raise an ambiguity in the specification of the ordered logit model. The appeal of the logistic distribution is largely its mathematical convenience, though the slightly thicker tails might lend it some additional utility. However, the mixture of a logistic disturbance with a normally distributed random effect is a bit unnatural. The Butler and Moffitt method does not extend readily to integrating the logistic distribution. However, the simulation method can easily be so adapted. The simulated ordered logit model is obtained by using the logistic cdf,  $\Lambda(\cdot)$  rather than the normal,  $\Phi(\cdot)$  in the function. Draws from the desired distribution are simply obtained by the appropriate transformation of draws,  $U_{ir}$ , from the standard uniform,  $U(0,1)$ ;  $\Phi^{-1}(U_{ir})$  for simulation from the normal, or  $\log[U_{ir}/(1-U_{ir})]$  for the logistic. The optimization process is the same for the two cases. The deeper question would seem to be whether the logistic/logistic model is a reasonable one in the abstract, compared to the more commonly used normal/normal.

### 6.3 Testing for Random or Fixed Effects: A Variable Addition Test

A natural question is whether there is a test one can use to determine whether fixed or random effects should be the preferred model. Since the models are not nested, no simple test based on the likelihood function is available. A counterpart to the Hausman (1978) test for the linear model seems desirable, however, unlike the linear case, the fixed effects estimator for this nonlinear model is inconsistent even when it is the appropriate estimator (due to the incidental parameters problem). If one is going to base any test on the estimator of the fixed effects model, it would appear to be necessary to use one of the modified approaches, by Das and van Soest (2000) or Frijters et al. (2004), or any of the individual implied binary choice models, any of which will produce a consistent estimator of  $\beta$  under the hypothesis that the fixed effects model is appropriate. As such, this will force the fixed effects benchmark in the test to rely on the ordered logit model estimates, say  $\hat{\beta}_{FE,logit}$ . Frijters et al. (2004) argue that the alternative estimator based on a random effects probit specification should estimate a multiple of the same coefficient vector, so the working hypothesis would be  $\hat{\beta}_{FE,logit} = \alpha \hat{\beta}_{RE,probit}$ . They then propose a type of likelihood ratio test based on computation of the log likelihood functions for the two models. There are a number of problems with this approach, not least of which is that if the working hypothesis is true, it is necessary to estimate  $\alpha$ . However, the models are not nested, the parameters must necessarily be based on different sized samples and it is unclear what one should use for the degrees of freedom of the test if it were valid – the authors suggest  $K$ , the number of parameters in the model, but neither log likelihood forces  $K$  constraints on the other; the degrees of freedom for the  $LR$  test is the reduction in the number of dimensions of the parameter space. In this instance, the parameter space has  $K$  dimensions under both null and alternative. D’Addio, Eriksson and Frijters (2007) estimated a fixed effects ordered logit model and a random effects ordered probit model for “job satisfaction” for data from the European Community Household Panel and found that the fixed effects model was the preferred specification.

No other clearly appropriate procedure has been proposed. This problem is common to other nonlinear models. One strategy does suggest itself, based on the logic of the variable addition test [Wu (1973) and Baltagi (2007)]. In the random effects model to which we added the group means of the variables, the ostensible purpose of the variable addition was to account for correlation between the common effect,  $u_i$ , and the regressors. With that correlation present, the appropriate approach is fixed effects. Without that correlation, the random effects model is appropriate. Thus, while conceding that the power of the test is completely unknown at this point, we propose a simple likelihood ratio – variable addition test of the joint significance of the group means in the expanded random effects model.

Estimates of the fixed and random effects models are shown in Tables 26-28. For our estimated models we have  $\log L = -32656.89$  for the random effects model (Table 27) and  $-32588$  for the RE model with the group means added (Table 28). The likelihood ratio statistic for the hypothesis that the coefficients on the means are all zero is twice the difference, or 137.00, with 5 degrees of freedom. The hypothesis is decisively rejected, so we conclude that the fixed effects model is the preferred specification. Unfortunately, this now raises the question of how to fit the model. The average group size is less than 5. The results in Figure 14 suggest that the bias in the full MLE is as much as 30%. The results in Table 28 may be the appropriate ones.

**Table 26 Fixed Effects Ordered Logit Models**

Ordered Probability Model		FIXED EFFECTS OrdPrb Model	
Number of observations	27326	Number of observations	27326
Log likelihood function	-35853.13	Log likelihood function	-28818.86
Number of parameters	9	Number of parameters	5264
Info. Criterion: AIC =	2.62476	Info. Criterion: AIC =	2.49454
Restricted log likelihood	-36734.32	Unbalanced panel has 7293 individuals	
Underlying probabilities based on Logit		2037 groups with inestimable a(i)	

Pooled Estimates					
Variable	Coefficient	Standard Error	b/St.Er.	P[ Z >z]	Mean of X
-----+Index function for probability					
Constant	3.67149***	.08245660	44.526	.0000	
AGE	-.03546***	.00114889	-30.868	.0000	43.525690
EDUC	.06248***	.00506927	12.325	.0000	11.320631
INCOME	.45921***	.06715557	6.838	.0000	.3520836
MARRIED	.03593	.02976386	1.207	.2274	.7586182
KIDS	.09708***	.02651122	3.662	.0003	.4027300
-----+Threshold parameters for index					
Mu(1)	2.16706***	.01487458	145.689	.0000	
Mu(2)	4.35141***	.01500920	289.916	.0000	
Mu(3)	5.18118***	.01897982	272.983	.0000	

Full Maximum Likelihood Fixed Effects					
Variable	Coefficient	Standard Error	b/St.Er.	P[ Z >z]	Mean of X
-----+Index function for probability					
AGE	-.12834***	.00565661	-22.688	.0000	44.010737
EDUC	.01818	.05390774	.337	.7360	11.285486
INCOME	.49017***	.14415478	3.400	.0007	.3494867
MARRIED	.10852	.08234975	1.318	.1876	.7717663
KIDS	-.15489***	.05767212	-2.686	.0072	.4104794
-----+Threshold parameters for index					
MU(1)	3.55863***	.04903893	72.568	.0000	
MU(2)	7.15964***	.06023764	118.857	.0000	
MU(3)	8.51890***	.06462540	131.820	.0000	

Note: \*\*\*, \*\*, \* = Significance at 1%, 5%, 10% level.

Conditional Fixed Effects Logit, Binary: Healthy = 1(Health > 2)					
Variable	Mean	Std.Dev.	Minimum	Maximum	Cases
HEALTHY	.228830	.420087	.000000	1.00000	27326

Variable	Coefficient	Standard Error	b/St.Er.	P[ Z >z]
AGE	-.17204***	.00861180	-19.977	.0000
EDUC	.02126	.07518524	.283	.7773
INCOME	.49311**	.20894735	2.360	.0183
MARRIED	.18060	.11474008	1.574	.1155
KIDS	-.05835	.08169051	-.714	.4751

**Table 27 Random Effects Ordered Logit Models – Quadrature and Simulation**

Random Effects Ordered Prob. Model		Random Coefficients OrdProbs Model	
Number of observations	27326	Number of observations	27326
Log likelihood function	-32656.89	Log likelihood function	-32669.96
Info. Criterion: AIC =	2.39090	Info. Criterion: AIC =	2.39186
Info. Criterion: BIC =	2.39391	Info. Criterion: BIC =	2.39486
Info. Criterion:HQIC =	2.39187	Info. Criterion:HQIC =	2.39283
Unbalanced panel has 7293 individuals		Unbalanced panel has 7293 individuals	
-----+-----			
Quadrature based estimation			
-----+-----			
Variable	Coefficient	Standard Error	b/St.Er.  P[ Z >z]  Mean of X
-----+-----			
+-----+Index function for probability			
Constant	5.82480***	.16903183	34.460 .0000
AGE	-.06017***	.00209930	-28.660 .0000 43.525690
EDUC	.08299***	.01128254	7.355 .0000 11.320631
INCOME	.26636***	.09503935	2.803 .0051 .3520836
MARRIED	.12875***	.04732264	2.721 .0065 .7586182
KIDS	.01476	.03964475	.372 .7097 .4027300
+-----+Threshold parameters for index model			
Mu(01)	3.02273***	.03576261	84.522 .0000
Mu(02)	6.28777***	.04471783	140.610 .0000
Mu(03)	7.45137***	.04732226	157.460 .0000
+-----+Std. Deviation of random effect			
Sigma	1.79351***	.02423137	74.016 .0000
-----+-----			
Simulation based estimation: 100 Halton draws			
-----+-----			
+-----+Means for random parameters			
Constant	5.78689***	.09391702	61.617 .0000
+-----+Nonrandom parameters			
AGE	-.05944***	.00123424	-48.162 .0000 43.525690
EDUC	.08378***	.00541280	15.478 .0000 11.320631
INCOME	.25495***	.06941556	3.673 .0002 .3520836
MARRIED	.12251***	.03049235	4.018 .0001 .7586182
KIDS	.01577	.02739330	.576 .5648 .4027300
+-----+Scale parameters for dists. of random parameters			
Constant	1.81125***	.01529676	118.407 .0000
+-----+Threshold parameters for probabilities			
MU(1)	3.01553***	.03279390	91.954 .0000
MU(2)	6.28238***	.04054372	154.953 .0000
MU(3)	7.44468***	.04297824	173.220 .0000
-----+-----			
Note: ***, **, * = Significance at 1%, 5%, 10% level.			
-----+-----			

**Table 28 Random Effects Model with Mundlak Correction**

Random Effects Ordered Probability Model					
Log likelihood function		-32588.39			
Number of parameters		15			
Info. Criterion: AIC =		2.38625			
Variable	Coefficient	Standard Error	b/St.Er.	P[ Z >z]	Mean of X
+-----+Index function for probability					
Constant	5.01093***	.18899294	26.514	.0000	
AGE	-.10573***	.00468085	-22.587	.0000	43.525690
EDUC	.02040	.05476626	.373	.7095	11.320631
INCOME	.38927***	.12115272	3.213	.0013	.3520836
MARRIED	.09947	.07063476	1.408	.1591	.7586182
KIDS	-.12489**	.05066166	-2.465	.0137	.4027300
AGEBAR	.05909***	.00530378	11.140	.0000	43.525690
EDUCBAR	.06300	.05588876	1.127	.2596	11.320631
INCBAR	.62547***	.20296625	3.082	.0021	.3520836
MARRBAR	-.11753	.09893186	-1.188	.2348	.7586182
KIDSBAR	.35591***	.08707387	4.087	.0000	.4027300
+-----+Threshold parameters for index model					
Mu(01)	3.02621***	.03570022	84.767	.0000	
Mu(02)	6.30011***	.04476340	140.742	.0000	
Mu(03)	7.46994***	.04747553	157.343	.0000	
+-----+Std. Deviation of random effect					
Sigma	1.79092***	.02412969	74.220	.0000	
Note: ***, **, * = Significance at 1%, 5%, 10% level.					

## 6.4 Extending Parameter Heterogeneity Models to Ordered Choices

Based on the results of the previous sections, the extension of the models with parameter heterogeneity involves only a minor change in the log likelihood and essentially none in the interpretation of the model. For example, in the random parameters model, the heterogeneity in the parameters is the same as in the random effect – it is useful to view the random effects model as a random parameters model in which only the constant term is random. The more general model is

$$\log L_i = \log \int_{\mathbf{w}_i} \prod_{t=1}^{T_i} \left[ \Phi(\mu_j - \beta' \mathbf{x}_i - (\mathbf{L}\mathbf{D}\mathbf{w}_i)' \mathbf{x}_i) - \Phi(\mu_{j-1} - \beta' \mathbf{x}_i - (\mathbf{L}\mathbf{D}\mathbf{w}_i)' \mathbf{x}_i) \right] F(\mathbf{w}_i) d\mathbf{w}_i.$$

The log likelihood for the sample is once again the sum over the  $N$  joint observations. The integration can now be replaced with a simulation over  $R$  draws from the multivariate standard normal population. The simulated log likelihood is, then

$$\log L = \sum_{i=1}^N \log \frac{1}{R} \sum_{r=1}^R \prod_{t=1}^{T_i} \left[ \Phi(\mu_j - \beta' \mathbf{x}_{it} - (\mathbf{L}\mathbf{D}\mathbf{w}_{ir})' \mathbf{x}_{it}) - \Phi(\mu_{j-1} - \beta' \mathbf{x}_{it} - (\mathbf{L}\mathbf{D}\mathbf{w}_{ir})' \mathbf{x}_{it}) \right].$$

The generalized ordered choice model (3) and the latent class model are handled similarly. For the first,

$$\begin{aligned} \log L_S(\beta, \Delta, \alpha, \delta, \gamma, \mathbf{L}, \mathbf{D}, \sigma, \tau) = \\ \sum_{i=1}^n \log \frac{1}{R} \sum_{r=1}^R \prod_{t=1}^{T_i} \left( F \left[ \frac{\mu_{ij,r} - \beta'_{i,r} \mathbf{x}_{it}}{\sqrt{\exp(\gamma' \mathbf{h}_i + \tau e_{i,r})}} \right] - F \left[ \frac{\mu_{i,j-1,r} - \beta'_{i,r} \mathbf{x}_{it}}{\sqrt{\exp(\gamma' \mathbf{h}_i + \tau e_{i,r})}} \right] \right) \end{aligned}$$

As before, the structure assumes that the heterogeneity is constant through time. For the latent class model, the appropriate log likelihood function is

$$\log L = \sum_{i=1}^N \log \left\{ \sum_{q=1}^Q \frac{\exp(\theta_q + \delta'_q \mathbf{z}_i)}{\sum_{q=1}^Q \exp(\theta_q + \delta'_q \mathbf{z}_i)} \left( \prod_{t=1}^{T_i} \sum_{j=0}^J m_{ij} \left[ \Phi(\mu_{j,q} - \beta'_q \mathbf{x}_{it}) - \Phi(\mu_{j-1,q} - \beta'_q \mathbf{x}_{it}) \right] \right) \right\}$$

The counterpart to the assumption of time invariant heterogeneity is the assumption that the class membership is the same in every period.

Random parameters and latent class estimates for the health care model are shown in Tables 29-31. The latent class model is fit with the full panel data set in Table 30, then with the cross section used previously (4,483 observations) in Table 31. The estimates are relatively stable across the two samples. However, the benefit from the larger sample is clearly visible in the much smaller standard errors in Table 30.

**Table 29 Random Parameters Ordered Logit Model**

Variable	Coefficient	Standard Error	b/St.Er.	P[ Z >z]	Mean of X
Random Coefficients Ordered Choice Model					
Ordered LOGIT probability model					
Dependent variable	HEALTH				
Number of observations	27326				
Log likelihood function	-32895.56				
Number of parameters	15				
Info. Criterion: AIC =	2.40874				
Info. Criterion: BIC =	2.41325				
Info. Criterion:HQIC =	2.41019				
Unbalanced panel has	7293 individuals.				
LHS variable = values	0,1,..., 4				
Simulation based on	20 Halton draws				
-----					
+-----+Means for random parameters					
Constant	5.49422***	.08843554	62.127	.0000	
AGE	-.05772***	.00118066	-48.892	.0000	43.525690
EDUC	.09802***	.00530070	18.491	.0000	11.320631
INCOME	.20420***	.06745335	3.027	.0025	.3520836
MARRIED	.15823***	.02897652	5.461	.0000	.7586182
KIDS	-.00095	.02670638	-.036	.9717	.4027300
+-----+Scale parameters for dists. of random parameters					
Constant	.03922***	.01226752	3.197	.0014	
AGE	.02561***	.00028746	89.092	.0000	
EDUC	.10451***	.00115944	90.135	.0000	
INCOME	.04246	.02922949	1.453	.1463	
MARRIED	.28916***	.01322763	21.861	.0000	
KIDS	.55735***	.01825860	30.525	.0000	
+-----+Threshold parameters for probabilities					
MU(1)	2.96883***	.03154534	94.113	.0000	
MU(2)	6.17657***	.03942745	156.656	.0000	
MU(3)	7.31328***	.04195129	174.328	.0000	
-----					
Pooled					
Log likelihood function	-35853.13				
-----					
+-----+Index function for probability					
Constant	3.67149***	.08245660	44.526	.0000	
AGE	-.03546***	.00114889	-30.868	.0000	43.525690
EDUC	.06248***	.00506927	12.325	.0000	11.320631
INCOME	.45921***	.06715557	6.838	.0000	.3520836
MARRIED	.03593	.02976386	1.207	.2274	.7586182
KIDS	.09708***	.02651122	3.662	.0003	.4027300
+-----+Threshold parameters for index					
Mu(1)	2.16706***	.01487458	145.689	.0000	
Mu(2)	4.35141***	.01500920	289.916	.0000	
Mu(3)	5.18118***	.01897982	272.983	.0000	
-----					
Note: ***, **, * = Significance at 1%, 5%, 10% level.					

**Table 30 Panel Data Latent Class Ordered Logit Model**

Latent Class / Panel OrdProbs Model						
Dependent variable	HEALTH					
Number of observations	27326					
Log likelihood function	-32639.79					
Info. Criterion: AIC =	2.39148					
Unbalanced panel has	7293 individuals.					
LHS variable = values	0,1,..., 4					
Variable	Coefficient	Standard Error	b/St.Er.	P[ Z >z]	Mean of X	
-----+Model parameters for latent class 1						
Constant	6.19492***	.27781573	22.299	.0000		
AGE	-.03605***	.00258199	-13.961	.0000	43.525690	
EDUC	.05712***	.01386691	4.119	.0000	11.320631	
INCOME	-.44011***	.13235483	-3.325	.0009	.3520836	
MARRIED	-.01197	.06038537	-.198	.8428	.7586182	
KIDS	.00839	.05659850	.148	.8821	.4027300	
MU(1)	2.28697***	.18388040	12.437	.0000		
MU(2)	4.63902***	.19146805	24.229	.0000		
MU(3)	5.66260***	.19099199	29.648	.0000		
-----+Model parameters for latent class 2						
Constant	2.66795***	.17559502	15.194	.0000		
AGE	-.04862***	.00257330	-18.894	.0000	43.525690	
EDUC	.06875***	.01098089	6.261	.0000	11.320631	
INCOME	.72058***	.14793773	4.871	.0000	.3520836	
MARRIED	.22071***	.06063146	3.640	.0003	.7586182	
KIDS	.02336	.05749856	.406	.6845	.4027300	
MU(1)	2.64040***	.04168310	63.345	.0000		
MU(2)	5.01676***	.08258391	60.747	.0000		
MU(3)	5.53435***	.10119982	54.687	.0000		
-----+Model parameters for latent class 3						
Constant	6.33422***	.23744548	26.677	.0000		
AGE	-.05911***	.00243588	-24.266	.0000	43.525690	
EDUC	.10674***	.00999348	10.681	.0000	11.320631	
INCOME	.26003**	.13089116	1.987	.0470	.3520836	
MARRIED	.13967**	.05446710	2.564	.0103	.7586182	
KIDS	-.00245	.04831235	-.051	.9595	.4027300	
MU(1)	3.66469***	.15664245	23.395	.0000		
MU(2)	7.23433***	.16902975	42.799	.0000		
MU(3)	8.68613***	.18133685	47.901	.0000		
-----+Estimated prior probabilities for class membership						
ONE_1	-.59260***	.10747644	-5.514	.0000		
FEMALE_1	-.03111	.08931083	-.348	.7276		
HANDDU_1	-.72480***	.16734906	-4.331	.0000		
WORKIN_1	-.06869	.09598755	-.716	.4742		
ONE_2	-.74731***	.10263543	-7.281	.0000		
FEMALE_2	.22391***	.08683560	2.579	.0099		
HANDDU_2	1.12965***	.10813499	10.447	.0000		
WORKIN_2	-.30028***	.09058678	-3.315	.0009		
ONE_3	.000***	.....(Fixed Parameter).....				
FEMALE_3	.000***	.....(Fixed Parameter).....				
HANDDU_3	.000***	.....(Fixed Parameter).....				
WORKIN_3	.000***	.....(Fixed Parameter).....				
Prior class probabilities at data means for LCM variables						
Class 1	Class 2	Class 3				
.22339	.27489	.50172				
-----+Note: ***, **, * = Significance at 1%, 5%, 10% level.						

**Table 31 Cross Section Ordered Logit Model**

Latent Class / Panel OrdProbs Model					
Dependent variable	HEALTH				
Number of observations	4483				
Log likelihood function	-5743.560				
Info. Criterion: AIC =	2.57799				
Sample is	1 pds and 4483 individuals.				
Ordered LOGIT probability model					
LHS variable = values 0,1,..., 4					
Variable	Coefficient	Standard Error	b/St.Er.	P[ Z >z]	Mean of X
-----+Model parameters for latent class 1					
Constant	4.84807***	1.10567170	4.385	.0000	
AGE	-.06426***	.01507410	-4.263	.0000	43.440107
EDUC	.09308	.07057360	1.319	.1872	11.418086
INCOME	.53644	.58467594	.918	.3589	.3487401
MARRIED	-.80972	1.27735593	-.634	.5261	.7521749
KIDS	.17360	.46003096	.377	.7059	.3794334
MU(1)	1.66526***	.49271119	3.380	.0007	
MU(2)	4.62309***	.58324699	7.926	.0000	
MU(3)	6.68726***	1.17391562	5.697	.0000	
-----+Model parameters for latent class 2					
Constant	4.89283***	1.41045138	3.469	.0005	
AGE	-.06472***	.02418431	-2.676	.0074	43.440107
EDUC	.01884	.09496112	.198	.8427	11.418086
INCOME	.86153	.91919502	.937	.3486	.3487401
MARRIED	1.13731	2.02882846	.561	.5751	.7521749
KIDS	-.41770	.64784855	-.645	.5191	.3794334
MU(1)	1.17500	.76043717	1.545	.1223	
MU(2)	4.69223***	1.17503658	3.993	.0001	
MU(3)	6.09581***	1.33046471	4.582	.0000	
-----+Model parameters for latent class 3					
Constant	4.71508***	1.47524161	3.196	.0014	
AGE	-.01366	.02134897	-.640	.5222	43.440107
EDUC	-.05643	.07540252	-.748	.4542	11.418086
INCOME	.26151	.68196672	.383	.7014	.3487401
MARRIED	.15012	.41137976	.365	.7152	.7521749
KIDS	-.30811	.30457558	-1.012	.3117	.3794334
MU(1)	4.04236**	1.89743084	2.130	.0331	
MU(2)	4.53578***	1.29161229	3.512	.0004	
MU(3)	4.53579***	1.23514521	3.672	.0002	
-----+Estimated prior probabilities for class membership					
ONE_1	.27149	1.70617793	.159	.8736	
FEMALE_1	.18836	.31092443	.606	.5446	
HANDDU_1	-.36330	.36481811	-.996	.3193	
WORKIN_1	.63905*	.38201305	1.673	.0944	
ONE_2	.50734	1.59323236	.318	.7502	
FEMALE_2	-.16926	.38324861	-.442	.6588	
HANDDU_2	-.52187	.39614098	-1.317	.1877	
WORKIN_2	.17198	.48670134	.353	.7238	
ONE_3	.000	.....(Fixed Parameter).....			
FEMALE_3	.000	.....(Fixed Parameter).....			
HANDDU_3	.000	.....(Fixed Parameter).....			
WORKIN_3	.000	.....(Fixed Parameter).....			
Prior class probabilities at data means for LCM variables					
	Class 1	Class 2	Class 3	Class 4	Class 5
	.44728	.34178	.21094	.00000	.00000
-----					
Note: ***, **, * = Significance at 1%, 5%, 10% level.					

## 7 Extensions

The preceding sections have examined the more or less standard approaches to modeling ordered data, beginning with the most basic model and ending with various specifications that accommodate observed and unobserved heterogeneity in panel data. In what follows, we examine some recent extensions of the model that include modifications to the basic structure and additions to it that occasionally mandate multiple equation frameworks. It will emerge shortly that most of these extensions do not fit comfortably into the ordered logit framework. At this point, it will prove convenient to drop the distinction between the probit and logit models, and focus attention, as in the received literature, on the ordered probit model.

### 7.1 Dynamic Models

Dynamic effects in ordered choice models have been introduced in two settings. In the pure time series applications in which researchers have examined asset price movements, interest rate changes and monetary policy, the focus is on inertia, and takes the form of an autoregressive model in the latent variable regression. The Czado, Heyn and Müller (2005) and Müller and Czado (2005) study of migraine headache severity is also presented in this framework, though their study can be usefully viewed as falling somewhere between the time series analysis of, e.g., Eichengreen et al.'s (1985) study of bank rate policy and the recent panel data studies, e.g., of health satisfaction. In panel data settings, such as Contoyannis, Jones and Rice (2004), the model is directed at state dependence, and, instead, takes the form of lagged effects in the observed variables. We will examine each of these in a bit more detail.

A natural form of the ordered probit model with lagged effects is suggested by Girard and Parent (2001),

$$\begin{aligned} y_t^* &= \boldsymbol{\beta}'\mathbf{x}_t + \varepsilon_t, \\ \varepsilon_t &= \rho\varepsilon_{t-1} + u_t \\ y_t &= j \text{ if } \mu_{j-1} < y_t^* \leq \mu_j \end{aligned}$$

with the usual restrictions. Estimation is carried using a Gibbs sampler (MCMC) and using Albert and Chib's data augmentation method; the values  $y_t^*$  as well as the initial value,  $y_0^*$  are treated as nuisance parameters to be included with  $\boldsymbol{\beta}$ ,  $\boldsymbol{\mu}$  and  $\rho$  for posterior analysis.

Eichengreen, Watson and Grossman (1985) examined the Bank Rate (BR) adjustment policies of the Bank of England over a period of 328 weeks. The structural model is

$$\begin{aligned} \text{Prob}[\Delta BR_t = -50 | J_t] &= P_{1t}(J_t) \\ \text{Prob}[\Delta BR_t = 0 | J_t] &= P_{2t}(J_t) \\ \text{Prob}[\Delta BR_t = 100 | J_t] &= P_{3t}(J_t), t = 1, \dots, T \end{aligned}$$

where the adjustment rates are in basis points and  $J_t$  is an information set that contains current and lagged values of exogenous variables  $\mathbf{x}_t$  and the entire preceding history of bank rates,  $BR_s$ ,  $s = 1, \dots, t-1$ . An underlying regression is specified for the "change in an unobserved "underlying" bank rate,

$$\Delta BR_t^* = \boldsymbol{\beta}'\mathbf{x}_t + \varepsilon_t, \quad \varepsilon_t | J_t \sim N[0, \sigma^2].$$

The observed Bank Rate changes when it is too far from  $BR_t^*$  according to the rule,

$$\begin{aligned}\Delta BR_t &= -50 \text{ if } BR_t^* < BR_{t-1} - \alpha_L \\ \Delta BR_t &= 0 \text{ if } BR_{t-1} - \alpha_L < BR_t^* < BR_{t-1} + \alpha_U \\ \Delta BR_t &= 100 \text{ if } BR_t^* > BR_{t-1} + \alpha_U.\end{aligned}$$

Thus, the rule is that the observed rate decreases by 50 basis points if  $BR_t^*$  is “appreciably” less than  $BR_{t-1}$  and increases by 100 basis points if  $BR_t^*$  is appreciably greater than  $BR_{t-1}$ . Appreciably is defined by the unknown threshold values,  $\alpha_L$  and  $\alpha_U$ . The authors note, the model resembles a familiar ordered probit model, but differs in at least two major respects. First, although the structural equations describe the changes in  $BR$ , the inequalities that invoke the similarity with the ordered probit model are defined in the levels of  $BR$ , not changes. Thus, there are stochastic dynamics in  $BR_t$ . Second, since the lagged value of the observed time series appear in the model definition, the identification of the model parameters must be developed in detail. It does not follow from simple examination of the specification as it does in the conventional model. The likelihood function (see their pp. 741-744) is markedly more complicated than that we have examined so far. Among the most challenging aspects is that because of the autoregressive nature of the random components in the model, the time series must be treated as a single  $T$ -variate observation. That implies integration of a  $T$  ( $=328$ ) variate normal integral. A strategy is devised in the paper. Eichengreen et al.’s (1985) study has provided the foundation for a number of subsequent studies of bank policy, including Genberg and Gerlach (2004) and Basu and de Jong (2006).

A somewhat simpler form of the ordered probit model has been used to analyze movements in stock prices when the movements of an underlying continuous price variable are expressed in discrete units (“ticks”). Tsay (2002) presents the following general characterization of an application: Define  $y_{it}^*$  to be the unobservable true price change of an asset, so that

$$y_{it}^* = P_{it}^* - P_{i,t-1}^*,$$

where  $P_{it}^*$  is the virtual price of the asset at time  $t$ . The ordered probit model derives from the assumed structure

$$\begin{aligned}y_{it}^* &= \boldsymbol{\beta}'\mathbf{x}_{it} + \varepsilon_{it} \\ E[\varepsilon_{it}|\mathbf{x}_{it}, \mathbf{w}_{it}] &= 0 \\ \text{Var}[\varepsilon_{it}|\mathbf{x}_{it}, \mathbf{w}_{it}] &= \sigma^2(\mathbf{w}_{it}) \\ \varepsilon_{it}|\mathbf{x}_{it}, \mathbf{w}_{it} &\sim N[0, \sigma^2(\mathbf{w}_{it})],\end{aligned}$$

where  $\mathbf{x}_{it}$  might contain the exogenously determined information available at time  $t-1$  and  $\mathbf{w}_{it}$  is conditioning data such as the time interval of the change as well as “some conditional[ly] heteroscedastic variables.” If the observed price change is restricted to a fixed set of intervals, then an ordered probit model emerges;

$$y_{it} = s_j \text{ if } \alpha_{j-1} < y_{it}^* \leq \alpha_j, j = 1, \dots, J$$

What follows is a familiar ordered probit model, distinguished from our earlier model by the assumed heteroscedasticity of  $\varepsilon_{it}$ . Tsay describes in detail an early study of more than 100 stocks by Hausman, Lo and MacKinlay (1992). Hausman et al. describe three features of the American stock market that motivate their treatment: First, stock prices were stated at the time (no longer)

in discrete, 1/8 dollar units, so the true continuous variable could not be measured. Second, the timing of transactions can be irregular and random, which makes discrete time modeling problematic. Third, received models have not adequately accounted for the correlations between price changes and other economic variables – these are captured in the latent regression equation in the ordered probit model.

Czado, Heyn and Müller (2005) also used a time series model with dynamics in the latent variable to study the reported severity of migrain headaches reported in the diary of a single patient. The underlying variable, severity of the headache in interval  $t$ , is modeled

$$y_t^* = \boldsymbol{\beta}'\mathbf{x}_t + \gamma y_{t-1}^* + \varepsilon_t$$

The observed severity is recorded on a scale 0,1,...,5, four times per day over a period of 268 days. The regressor variable includes such variables as weather conditions and day of the week. The application is a pure time series model. As in the Eighengreen et al. study, the dynamics greatly complicate the estimation process. A customized form of Markov Chain Monte Carlo (Bayesian) estimation method for this model is presented in Müller and Czado (2005).

The autoregressive models examined so far are natural specifications for the observed outcomes. Contoyannis, Jones and Rice (2004) examined self assessed health status in the British Household Panel Survey (BHPS). The measure of health status is reported with values 1,...,5. Individuals have a general tendency to repeat the same value unless other factors change. The common effects regression suggested to account for this state dependence is

$$h_{it}^* = \boldsymbol{\beta}'\mathbf{x}_{it} + \sum_{j=1}^5 \gamma_j m_{j,i,t-1} + \alpha_i + \varepsilon_{it}$$

where  $\alpha_i$  is a fixed effect and

$$m_{i,j,t-1} = 1 \text{ if } y_{i,t-1} = j \text{ and } 0 \text{ otherwise.}$$

Ostensibly, a familiar ordered probit model applies;

$$h_{it} = j \text{ iff } \mu_{j-1} < y_{it}^* \leq \mu_j.$$

Initially, it is proposed to treat this as a random effects model using the method of Butler and Moffitt (1982). In order to accommodate possible correlation between  $\alpha_i$  and the (means of the) other variables and to handle the problem of the initial conditions [Heckman (1981)], they employ the Mundlak (1978) device in:

$$\alpha_i = \alpha_0 + \sum_{j=1}^5 \alpha_j m_{i,1,j} + \boldsymbol{\theta}'\bar{\mathbf{x}}_i + u_i.$$

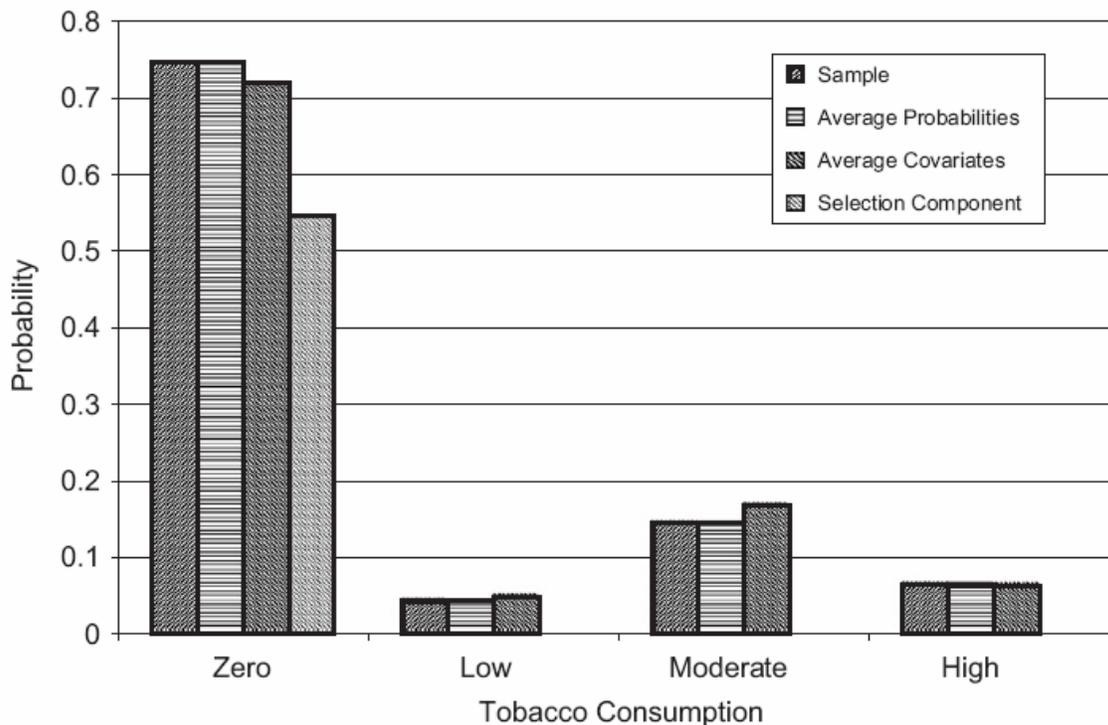
where  $u_i \sim N[0, \sigma^2]$ . Inserting this equation into the latent regression provides their ordered probit model,

$$h_{it}^* = \boldsymbol{\beta}'\mathbf{x}_{it} + \sum_{j=1}^5 \gamma_j m_{j,i,t-1} + \alpha_0 + \sum_{j=1}^5 \alpha_j m_{i,1,j} + \boldsymbol{\theta}'\bar{\mathbf{x}}_i + u_i + \varepsilon_{it}.$$

(A few normalizations, such as removal of a redundant constant term, are needed to secure identification of the parameters.) A final adjustment to the model based on a procedure devised by Wooldridge (2002) is used to account for the rather substantial attrition over the 8 waves of their panel.

## 7.2 Inflation Models

Harris and Zhao (2007) analyzed a sample of 28,813 Australian individuals' responses to the question "How often do you now smoke cigarettes, pipes or other tobacco products?" [Data are from the Australian National Drug Strategy Household Survey, NDSHS (2001).] Responses were "zero, low, moderate, high," coded 0,1,2,3. Figure 15 below reproduces their Figure 3 (page 1095). The leftmost bar of each set shows the sample histogram. The spike at zero shows a considerable excess of zeros compared to what might be expected in an ordered choice model. The authors reason that there are numerous explanations for a zero response: "genuine nonsmokers, recent quitters, infrequent smokers who are not currently smoking and potential smokers who might smoke when, say, the price falls." It is also possible that the zero response includes some individuals who prefer to identify themselves as noonsmokers. The question is ambiguously worded, but arguably, the group of interest is the genuine nonsmokers. This suggests a type of latent class arrangement in the population. There are (arguably) two types of zeros, the one of interest, and another type generated by the appearance of the respondent in the latent class of people who respond zero when another response would actually be appropriate. The end result is an inflation of the proportion of zero responses in the data. The "Zero Inflation" model is proposed to accommodate this failure of the base case model.



**Figure 15 Tobacco Consumption Survey and Model Results**

Zero inflation as a formal model to explain data such as these originates in Lambert's (1992) study of quality control in industry. Sampling for defectives in a production process can produce two types of zeros (per unit of time). The process may be under control, or it may be out of control and the observer happens to draw zero defectives in a particular sample. This inflates the number of zeros in a sample beyond what would be expected by a count model such as the Poisson model – the modification named the ZIP (zero inflated) or ZAP (zero altered) Poisson model. [See also Heilbron (1994), Hinde et al. (1998), Mullahy (1997) and Greene (1994).]

Harris and Zhao proposed the following zero inflated ordered probit (ZIOP) model:

*Participation equation:*

Regime 0 for nonparticipation (nonsmoker), Regime 1 for participation

$$r^* = \boldsymbol{\alpha}'\mathbf{z} + u, u \sim N[0,1]$$

$$r = 1 \text{ if } r^* > 0, 0 \text{ otherwise}$$

$$\text{Prob}(r = 1 | \mathbf{z}) = \Phi(\boldsymbol{\alpha}'\mathbf{z}).$$

*Activity equation*

$$y^* = \boldsymbol{\beta}'\mathbf{x} + \varepsilon, \varepsilon \sim N[0,1], \text{ independent of } u,$$

$$y = j \text{ if } \mu_{j-1} < y^* \leq \mu_j, j = 0, 1, \dots, J.$$

(At the risk of some confusion below, we have modified Harris and Zhao's notation to conform to the conventions we have used up to this point.) Thus, a standard probit model governs participation and our familiar ordered probit model governs "true" activity. The observed activity level, however, is not  $y$ . It is

$$y_o = r \times y.$$

A nonparticipant reports a zero as well as some participants. Thus, the zero outcome occurs when  $r = 0$  and when  $r = 1$  and  $y = 0$ . Therefore, the zero outcome is inflated by the  $r = 0$  regime. The applicable probabilities for the observed outcome are

$$\text{Prob}(y_o = 0 | \mathbf{x}, \mathbf{z}) = \text{Prob}(r = 0 | \mathbf{z}) + \text{Prob}(r = 1 | \mathbf{z}) \times \text{Prob}(y = 0 | \mathbf{x}, r = 1)$$

$$\text{Prob}(y_o = j | \mathbf{x}, \mathbf{z}) = \text{Prob}(r = 1 | \mathbf{z}) \times \text{Prob}(y = j | \mathbf{x}, r = 1).$$

Note at this point, by dint of the independence of  $\varepsilon$  and  $u$ ,  $\text{Prob}(y = 0 | \mathbf{x}, r = 1) = \text{Prob}(y = 0 | \mathbf{x})$ . We will relax this assumption later.

With the assumption of normality of  $\varepsilon$  and  $u$ , the associated probabilities are obtained from those of the binary probit model and the ordered probit model;

$$\text{Prob}(y_o = 0 | \mathbf{x}, \mathbf{z}) = [1 - \Phi(\boldsymbol{\alpha}'\mathbf{z})] + \Phi(\boldsymbol{\alpha}'\mathbf{z}) \times \Phi(0 - \boldsymbol{\beta}'\mathbf{x})$$

$$\text{Prob}(y_o = j | \mathbf{x}, \mathbf{z}) = \Phi(\boldsymbol{\alpha}'\mathbf{z}) \times [\Phi(\mu_j - \boldsymbol{\beta}'\mathbf{x}) - \Phi(\mu_{j-1} - \boldsymbol{\beta}'\mathbf{x})], j = 1, \dots, J$$

with the same normalization as earlier,  $\mu_{-1} = -\infty$ ,  $\mu_0 = 0$ ,  $\mu_J = +\infty$ . The log likelihood function is built up as the sum of the logs of the probabilities of the observed outcomes.

An extension which would seem to be appropriate for this application is to allow the unobserved effects in the participation equation and the activity equation to be correlated (producing a ZIOPC model). Thus, we now have

$$\begin{pmatrix} u \\ \varepsilon \end{pmatrix} \sim N \left[ \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \right].$$

The correlation coefficient,  $\rho$ , is now an additional parameter to be estimated. With this modification, we no longer have  $\text{Prob}(y = 0 | \mathbf{x}, r = 1) = \text{Prob}(y = 0 | \mathbf{x})$ ; the former is now a

probability from the bivariate normal distribution. The probabilities of the observed outcomes become

$$\text{Prob}(y_o = 0 \mid \mathbf{x}, \mathbf{z}) = [1 - \Phi(\boldsymbol{\alpha}'\mathbf{z})] + \Phi_2(\boldsymbol{\alpha}'\mathbf{z}, -\boldsymbol{\beta}'\mathbf{x}, -\rho)$$

$$\text{Prob}(y_o = j \mid \mathbf{x}, \mathbf{z}) = \Phi_2(\boldsymbol{\alpha}'\mathbf{z}, \mu_j - \boldsymbol{\beta}'\mathbf{x}, -\rho) - \Phi_2(\boldsymbol{\alpha}'\mathbf{z}, \mu_{j-1} - \boldsymbol{\beta}'\mathbf{x}, -\rho), j = 1, \dots, J$$

where  $\Phi_2(\dots)$  denotes the probability of a joint event from the bivariate normal cdf. This modification drastically alters the partial effects in the model. To organize these in a convenient fashion, we adopt the authors' device. Let  $\mathbf{x}^* = (\mathbf{x}_o, \mathbf{x}_c, \mathbf{z}_o)$  so that  $\mathbf{x}_o$  is variables in  $\mathbf{x}$  that are not also in  $\mathbf{z}$ ,  $\mathbf{x}_c$  is variables that are in both  $\mathbf{x}$  and  $\mathbf{z}$ , and  $\mathbf{z}_o$  is variables in  $\mathbf{z}$  that are not in  $\mathbf{x}$ . By rearranging and reordering the parameter vectors,  $\boldsymbol{\beta}$  and  $\boldsymbol{\alpha}$  into  $\boldsymbol{\beta}^* = (\boldsymbol{\beta}_o, \boldsymbol{\beta}_c, \mathbf{0})$  and  $\boldsymbol{\alpha}^* = (\mathbf{0}, \boldsymbol{\alpha}_c, \boldsymbol{\alpha}_o)$ , then  $\boldsymbol{\beta}'\mathbf{x} = \boldsymbol{\beta}^*\mathbf{x}^*$  and  $\boldsymbol{\alpha}'\mathbf{z} = \boldsymbol{\alpha}^*\mathbf{x}^*$ . We can thus obtain the partial effects by differentiating with respect to  $\mathbf{x}^*$  and obtaining the needed decomposition. Then, with this in place,

$$\begin{aligned} \frac{\partial \text{Prob}(y_o = 0 \mid \mathbf{x}^*)}{\partial \mathbf{x}^*} &= \left[ \Phi \left( \frac{-\boldsymbol{\beta}'\mathbf{x} + \rho\boldsymbol{\alpha}'\mathbf{z}}{\sqrt{1-\rho^2}} \right) - 1 \right] \phi(\boldsymbol{\alpha}'\mathbf{z})\boldsymbol{\alpha}^* - \Phi \left( \frac{\boldsymbol{\alpha}'\mathbf{z} - \rho\boldsymbol{\beta}'\mathbf{x}}{\sqrt{1-\rho^2}} \right) \phi(\boldsymbol{\beta}'\mathbf{x})\boldsymbol{\beta}^* \\ \frac{\partial \text{Prob}(y_o = j \mid \mathbf{x}^*)}{\partial \mathbf{x}^*} &= \left[ \Phi \left( \frac{\mu_j - \boldsymbol{\beta}'\mathbf{x} + \rho\boldsymbol{\alpha}'\mathbf{z}}{\sqrt{1-\rho^2}} \right) - \Phi \left( \frac{\mu_{j-1} - \boldsymbol{\beta}'\mathbf{x} + \rho\boldsymbol{\alpha}'\mathbf{z}}{\sqrt{1-\rho^2}} \right) \right] \phi(\boldsymbol{\alpha}'\mathbf{z})\boldsymbol{\alpha}^* \\ &\quad + \left[ \phi(\mu_{j-1} - \boldsymbol{\beta}'\mathbf{x})\Phi \left( \frac{\boldsymbol{\alpha}'\mathbf{z} + \rho(\mu_{j-1} - \boldsymbol{\beta}'\mathbf{x})}{\sqrt{1-\rho^2}} \right) - \phi(\mu_j - \boldsymbol{\beta}'\mathbf{x})\Phi \left( \frac{\boldsymbol{\alpha}'\mathbf{z} + \rho(\mu_j - \boldsymbol{\beta}'\mathbf{x})}{\sqrt{1-\rho^2}} \right) \right] \boldsymbol{\beta}^* \end{aligned}$$

These results are likely to bear little resemblance to the raw coefficients, particularly for variables which appear in both equations.

Testing the null hypothesis of the ZIOP model against the alternative of the ZIOPC model is a simple test of the hypothesis that  $\rho$  equals zero. This can be done using a Wald ( $t$ ) test or a likelihood ratio test. Testing for the inflation effects is more complicated however. The obvious restriction,  $\boldsymbol{\alpha} = \mathbf{0}$ , does not remove the inflation effect; it makes the regime probabilities both equal to one half. What is needed to remove the inflation effect is  $\boldsymbol{\alpha}'\mathbf{z} \rightarrow \infty$ , which cannot be imposed. The hypotheses are not nested. Greene (1994) proposed using the Vuong (1989) test for this hypothesis. Denote the probability for the observed outcome from the inflation model as  $f_I(y_o, r \mid \mathbf{x}, \mathbf{z})$  and that for the uninflated model as  $f_U(y_o, r \mid \mathbf{x}, \mathbf{z})$ . Then,

$$m_i = \log \left( \frac{f_I(y_{o,i}, r_i \mid \mathbf{x}_i, \mathbf{z}_i)}{f_U(y_{o,i}, r_i \mid \mathbf{x}_i, \mathbf{z}_i)} \right).$$

The test statistic is

$$V = \frac{\sqrt{N} (1/N) \sum_{i=1}^N m_i}{\sqrt{(1/N) \sum_{i=1}^N (m_i - \bar{m})^2}} = \frac{\sqrt{N} \bar{m}}{s_m}$$

The limiting distribution of  $V$  under the null hypothesis of no difference is  $N(0,1)$ . The test is directional. Large positive values favor the inflation model; large negative values favor the uninflated model. The inconclusive region for a 5% significance level would be  $(-1.96, +1.96)$ .

Given the greater number of parameters in the inflation model, it will be rare for  $V$  to be strongly negative. It will often strongly favor the larger model.

Brooks, Harris and Spencer (2007) applied the same style of analysis to the policy decisions of the members of the Bank of England Monetary Policy Committee. In this study, the participation equation is a decision to adjust monetary policy (at all). The activity equation is whether rates should decrease ( $y_o = 0$ ), stay the same ( $y_o = 1$ ) or increase ( $y_o = 2$ ). (The model of Eichengreen, Watson and Grossman (1985) is developed on this logic as well.) In this case, the no change result can occur because of a decision not to change rates, or by an inclination to change rates followed later by a decision not to. Thus, the model produces “one inflation.”

### 7.3 Multiple Equations

A multiple equation specification for, say,  $M$  ordered choices is a natural extension of the model. The extension is based on a seemingly unrelated regressions (SUR) model for the latent regressions:

$$\begin{aligned}
 y_{i,1}^* &= \beta_1' \mathbf{x}_{i,1} + \varepsilon_{i,1}, y_{i,1} = j \text{ if } \mu_{j-1,1} < y_{i,1}^* < \mu_{j,1}, \varepsilon_{i,1} \sim N[0,1], \\
 \dots \\
 y_{i,M}^* &= \beta_M' \mathbf{x}_{i,M} + \varepsilon_{i,M}, y_{i,M} = j \text{ if } \mu_{j-1,M} < y_{i,M}^* < \mu_{j,M}, \varepsilon_{i,M} \sim N[0,1], \\
 (\varepsilon_{i,1}, \dots, \varepsilon_{i,M}) &\sim N[\mathbf{0}, \mathbf{R}],
 \end{aligned}$$

where  $\mathbf{R}$  is the unrestricted correlation matrix of the random terms. In principle, this is a straightforward extension of the single variable model. The estimation is substantially complicated because of the amount of computation involved. In the one variable case, the probability is the area under the univariate normal density bounded by two points on a line, which requires two function evaluations of the univariate normal cdf. For two dimensions, the probability is the area under the bivariate normal surface bounded by a rectangle, which, in general, requires four function evaluations of the bivariate normal integral. For three dimensions, it requires eight function evaluations of the trivariate normal integral. And so on. The amount of computation rises with  $2^M$ . Moreover, the computation of the integrals, themselves, is cumbersome. For one dimension, the typical library routine computation of the normal integral involves evaluation of a ratio of two fourth or fifth order polynomials. The bivariate normal integral must typically be done using quadrature. [See, e.g., Drezner (1978).] For three dimensions or higher, the computation must be done by simulation, which will (with current technology) involve a formidable amount of computing. This model, even for only two dimensions, does not lend itself conveniently to the ordered logit form, and the received applications use the ordered probit model exclusively. [See, however, Dardanomi and Forcina (2004), who do obtain some analytical results for a multivariate ordered logit model.]

#### 7.3.1 Bivariate Ordered Probit Models

The two equation case has dominated the received applications, largely because of the practical difficulty of evaluating the higher order normal integrals needed to estimate the models.

For two outcomes, we have

$$\begin{aligned}
y_{i,1}^* &= \boldsymbol{\beta}_1' \mathbf{x}_{i,1} + \varepsilon_{i,1}, y_{i,1} = j \text{ if } \mu_{j-1} < y_{i,1}^* < \mu_j, j = 0, \dots, J_1, \\
y_{i,2}^* &= \boldsymbol{\beta}_2' \mathbf{x}_{i,2} + \varepsilon_{i,2}, y_{i,2} = j \text{ if } \delta_{j-1} < y_{i,2}^* < \delta_j, j = 0, \dots, J_2, \\
\begin{pmatrix} \varepsilon_{i,1} \\ \varepsilon_{i,2} \end{pmatrix} &\sim N \left[ \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \right].
\end{aligned}$$

The joint probability for  $y_{i,1} = j$  and  $y_{i,2} = k$  is

$$\begin{aligned}
\text{Prob}(y_{i,1} = j, y_{i,2} = k | \mathbf{x}_{i,1}, \mathbf{x}_{i,2}) &= \\
&\left[ \Phi_2[(\mu_j - \boldsymbol{\beta}_1' \mathbf{x}_{i,1}), (\delta_k - \boldsymbol{\beta}_2' \mathbf{x}_{i,2}), \rho] \right] - \left[ \Phi_2[(\mu_j - \boldsymbol{\beta}_1' \mathbf{x}_{i,1}), (\delta_{k-1} - \boldsymbol{\beta}_2' \mathbf{x}_{i,2}), \rho] \right] \\
&\left[ -\Phi_2[(\mu_{j-1} - \boldsymbol{\beta}_1' \mathbf{x}_{i,1}), (\delta_k - \boldsymbol{\beta}_2' \mathbf{x}_{i,2}), \rho] \right] - \left[ -\Phi_2[(\mu_{j-1} - \boldsymbol{\beta}_1' \mathbf{x}_{i,1}), (\delta_{k-1} - \boldsymbol{\beta}_2' \mathbf{x}_{i,2}), \rho] \right]
\end{aligned}$$

These are the probabilities that enter the log likelihood for a maximum likelihood estimator of the parameters.

Partial effects for this model will be complicated functions of the parameters regardless of how they are defined. But, for a bivariate model, such as this one, even what margin is of interest is not obvious. Derivatives of the bivariate probability,  $\text{Prob}(y_{i,1} = j, y_{i,2} = k | \mathbf{x}_{i,1}, \mathbf{x}_{i,2})$  might well not correspond to a useful experiment. One might, instead, wish to compute the derivatives of the conditional probability,

$$\text{Prob}(y_{i,1} = j | y_{i,2} = k, \mathbf{x}_{i,1}, \mathbf{x}_{i,2}) = \frac{\text{Prob}(y_{i,1} = j, y_{i,2} = k | \mathbf{x}_{i,1}, \mathbf{x}_{i,2})}{\text{Prob}(y_{i,2} = k | \mathbf{x}_{i,2})}$$

The denominator would be computed using the marginal, univariate ordered probit model. In either case, the computation will be based on a common result. For convenience, we drop the observation subscript and define the variables,

$$A_L = \mu_{j-1} - \boldsymbol{\beta}_1' \mathbf{x}_1, A_U = \mu_j - \boldsymbol{\beta}_1' \mathbf{x}_1, B_L = \delta_{k-1} - \boldsymbol{\beta}_2' \mathbf{x}_2, B_U = \delta_k - \boldsymbol{\beta}_2' \mathbf{x}_2$$

where subscripts “L” and “U” refer to “lower” and “upper,” respectively. Then, the bivariate probability is

$$\text{Prob}(y_{i,1} = j, y_{i,2} = k | \mathbf{x}_{i,1}, \mathbf{x}_{i,2}) = \left[ \Phi_2[A_U, B_U, \rho] - \Phi_2[A_L, B_U, \rho] \right] - \left[ \Phi_2[A_U, B_L, \rho] - \Phi_2[A_L, B_L, \rho] \right]$$

and the marginal univariate probability is

$$\text{Prob}(y_2 = k) = \Phi(B_U) - \Phi(B_L).$$

Computing partial effects from either viewpoint will require the result

$$\frac{\partial \Phi_2(A, B, \rho)}{\partial A} = \phi(A) \Phi \left( \frac{B - \rho A}{\sqrt{1 - \rho^2}} \right).$$

(The result is symmetric in  $A$  and  $B$ .) Collecting results, then

$$\frac{\partial \text{Prob}(y_1 = j, y_2 = k | \mathbf{x}_1, \mathbf{x}_2)}{\partial \mathbf{x}_1} = \begin{pmatrix} \phi(A_U)\Phi\left(\frac{B_U - \rho A_U}{\sqrt{1 - \rho^2}}\right) - \phi(A_L)\Phi\left(\frac{B_U - \rho A_L}{\sqrt{1 - \rho^2}}\right) \\ -\phi(A_U)\Phi\left(\frac{B_L - \rho A_U}{\sqrt{1 - \rho^2}}\right) + \phi(A_L)\Phi\left(\frac{B_L - \rho A_L}{\sqrt{1 - \rho^2}}\right) \end{pmatrix} (-\beta_1)$$

$$\frac{\partial \text{Prob}(y_1 = j, y_2 = k | \mathbf{x}_1, \mathbf{x}_2)}{\partial \mathbf{x}_2} = \begin{pmatrix} \phi(B_U)\Phi\left(\frac{A_U - \rho B_U}{\sqrt{1 - \rho^2}}\right) - \phi(B_L)\Phi\left(\frac{A_U - \rho B_L}{\sqrt{1 - \rho^2}}\right) \\ -\phi(B_U)\Phi\left(\frac{A_L - \rho B_U}{\sqrt{1 - \rho^2}}\right) + \phi(B_L)\Phi\left(\frac{A_L - \rho B_L}{\sqrt{1 - \rho^2}}\right) \end{pmatrix} (-\beta_2)$$

If any variables appear in both equations, the effects are added. For the conditional probabilities,

$$\frac{\partial \text{Prob}(y_1 = j, y_2 = k | \mathbf{x}_1, \mathbf{x}_2) / \text{Prob}(y_2 = k | \mathbf{x}_2)}{\partial \mathbf{x}_1} = \frac{\partial \text{Prob}(y_1 = j, y_2 = k | \mathbf{x}_1, \mathbf{x}_2) / \partial \mathbf{x}_1}{\text{Prob}(y_2 = k | \mathbf{x}_2)}$$

$$\frac{\partial \text{Prob}(y_1 = j, y_2 = k | \mathbf{x}_1, \mathbf{x}_2) / \text{Prob}(y_2 = k | \mathbf{x}_2)}{\partial \mathbf{x}_2} = \frac{\partial \text{Prob}(y_1 = j, y_2 = k | \mathbf{x}_1, \mathbf{x}_2) / \partial \mathbf{x}_2}{\text{Prob}(y_2 = k | \mathbf{x}_2)} - \text{Prob}(y_1 = j | y_2 = k, \mathbf{x}_1, \mathbf{x}_2) \frac{\phi(B_U) - \phi(B_L)}{\text{Prob}(y_2 = k | \mathbf{x}_2)} (-\beta_2)$$

As before, if variables appear in both equations, the two components are added. Before examining the applications of the model in detail, it is useful to look more closely at some special cases.

An admittedly trivial extension is the bivariate model in which  $\rho$  equals zero. In this instance, the bivariate model becomes a pair of univariate models. We mention this case at this point, as chronologically, the second application of the bivariate ordered probit model, Gustaffson and Stafford (1992), used this model to study child care subsidies and labor supply behavior for a sample of Swedish mothers. The hypothesis of uncorrelated equations is easily testable in this setting using either a likelihood ratio test or the Wald statistic ( $t$  ratio) associated with the estimate of  $\rho$ . Butler and Chatterjee (1995) consider other tests of the model specification, normality and exogeneity of the right hand sides, using GMM rather than maximum likelihood estimation. (They apply their methods to the study of dogs/television ownership noted below.)

### 7.3.2 Polychoric Correlation

The *polychoric correlation coefficient* is computed for a pair of discrete ordered variables, such as  $y_{i,1}$  and  $y_{i,2}$  above. The theory behind the computation is that  $y_{i,1}$  and  $y_{i,2}$  are censored versions of underlying, bivariate normally distributed variables, again, precisely as  $y_{i,1}$  and  $y_{i,2}$  above are obtained. The polychoric correlation coefficient is an estimator of the correlation coefficient in the underlying bivariate normal distribution. The best known method of computing the coefficient for grouped data (in the form of contingency tables), is due to Olssen (1979, 1980). [See, also, Ronning (1990) and Ronning and Kukuk (1996).] The development above suggests a counterpart for how to compute the coefficient when the data are individually

measured. If the two equations in the bivariate model have only their constant terms, and no regressors, then precisely the suggested underlying model emerges.

$$\begin{aligned}
 y_{i,1}^* &= \beta_1 + \varepsilon_{i,1}, y_{i,1} = j \text{ if } \mu_{j-1} < y_{i,1}^* < \mu_j, j = 0, \dots, J_1, \\
 y_{i,2}^* &= \beta_2 + \varepsilon_{i,2}, y_{i,2} = j \text{ if } \delta_{j-1} < y_{i,2}^* < \delta_j, j = 0, \dots, J_2, \\
 \begin{pmatrix} \varepsilon_{i,1} \\ \varepsilon_{i,2} \end{pmatrix} &\sim N \left[ \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \right].
 \end{aligned}$$

Thus, the implied algorithm, which has been built into modern software such as *NLOGIT*, *Stata* and *SAS*, is simply to fit a bivariate ordered probit model which has only constant terms in the two equations. [See, as well, Calhoun (1986, 1995) for further discussion of computer programs.] Returning to the regression model, it follows that the correlation coefficient in the bivariate ordered probit regression model can be interpreted as the conditional (on  $\mathbf{x}_{i,1}$  and  $\mathbf{x}_{i,2}$ ) polychoric correlation coefficient.

### 7.3.3 Semi-Ordered Bivariate Probit Model

A second interesting special case arises if one of the variables is binary;

$$\begin{aligned}
 y_{i,1}^* &= \beta_1' \mathbf{x}_{i,1} + \varepsilon_{i,1}, y_{i,1} = 0 \text{ if } y_{i,1}^* < 0, \text{ and } y_{i,1} = 1 \text{ if } y_{i,1}^* > 0, \\
 y_{i,2}^* &= \beta_2' \mathbf{x}_{i,2} + \varepsilon_{i,2}, y_{i,2} = j \text{ if } \delta_{j-1} < y_{i,2}^* < \delta_j, j = 0, \dots, J_2, \\
 \begin{pmatrix} \varepsilon_{i,1} \\ \varepsilon_{i,2} \end{pmatrix} &\sim N \left[ \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \right].
 \end{aligned}$$

This case (like the previous one) does not mandate any special modification of the likelihood function. The appropriate terms can be obtained directly from the earlier general result. This particular form has appeared in a number of applications, under the name “Bivariate Semi-Ordered Probit Model.” Weiss (1993) used this model to examine the extent of injuries in motorcycle injuries, with the binary variable is helmet use. Armstrong and McVicar (2000) used this form to examine the relationship between education and vocational training for a sample of Irish youth. McVicar and McKee (2002), using the same model, studied the two variables, vocational attainment (ordered) and working part time during education (binary), also for a sample of Irish youth. In this study, the education achievement is a four level exam measure.

### 7.3.4 Applications of the Bivariate Ordered Probit Model

The first application of the bivariate ordered probit model is Calhoun (1991, 1994) who examined the joint distribution of “Desired Family Size” (DFS) and “Children Ever Born” (CEB). In a followup analyze, he used CEB to truncate DFS, to eliminate unwanted children, then reexamined the model with this form of truncation. In an application to descriptions of criminal behavior and subsequent labor market experience, Nagin and Waldfogel (1995) examined the job market performance of young British offenders at ages 17 and 19. In a related analysis, Paternoster and Brame (1998) examined “self control” and “criminal behavior” in a study in criminology [See, also, comments in Britt (2000).] Butler and Chatterjee (1997), in their contribution to pet econometrics, analyzed the joint ownership of dogs and televisions. This is one of several studies in which authors used the bivariate ordered probit to model variables that

arguably should be analyzed as counts (with something like a Poisson regression model. However, the bivariate Poisson regression model remains to be well developed. [See, also, Sanko et al. (2004) who looked at ownership of cars and motorcycles.] The ordered probit model has been modified for use in contingent valuation studies, in which survey respondents express their preferences with a range of values rather than a point. Kuriama et al. (1998) used a contingent valuation study to examine consumers' preferences for a world heritage site in Japan. The ordered probit study follows a Vote/No Vote choice, and so has elements of the semiordered bivariate probit model described earlier as well. In two very natural application, Kohler and Rodgers (1999) studied the motivation to have children in a survey of pairs of twins. Christensen et al. (2003) also examined twins, in their case, seeking a genetic effect on fertility. Biswas and Das (2002) examined an epidemiologic study of diabetic retinopathy. Separate equations are specified for the right and left eye severity of the disease (coded 0 to 4). This is one of only a few Bayesian applications. [Biswas and Das benchmarked their study against an earlier analysis of the same data by Kim (1995). It is surprising that they did not use Kim's estimates in their priors. This seems like a natural application of Bayesian updating.] A variety of other applications have appeared, most since 2000, in economics, finance and transportation research. Table 32 lists some of the recent applications. (Full citations appear in the references list.)

**Table 32 Applications of Bivariate Ordered Probit Since 2000**

Year	Authors	Application
2000	Magee, et al.	Correlation between husband's and wife's education:
2002	Lawrence and Palmer	Views on health care reform,
2004	Bedi and Tunali	Participation in land and labor contracts in turkish agriculture
2004	Dupor et al.	Federal Reserve Open Market Committee: Bias announcement (ease, neutral, tighten) and magnitude of next meeting adjustment (-25, 25/0, 0, 0/25, 25+)
2005	Dueker et al.	Job restrictions of nurses:
2005	Filer and Honig	Pensions and retirement behavior,
2006	Adams	University and internal cost allocations of R&D expenditure
2006	Scott and Axhausen	Interactions between cars and season tickets,
2006	Scotti	Bivariate Model of Fed and European Central Bank main policy rates
2007	Mitchell and Weale	Accuracy of expectations about financial circumstances in the British Household Panel Survey

### 7.3.5 A Panel Data Version of the Bivariate Ordered Probit Model

Since it is a two equation model, it is unclear how common heterogeneity effects should enter the bivariate model. [See, e.g., Verbeek (1990), Verbeek and Nijman (1992) and Zabel (1992) for a similar exchange in the context of the sample selection model.] Generically, a bivariate model with time invariant random effects might appear

$$y_{it,1}^* = \beta_1' \mathbf{x}_{it,1} + \varepsilon_{it,1} + u_{1,i}; \quad y_{it,1} = j \text{ if } \mu_{j-1} < y_{it,1}^* < \mu_j, j = 0, \dots, J_1,$$

$$y_{it,2}^* = \beta_2' \mathbf{x}_{it,2} + \varepsilon_{it,2} + u_{2,i}; \quad y_{it,2} = j \text{ if } \delta_{j-1} < y_{it,2}^* < \delta_j, j = 0, \dots, J_2,$$

$$\begin{pmatrix} \varepsilon_{it,1} \\ \varepsilon_{it,2} \end{pmatrix} \sim N \left[ \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \right]; \quad \begin{pmatrix} u_{i,1} \\ u_{i,2} \end{pmatrix} \sim N \left[ \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{pmatrix} \right]$$

Computation of the parameters in this model would involve integration over both bivariate normal integrals. The approach used by Riphahn, Wambach and Million (2003) for a bivariate Poisson model with two random effects suggests an approach. Conditioned on the random effects, the likelihood function is

$$L | \mathbf{u}_1, \mathbf{u}_2 = \prod_{i=1}^N \prod_{t=1}^{T_i} \sum_{j=0}^{J_1} \sum_{k=0}^{J_2} m_{it,j} n_{it,k} \left\{ \begin{array}{l} \left[ \begin{array}{l} \Phi_2[(\mu_j - \beta'_1 \mathbf{x}_{it,1} - u_{i1}), (\delta_k - \beta'_2 \mathbf{x}_{it,2} - u_{i2}), \rho] \\ -\Phi_2[(\mu_{j-1} - \beta'_1 \mathbf{x}_{it,1} - u_{i1}), (\delta_k - \beta'_2 \mathbf{x}_{it,2} - u_{i2}), \rho] \end{array} \right] \\ \left[ \begin{array}{l} \Phi_2[(\mu_j - \beta'_1 \mathbf{x}_{it,1} - u_{i1}), (\delta_{k-1} - \beta'_2 \mathbf{x}_{it,2} - u_{i2}), \rho] \\ -\Phi_2[(\mu_{j-1} - \beta'_1 \mathbf{x}_{it,1} - u_{i1}), (\delta_{k-1} - \beta'_2 \mathbf{x}_{it,2} - u_{i2}), \rho] \end{array} \right] \end{array} \right\}$$

where  $m_{it,j} = 1$  if  $y_{it,1} = j$  and 0 otherwise and  $n_{it,k} = 1$  if  $y_{it,2} = k$  and 0 otherwise. To obtain a form of the likelihood function we can use for estimation, it is necessary to eliminate the unobserved random effects. We use a Cholesky decomposition of the covariance matrix to write

$$\begin{pmatrix} u_{i,1} \\ u_{i,2} \end{pmatrix} = \begin{bmatrix} \gamma_{11} & 0 \\ \gamma_{21} & \gamma_{22} \end{bmatrix} \begin{pmatrix} v_{i,1} \\ v_{i,2} \end{pmatrix}$$

where  $(v_{i,1}, v_{i,2})$  are independent  $N(0,1)$  variables. It follows that  $\gamma_{11}^2 = \sigma_1^2$ ,  $\gamma_{21}^2 + \gamma_{22}^2 = \sigma_2^2$  and  $\gamma_{21}\gamma_{22} = \sigma_{12}$ . The specific probabilities above with this substitution become

$$\text{Prob}(y_{it,1} = j, y_{it,2} = k | u_{i1}, u_{i2}, \mathbf{x}_{it,1}, \mathbf{x}_{it,2}) = \left\{ \begin{array}{l} \left[ \begin{array}{l} \Phi_2[(\mu_j - \beta'_1 \mathbf{x}_{it,1} - \gamma_{11}v_{i1}), (\delta_k - \beta'_2 \mathbf{x}_{it,2} - \gamma_{21}v_{i1} - \gamma_{22}v_{i2}), \rho] \\ -\Phi_2[(\mu_{j-1} - \beta'_1 \mathbf{x}_{it,1} - \gamma_{11}v_{i1}), (\delta_k - \beta'_2 \mathbf{x}_{it,2} - \gamma_{21}v_{i1} - \gamma_{22}v_{i2}), \rho] \end{array} \right] \\ \left[ \begin{array}{l} \Phi_2[(\mu_j - \beta'_1 \mathbf{x}_{it,1} - \gamma_{11}v_{i1}), (\delta_{k-1} - \beta'_2 \mathbf{x}_{it,2} - \gamma_{21}v_{i1} - \gamma_{22}v_{i2}), \rho] \\ -\Phi_2[(\mu_{j-1} - \beta'_1 \mathbf{x}_{it,1} - \gamma_{11}v_{i1}), (\delta_{k-1} - \beta'_2 \mathbf{x}_{it,2} - \gamma_{21}v_{i1} - \gamma_{22}v_{i2}), \rho] \end{array} \right] \end{array} \right\}$$

The unconditional log likelihood is obtained by integrating out the random effects. This step has been simplified by the Cholesky decomposition, since the bivariate integration involves independent standard normals. This could be done using nested Hermite quadratures or simulation. The latter is likely to be simpler and faster. The simulated log likelihood function is

$$\log L_S = \sum_{i=1}^N \log \frac{1}{R} \sum_{r=1}^R \prod_{t=1}^{T_i} \sum_{j=0}^{J_1} \sum_{k=0}^{J_2} m_{it,j} n_{it,k} \left\{ \begin{array}{l} \left[ \begin{array}{l} \Phi_2[(\mu_j - \beta'_1 \mathbf{x}_{it,1} - \gamma_{11}v_{i1,r}), (\delta_k - \beta'_2 \mathbf{x}_{it,2} - \gamma_{21}v_{i1,r} - \gamma_{22}v_{i2,r}), \rho] \\ -\Phi_2[(\mu_{j-1} - \beta'_1 \mathbf{x}_{it,1} - \gamma_{11}v_{i1,r}), (\delta_k - \beta'_2 \mathbf{x}_{it,2} - \gamma_{21}v_{i1,r} - \gamma_{22}v_{i2,r}), \rho] \end{array} \right] \\ \left[ \begin{array}{l} \Phi_2[(\mu_j - \beta'_1 \mathbf{x}_{it,1} - \gamma_{11}v_{i1,r}), (\delta_{k-1} - \beta'_2 \mathbf{x}_{it,2} - \gamma_{21}v_{i1,r} - \gamma_{22}v_{i2,r}), \rho] \\ -\Phi_2[(\mu_{j-1} - \beta'_1 \mathbf{x}_{it,1} - \gamma_{11}v_{i1,r}), (\delta_{k-1} - \beta'_2 \mathbf{x}_{it,2} - \gamma_{21}v_{i1,r} - \gamma_{22}v_{i2,r}), \rho] \end{array} \right] \end{array} \right\}$$

A fixed effects model might be considered as an alternative, however this would have several drawbacks: Two full sets of effects must be estimated. As usual, the fixed effects preclude time invariant variables in either equation. Though it remains to be established, it seems likely that the force of the incidental parameters problem (small  $T$  bias) would operate here as well. The

Mundlak (1978) device of including the group means of the time varying variables in the equations might be a useful middle ground.

### 7.3.6 Trivariate Ordered Probit Model

As noted earlier, for practical reasons, the bivariate probit is more or less the dimensional limit of the applications of the multivariate ordered probit model. Nonetheless, there have been a handful of applications of the trivariate probit model. Two in the area of transportation research that focus on joint determination of activity and travel model are Scott and Kanaroglou (2001) and Buliung (2005). Genius, Pantzios and Tzouvelakis (2005) estimate a “trivariate semi-ordered probit model.” In their application to organic farming in Greece, two of the three equations, contact with an extension agent and use of other sources of information, are binary, while the land adoption decision (none, part, full) has three outcomes. Crouchley (2005) is a methodology study of statistical modeling.

### 7.3.7 Models of Sample Selection with an Ordered Probit Selection Rule

The models of sample selectivity in this area are built as extensions of Heckman’s (1979) canonical model,

*Probit Participation Equation*

$$z_i^* = \alpha'w_i + u_i$$

$$z_i = 1[z_i^* > 0]$$

*Regression Activity Equation*

$$y_i^* = \beta'x_i + \varepsilon_i$$

$$(\varepsilon_i, u_i) \sim N[(0,0), (1, \rho\sigma_\varepsilon, 1)]$$

*Observation:* For observations with  $z_i = 1$ ,

$$\begin{aligned} E[y_i^* | x_i, w_i, z_i = 1] &= \beta'x_i + (\rho\sigma_\varepsilon)[\phi(\alpha'w_i)/\Phi(\alpha'w_i)] \\ &= \beta'x_i + \theta\lambda_i. \end{aligned}$$

Estimation of the regression equation by least squares while ignoring the selection issue produces biased and inconsistent estimators of all the model parameters. Estimation of this model by two step methods is documented in a voluminous literature, including Heckman (1979) and Greene (2008a). The two step method involves estimating  $\alpha$  first in the participation equation using an ordinary probit model, then computing an estimate of  $\lambda_i$ ,  $\hat{\lambda}_i = \phi(\hat{\beta}'x_i)/\Phi(\hat{\beta}'x_i)$ , for each individual in the selected sample. At the second step, an estimate of  $(\beta, \theta)$  is obtained by linear regression of  $y_i$  on  $x_i$  and  $\hat{\lambda}_i$ . Necessary corections to the estimated standard errors are described in Heckman (1979), Greene (1981, 2008b), and, in general terms, in Murphy and Topel (2002). As noted earlier, the binary probit model is a special case of the ordered probit model. The extension of the sample selection model would follow from replacing the participation equation with

*Ordered Probit Participation Equation*

$$z_i^* = \alpha'w_i + u_i$$

$$z_i = j \text{ if } \mu_{j-1} < z_i^* \leq \mu_j.$$

Then, the objective is to recast the conditional mean function,  $E[y_i^* | \mathbf{x}_i, \mathbf{w}_i, z_i = j]$  and determine an appropriate estimator and set of inference procedures. A typical application (several of those listed below) considers an “Educational Attainment” participation equation (*secondary, college, graduate*) and an outcome equation such as an earnings equation.

Garen (1984) builds directly on the Heckman model. He departs from a model in which

$$\begin{aligned} y_i | \mathbf{x}_i, z_i = 0 &= \boldsymbol{\beta}_0' \mathbf{x}_i + \varepsilon_{i0} \\ y_i | \mathbf{x}_i, z_i = 1 &= \boldsymbol{\beta}_1' \mathbf{x}_i + \varepsilon_{i1} \\ z_i^* &= \boldsymbol{\pi}_1' \mathbf{x}_i + \boldsymbol{\pi}_2' \mathbf{w}_i + u_i, \quad z_i = 1[z_i^* > 0], \end{aligned}$$

which is similar to the selection model shown above. [As stated, it is a “mover/stayer model.” See, e.g., Nakosteen and Zimmer (1980) and Greene (2008a, p. 888).] Garen’s suggestion from here suggests how to proceed if  $z_i$  is continuous – i.e., if  $z_i^*$  were the observation. He proposes to treat  $z_i$  as if it were observed in the form of integer values,  $1, \dots, n$ , noting that the continuous variable emerges as  $n \rightarrow \infty$ . There is, then a different regression equation for each value of  $z_i$ . What follows is an analysis of a transformed regression equation that is augmented with powers of  $z_i$  and products of  $z_i$  and  $\mathbf{x}_i$ . While not a sample selection treatment as such, this does point in the direction of a formal sample selection treatment based on the ordered probit model.

Terza (1987) develops the two step estimator for a regression model in which one of the regressors is generated by an ordered ordered probit model without regressors. The structural equations are equivalent to

$$\begin{aligned} y_i &= \boldsymbol{\beta}' \mathbf{x}_i + \theta q_i + \varepsilon_i \\ q_i^* &= \alpha + u_i \\ q_i &= j \text{ if } \mu_{j-1} < q_i^* \leq \mu_j. \\ (\varepsilon_i, u_i) &\sim N[(0, 0), (\sigma_\varepsilon^2, \rho\sigma_\varepsilon, 1)] \end{aligned}$$

It is convenient to define (once again)

$$m_{ij} = 1 \text{ if } q_i = j \text{ and } m_{ij} = 0 \text{ otherwise.}$$

Under these assumptions, Terza’s main result is

$$E[y_i | \mathbf{x}_i, m_{i0}, m_{i1}, \dots, m_{iJ}] = \boldsymbol{\beta}' \mathbf{x}_i + (\theta\rho) f_i$$

where

$$f_i = \sum_{j=0}^J m_{ij} \left( \frac{\phi(\mu_{j-1} - \alpha) - \phi(\mu_j - \alpha)}{\Phi(\mu_j - \alpha) - \Phi(\mu_{j-1} - \alpha)} \right).$$

[A similar result for the conditional mean of a doubly truncated variable appears in Maddala (1983, p. 366).] Terza goes on to propose a two step estimation procedure. The first step involves maximum likelihood estimation of  $(\alpha, \mu_{-1}, \mu_0, \mu_1, \dots, \mu_J)$ . This can be done (first noting that as usual,  $\mu_{-1} = -\infty$ ,  $\mu_1 = 0$  and  $\mu_J = \infty$ ) using only the sample proportions in the  $J+1$  cells. The model for  $q_i$  implies  $\text{Prob}(q_i > 0) = \Phi(\alpha)$ , so the estimator of  $\alpha$  is  $\Phi^{-1}(1-P_0)$ . Continuing,  $\text{Prob}(q_i$

$> 1) = \Phi(\mu_1 - \alpha)$  which suggests a method of moments estimator of  $\mu_1$  based on  $P_1$ , and so on. With these estimates in hand, he then proposes linear regression of  $\mathbf{y}$  on  $\mathbf{X}$  and  $\hat{\mathbf{f}}$  to estimate  $\boldsymbol{\beta}$  and  $(\theta\rho)$ . (A method of computing appropriate standard errors is presented later.) The use of the constructed regressor is a means to another end, consistent estimation of  $\boldsymbol{\beta}$ .

As Terza (1987) notes (p. 278) his model is not a correction for selection because the values of the dependent variable are observed for all observations. On the other hand, by a minor rearrangement of terms, the results are precisely what is needed for a model of sample selection. First, while retaining the ordered probit observation mechanism for  $q_i$ , replace the constant  $\alpha$  with the mean of the latent regression,  $\boldsymbol{\alpha}'\mathbf{w}_i$ . Second, we note that in the “selection on  $j$ ” case, we observe not  $(m_{i0}, m_{i1}, \dots, m_{iJ})$  in full, but only one of them. Terza’s results then imply

$$E[y_i | \mathbf{x}_i, \mathbf{w}_i, q_i = j] = E[y_i | \mathbf{x}_i, m_{ij} = 1] = \boldsymbol{\beta}'\mathbf{x}_i + \gamma \left( \frac{\phi(\mu_{j-1} - \boldsymbol{\alpha}'\mathbf{w}_i) - \phi(\mu_j - \boldsymbol{\alpha}'\mathbf{w}_i)}{\Phi(\mu_j - \boldsymbol{\alpha}'\mathbf{w}_i) - \Phi(\mu_{j-1} - \boldsymbol{\alpha}'\mathbf{w}_i)} \right).$$

This is the result needed to complete the sample selection model. The same two step method can now be applied. Terza’s method of computing corrected asymptotic standard errors is essentially unchanged.

Jimenez and Kugler (1987) appears to be the first formal application of the preceding sample selection model. The application is an earnings equation for the Bogota subsample of a 1979-1981 nationwide survey of graduates in Colombia. The selection mechanism is determined by participation in a vocational and technical training course (SENA), recorded as *none*, *short* or *long*. The authors derived the conditional mean function from first principles; the derivation follows naturally from earlier results in Maddala (1983), Garen (1984), Heckman (1979), Kenny et al. (1979), Lee and Trost (1978) and Trost and Lee (1978). Kao and Wu (1990) applied the same model to an analysis of bond yields in which the selection mechanism assigns bonds to risk classes by a rating agency. [See, as well, Acharya (1988) for a more elaborate development of the sample selection model.]

Frazis’s (1993) study is similar to Jimenez and Kugler. This study analyzes earnings of high school seniors from the National Longitudinal Study of the High School Class of 1972. A panel of seniors was interviewed in 1972, then again five times between 1973 and 1986. Frazis’s analysis departed from the basic framework in two ways. The earnings equation is

$$\log y = \boldsymbol{\beta}'\mathbf{x} + \sum_j \gamma_j S_j + \delta XS + \phi u + \lambda uS + \varepsilon$$

where  $y$  is earnings,  $\mathbf{x}$  is a vector of control variables,  $S_j$  is a set of dummy variables that equal one if the least level of schooling,  $j$ , is attained and zero otherwise,  $XS$  is interactions of the school attainment dummy variables with  $X$  and  $u$  represents “aspects of the ability to acquire human capital that are unobservable to the researcher.” Thus, since schooling level is the ordered selection mechanism, as stated, this model resembles a treatment effects model, and is also similar to Terza’s (1987) formulation. (Motivation for the parts of the equation are given in the paper.) However, note once again, that the observation will be conditioned not on all  $S_j$ , but only on the one that corresponds to the individual’s schooling level. Estimates of  $E[u|S_j]$  to serve as the proxy for  $u$  in the earnings equation are obtained by estimating the ordered probit model for schooling level and computing the conditional mean function given earlier. The estimating equation (fit by ordinary least squares) is obtained by replacing  $u$  in the equation above (in both places) with

$$\hat{U}_j = \left( \frac{\phi(\mu_{j-1} - \hat{\alpha}'_{j-1} \mathbf{w}_i) - \phi(\mu_j - \hat{\alpha}'_j \mathbf{w}_i)}{\Phi(\mu_j - \hat{\alpha}'_j \mathbf{w}_i) - \Phi(\mu_{j-1} - \hat{\alpha}'_{j-1} \mathbf{w}_i)} \right).$$

The second noteworthy point is that, as the author mentions in passing, the ordered probit model provides separate regression coefficients for each level of education. As he notes, this allows negative probabilities. A discussion of aspects of the data set that should prevent this is given.

Two remaining studies of sample selection with ordered probit selection mechanisms are Amel and Liang (1994, 1997) and Butler et al. (1994, 1998). In the first of these, the authors examine firm performance in the banking industry. The conditioning equation used depends on the setting. It depends on the amount of entry in the market; the authors describe small markets in which entry is described with a simple probit model, and large ones in which ordered probit and truncated Poisson models are used. Butler, Finegan and Siegfried (1998) [see, also Butler et al. (1994)] analyzed performance in economics courses. The selection mechanism is calculus proficiency measured by level of training across several possible courses.

Li and Tobias (2006a) replicated Butler et al. (1998) using a Bayesian method rather than two step least squares. The authors describe an “augmented likelihood function” for the model. With noninformative priors, they “virtually identically” replicated the original results, which suggests that the augmented likelihood function is not equal to the one given above. Technical details are not provided in the paper, but are promised in a no longer existing Iowa State University Economics Department working paper. [Li and Tobias (2006c).] The working paper is reincarnated under the same title in Li and Tobias (2006b). There the authors note that the dependent variable in the regression is actually a grade level, which is also discrete and ordered. The model in (2006b) is a treatment effects model in a triangular system with the outcome of the first ordered probit regression, in the form of a set of endogenous dummy variables, appearing on the right hand side of a second ordered outcome model, the grade attainment. [Sajaia (2008) is vaguely related to this, however, his treatment of the recursive model builds a simultaneous equations system in the latent regression, which seems difficult to motivate. This paper merely documents a *Stata* program, and does not provide detailed technical background.] The Li and Tobias model without the dummy variables (i.e., under a restriction that their coefficients are zero) would be the bivariate ordered probit model of Section 7.3.1, so it appears that the authors have rediscovered the MLE for the bivariate model, using a Gibbs sampling and MCMC algorithm rather than classical maximum likelihood. Technical details are omitted from the (2006b) paper, so it is difficult to discern how closely the results resemble each other, but one would expect them, with noninformative priors, to give roughly the same numerical results.

Missing from the preceding and from the received literature is a maximum likelihood estimator for the ordered probit sample selection model. One reason one might wish to consider an MLE as an alternative approach is that the two step estimators do not produce an estimator of  $\rho$ , which is likely to be an interesting parameter, for example, if one wished to test for “selectivity.” [In the basic case, there is a method of moments estimator of  $\rho$  available – See Heckman (1979) and Greene (1981). However, none has been derived for the ordered choice case. An analog to the estimator developed by Heckman (1979) would be straightforward. However, it will have the same shortcoming as the one in the basic model. As shown in Greene (1981), the estimator is not bounded by -1 and +1. Moreover, even when it does fall in the right range, no inference is possible. (This latter point is of minor consequence. In the original model above, inference is possible about  $\theta = \rho\sigma$  based on the OLS results, and  $\rho = 0$  is both necessary and sufficient for  $\theta = 0$ , as  $\sigma$  cannot be zero in a sensible model.)

The log likelihood for the original sample selection model (binary selection and linear regression) is given in Greene (2008a, eq. 24-33) and in Econometric Software (2007);

$$\log L = \sum_{z_i=0} \log \Phi(-\boldsymbol{\alpha}'\mathbf{w}_i) + \sum_{z_i=1} \log \left[ \frac{1}{\sigma_\varepsilon} \phi\left(\frac{y_i - \boldsymbol{\beta}'\mathbf{x}_i}{\sigma_\varepsilon}\right) \Phi\left(\frac{\rho(y_i - \boldsymbol{\beta}'\mathbf{x}_i)/\sigma_\varepsilon + \boldsymbol{\alpha}'\mathbf{w}_i}{\sqrt{1-\rho^2}}\right) \right].$$

This estimator, though apparently much less frequently used than the two step method, is available as a preprogrammed procedure in contemporary software such as *Stata* and *NLOGIT*. Note that it is a full information maximum likelihood estimator for all the parameters in the model. The estimator is not less robust than the two step estimator; both are fully parametric based on the bivariate normal distribution.

The counterpart for an ordered probit sample selection model will replace the term  $\Phi(\cdot)$  in the square brackets with

$$F_i = \Phi\left(\frac{\rho(y_i - \boldsymbol{\beta}'\mathbf{x}_i)/\sigma_\varepsilon - (\mu_j - \boldsymbol{\alpha}'\mathbf{w}_i)}{\sqrt{1-\rho^2}}\right) - \Phi\left(\frac{\rho(y_i - \boldsymbol{\beta}'\mathbf{x}_i)/\sigma_\varepsilon - (\mu_{j-1} - \boldsymbol{\alpha}'\mathbf{w}_i)}{\sqrt{1-\rho^2}}\right)$$

and the term  $\Phi(-\boldsymbol{\alpha}'\mathbf{w}_i)$  with

$$\text{Prob}(z_i \neq j | \mathbf{w}_i) = 1 - [\Phi(\mu_j - \boldsymbol{\alpha}'\mathbf{w}_i) - \Phi(\mu_{j-1} - \boldsymbol{\alpha}'\mathbf{w}_i)].$$

As stated, this is a conventional maximum likelihood estimator that produces the familiar properties consistency, asymptotic normality, etc. If the selection is “selection on a particular  $j$ ,” however, then no more than one of the threshold parameters will be estimable. Assuming that  $\boldsymbol{\alpha}$  contains a constant term, if selection is on  $j = 0$ , then the second probability becomes zero and  $\mu_0$  already equals zero. If selection is on  $j = 1$ , then  $\mu_0$  in the second probability is zero and the constant in  $\boldsymbol{\alpha}$  is identified, while in the first probability,  $\mu_1$  is estimable distinct from the constant in  $\boldsymbol{\alpha}$ . If selection is on  $j > 1$ , then the two probabilities have separate constant terms, but only two distinct constant terms are estimable. The first constant term estimates  $(\alpha_0 - \mu_j)$  and the second estimates  $(\alpha_0 - \mu_{j-1})$ .

Full information maximum likelihood based on the probabilities shown above should be a conventional, relatively straightforward exercise. However, there is a simplification that might prove useful. This (and the original model) is an ideal setting to employ the Murphy and Topel (2002) kind of two step estimator. As already seen, we can estimate the ordered probit model in isolation, using maximum likelihood. Let

$$\hat{A}_j = \hat{\mu}_j - \hat{\boldsymbol{\alpha}}'z_i, \quad \hat{A}_{j-1} = \hat{\mu}_{j-1} - \hat{\boldsymbol{\alpha}}'z_i.$$

Then, a two step approach can be used in which the log likelihood function maximized at the second step is

$$\log L = \sum_{z_i \neq j} \log [1 - (\Phi(\hat{A}_j) - \Phi(\hat{A}_{j-1}))] + \sum_{z_i = j} \log \left[ \frac{1}{\sigma_\varepsilon} \phi\left(\frac{y_i - \boldsymbol{\beta}'\mathbf{x}_i}{\sigma_\varepsilon}\right) \left\{ \Phi\left(\frac{\rho(y_i - \boldsymbol{\beta}'\mathbf{x}_i)/\sigma_\varepsilon - \hat{A}_j}{\sqrt{1-\rho^2}}\right) - \Phi\left(\frac{\rho(y_i - \boldsymbol{\beta}'\mathbf{x}_i)/\sigma_\varepsilon - \hat{A}_{j-1}}{\sqrt{1-\rho^2}}\right) \right\} \right].$$

Note that the first term is now an irrelevant constant, and the log likelihood function to be maximized is based only on the selected sample. This can be made even more convenient by

reparameterizing it with the Olsen (1978) reparameterization,  $\theta = 1/\sigma_\varepsilon$  and  $\boldsymbol{\gamma} = (1/\sigma_\varepsilon)\boldsymbol{\beta}$ . Now, the relevant log likelihood is

$$\log L^* = \sum_{z_i=j} \log \left[ \theta \phi(\theta y_i - \boldsymbol{\gamma}'\mathbf{x}_i) \left\{ \Phi \left( \frac{\rho(\theta y_i - \boldsymbol{\gamma}'\mathbf{x}_i) - \hat{A}_j}{\sqrt{1-\rho^2}} \right) - \Phi \left( \frac{\rho(\theta y_i - \boldsymbol{\gamma}'\mathbf{x}_i) - \hat{A}_{j-1}}{\sqrt{1-\rho^2}} \right) \right\} \right].$$

Finally, let  $\tau = \rho/\sqrt{1-\rho^2}$ . Then, the log likelihood simplifies a bit more to

$$\log L^* = \sum_{z_i=j} \log \theta \phi(\theta y_i - \boldsymbol{\gamma}'\mathbf{x}_i) + \sum_{z_i=j} \log \left\{ \Phi \left( \tau(\theta y_i - \boldsymbol{\gamma}'\mathbf{x}_i) - \sqrt{1+\tau^2} \hat{A}_j \right) - \Phi \left( \tau(\theta y_i - \boldsymbol{\gamma}'\mathbf{x}_i) - \sqrt{1+\tau^2} \hat{A}_{j-1} \right) \right\}$$

Once estimates of  $\theta$ ,  $\boldsymbol{\gamma}$  and  $\tau$  are in hand, estimates of the structural parameters,  $\sigma_\varepsilon$ ,  $\boldsymbol{\beta}$  and  $\rho$ , can be obtained by inverting the transformations. This approach has an additional benefit in that the range of  $\tau$  is unrestricted, while that of  $\rho$  must be restricted to  $(-1,+1)$  during estimation.

### 7.3.8 A Sample Selected Ordered Probit Model

The second case we consider, also absent from the literature heretofore, reverses the role of the regression and the ordered probit model. As an example, we might consider a model of educational attainment or performance in a training or vocational education program (e.g., low, median, high), with selection into the program as an observation mechanism. [Boes (2007) examines a related case, that of a treatment,  $D$  that acts as an endogenous dummy variable in the ordered outcome model.] The structural equations would be

*Selection Equation*

$$z^* = \boldsymbol{\alpha}'\mathbf{w} + u$$

$$z = 1[z^* > 0]$$

*Ordered Probit Outcome*

$$y^* = \boldsymbol{\beta}'\mathbf{x} + \varepsilon$$

$$y = j \text{ if } \mu_{j-1} < y^* \leq \mu_j.$$

*Observation Mechanism*

$$y, \mathbf{x} \text{ observed when } z = 1.$$

$$(\varepsilon, u) \sim N[(0,0), (1,\rho,1)]$$

In this situation, the “second step” model is nonlinear. The received literature contains many applications in which authors have “corrected for selectivity” by following the logic of the Heckman two step estimator, that is, by constructing  $\lambda_i = \phi(\boldsymbol{\alpha}'\mathbf{w}_i)/\Phi(\boldsymbol{\alpha}'\mathbf{w}_i)$  from an estimate of the probit selection equation and adding it to the outcome equation. [See, e.g., Greene (1994). Several other examples are provided in Greene (2008b).] This is, however, only appropriate in the linear model with normally distributed disturbances. An explicit expression, which does not involve an inverse Mills ratio, is given for the case in which the unconditional regression is  $E[y|\mathbf{x},\varepsilon] = \exp(\boldsymbol{\beta}'\mathbf{x} + \varepsilon)$  is given in Terza (1998). A template for nonlinear single index function models subject to selectivity is developed in Terza (1998) and Greene (2006, 2008a, Sec. 24.5.7).

Applications specifically to the Poisson regression appear in several places, including Greene (1995, 2005). The general case typically involves estimation either using simulation or quadrature to eliminate an integral involving  $u$  in the conditional density for  $y$ . Cases in which both variables are discrete, however, are somewhat simpler. A near parallel to the model above is the bivariate probit model with selection developed by Boyes, Hoffman and Low (1989) in which the outcome equation above would be replaced with a second probit model. [Wynand and van Praag (1981) proposed the bivariate probit/selection model, but used the two step approach rather than maximum likelihood.] The log likelihood function for the bivariate probit model is given in Boyes et al. (1989) and Greene (2008a, p. 896):

$$\log L = \sum_{z=0} \log \Phi(-\alpha'w) + \sum_{z=1,y=0} \log \Phi_2(-\beta'x, \alpha'w, -\rho) + \sum_{z=1,y=1} \log \Phi_2(\beta'x, \alpha'w, \rho)$$

A straightforward extension of the result provides the log likelihood for the ordered probit case,

$$\log L = \sum_{z=0} \log \Phi(-\alpha'w) + \sum_{z=1} \sum_{j=0}^J m_{ij} \log [\Phi_2(\mu_j - \beta'x, \alpha'w, \rho) - \Phi_2(\mu_{j-1} - \beta'x, \alpha'w, \rho)]$$

where  $m_{ij} = 1$  if  $y_i = j$ .

Table 32 presents estimates of a sample selection model. We have used the choice of *PUBLIC* insurance as the selection mechanism. About 87% of the sample choose the public insurance. We speculate that the factors underlying the motivation to purchase the insurance are also related to the response of health satisfaction. The full model is

$$\begin{aligned} PUBLIC_i^* &= \alpha_1 + \alpha_2 AGE_i + \alpha_3 EDUC_i + \alpha_4 HANDDUM_i + u_i, \\ PUBLIC_i &= 1[PUBLIC_i^* > 0] \\ HEALTH_i^* &= \beta'x_i + \varepsilon_i \\ HEALTH_i &= j \text{ if } \mu_{j-1} < HEALTH_i^* \leq \mu_j \\ (HEALTH_i, x_i) &\text{ observed when } PUBLIC_i = 1, \\ (u_i, \varepsilon_i) &\sim N_2[(0,1), (1,1,\rho)], \end{aligned}$$

using the same set of regressors as previously. The estimate of  $\rho$  suggests that the conjecture might be correct. On average, the factors that motivate insurance purchase seem also to motivate a higher response to the health satisfaction question.

**Table 32 Estimated Ordered Probit Sample Selection Model**

Variable	Coefficient	Standard Error	b/St.Er.	P[ Z >z]	Mean of X
Binomial Probit Model					
Dependent variable	PUBLIC				
Number of observations	4483				
Log likelihood function	-1471.427				
Restricted log likelihood	-1711.545				
Results retained for SELECTION model.					
-----					
Index function for probability					
Constant	3.59248***	.16511284	21.758	.0000	
AGE	-.00271	.00243682	-1.110	.2670	43.440107
EDUC	-.19666***	.00935786	-21.016	.0000	11.418086
HANDDUM	.28812***	.09802085	2.939	.0033	.1119786
Note: ***, **, * = Significance at 1%, 5%, 10% level.					
-----					
Predictions for Binary Choice Model. Predicted value is 1 when probability is greater than .500000, 0 otherwise.					
Actual Value	Predicted Value		Total Actual		
	0	1			
0	164 ( 3.7%)	408 ( 9.1%)	572 ( 12.8%)		
1	141 ( 3.1%)	3770 ( 84.1%)	3911 ( 87.2%)		
Total	305 ( 6.8%)	4178 ( 93.2%)	4483 (100.0%)		
-----					
Ordered Probit Model with Selection.					
Dependent variable	HEALTH				
Log likelihood function	-6496.032				
-----					
Variable	Coefficient	Standard Error	b/St.Er.	P[ Z >z]	Mean of X
Index function for probability					
Constant	2.23467***	.12704133	17.590	.0000	
AGE	-.01597***	.00163299	-9.780	.0000	43.440107
EDUC	-.03143***	.00924985	-3.398	.0007	11.418086
INCOME	.23843**	.09938389	2.399	.0164	.3487401
MARRIED	-.00934	.03862829	-.242	.8089	.7521749
KIDS	.05435	.03707325	1.466	.1427	.3794334
Threshold parameters for index					
Mu(1)	.96947***	.03943988	24.581	.0000	
Mu(2)	2.23993***	.05243459	42.718	.0000	
Mu(3)	2.70909***	.05470858	49.519	.0000	
Selection equation					
Constant	3.45116***	.16227409	21.267	.0000	
AGE	-.00535**	.00245191	-2.181	.0292	43.440107
EDUC	-.18037***	.00930000	-19.394	.0000	11.418086
HANDDUM	.67100***	.08032991	8.353	.0000	.1119786
Cor[u(probit),e(ordered probit)]					
Rho(u,e)	.80801***	.04519162	17.880	.0000	
Note: ***, **, * = Significance at 1%, 5%, 10% level.					

### 7.3.9 An Ordered Probit Model with Endogenous Treatment Effects

Munkin and Trivedi (2008) have analyzed a model that bears some connection to the selection model proposed in the previous section. The model extension considered involves a set of endogenous “treatment dummy variables.” That is,

$$y_i^* = \boldsymbol{\beta}'\mathbf{x}_i + \boldsymbol{\delta}'\mathbf{d}_i + \varepsilon_i$$
$$y_i = j \text{ if } \mu_{j-1} < y_i^* \leq \mu_j.$$

where  $y_i$  is a measure of medical service utilization (actually a count with excess zeros – the ordered choice model is used as an approximation). The additional vector of covariates,  $\mathbf{d}_i$ , is a set of dummy variables that is the outcome of a choice of treatments; one of  $M$  treatments is chosen and for that choice,  $d_{im} = 1$  and  $d_{im'} = 0$  for all others. (We have included all  $M$  treatments in  $\mathbf{d}_i$  for pedagogical convenience. In their analysis, one of the dummy variables is immediately dropped from the model since only  $M-1$  are needed to determine the observed outcome.) The treatment outcome is determined by a multinomial probit model of underlying utility across the choices. [See Train (2003) and the large number of sources cited by Munkin and Trivedi for discussion of the multinomial probit model.] The endogeneity of the treatment effects follows from the correlations between the random elements of the random utility equations in the choice model and the random term,  $\varepsilon_i$  in the ordered choice model. A Bayesian (MCMC) treatment is used to estimate the posterior means of the parameters

## 8 Semiparametric Estimators and Analyses

The foregoing has surveyed nearly all of the literature on ordered choice modeling. We have, of course, listed only a small fraction of the received applications. But, the full range of methodological developments has been presented, with a single remaining exception. As in many other areas of econometrics, a thread of the contemporary literature has explored the boundaries of the model that are circumscribed by the distributional assumptions. We have limited ourselves to ordered logit and probit models, while relaxing certain assumptions such as homoscedasticity within the boundaries of the parametric model. The last strand of literature to be examined is the development of estimators that extend beyond the parametric distributional assumptions. As a general proposition, it is useful to organize the overview around a few features of the model, scaling, the distribution of the disturbance, the functional form of the regression, and so on. In each of these cases, we can focus on applications that broaden the reach of the ordered choice model to less tightly specified settings.

There is a long, rich history of semiparametric and nonparametric analysis of binary choice modeling (far too long and rich to examine in depth in this already long survey) that begins in the 1970s, only a few years after analysis of individual binary data became a standard technique. The binary choice literature has two focal points, maximum score estimation [Manski (1975, 1985), Manski and Thompson (1985) and Horowitz (1992)] and the Klein and Spady (1993) kernel based semiparametric estimator for binary choice. (As noted, there is a huge number of other papers on the subject. We are making no attempt to survey this literature.) Some of the more recent developments build on these two (mainly on the second; MSCORE remains to provide a platform for analysis of ordered choices). Surprisingly, the formal extension of the binary choice models to what would seem to be the natural next step, ordered choice, takes place entirely since 2000.

To a very small extent, some of the developments already mentioned move the analysis in the direction of a semiparametric approach. Agresti (1999), for example, notes the extension of GEE methods [see Diggle, Liang and Zeger (1994)] to the ordered choice model. GEE modeling is based more strongly on conditional means and variances than on distributions, and can be viewed as a small step away from the maximum likelihood estimator. (The step is quite small; the formal distributional model is still assumed. One might surmise, however, that the GEE estimator has at least the potential to be robust to failures of the distributional assumption. This remains to be verified, however.) On the other hand, if the latent class model (LCM) that we examined in Section 5.2.6 is simply interpreted as a mixing model rather than as a latent grouping model, then the LCM certainly qualifies as a semiparametric approach. [See Heckman and Singer (1984) for example.] Likewise, the mixed ordered probit (random parameters) model can also be viewed as a samiparametric estimator; a continuous mixture of underlying distributions that does not adhere to a strict distributional assumption. [See, e.g., McFadden and Train (2001) for discussion of using continuous mixture models to approximate any underlying distribution.] (For the ordered choice model, to achieve full generality in this interpretation, we would want to allow the thresholds, as well as the regression slopes, to be random.)

The received literature on semiparametric (and semi-nonparametric and nonparametric) analysis of ordered choice models is fairly compact. We begin with a study by Chen and Khan (2003) that considers the ordered probit model in the presence of unknown (and not parameterized) heteroscedasticity. Lewbel (2000) goes a step beyond Chen and Khan in allowing the distribution to be unspecified as well. We will then examine Stewart's (2003) parameterized model that approximates an unknown distribution. Some general observations are collected in Section 8.5. This is not a complete enumeration of this thred of literature (though it is fairly close). Three studies not examined in detail below, but mentioned here are Coppejans (2007) and Klein and Sherman (2002), both of which develop consistent parameters, but are, at the same

time, focused somewhat more heavily on methodological aspects of estimation than the papers examined below.

## 8.1 Heteroscedasticity.

Chen and Khan (2003) propose a semiparametric estimator for the heteroscedastic *ordered probit* model,

$$\begin{aligned} y_i^* &= \alpha + \boldsymbol{\beta}'\mathbf{x}_i + \sigma(\mathbf{x}_i)\varepsilon_i \\ y_i &= j \text{ if } \mu_{j-1} < y_i^* \leq \mu_j, j = 0, 1, \dots, J. \end{aligned}$$

(we are adapting their application to our notation – theirs differs in several ways likely to produce ambiguities in the presentation). The issue is whether it is possible efficiently (by semiparametric standards) to estimate  $\boldsymbol{\beta}$ . Several normalizations are necessary to begin. As usual,  $\mu_{-1} = -\infty$  and  $\mu_J = +\infty$ . Since there is assumed to be a nonzero constant term,  $\alpha$ ,  $\mu_0 = 0$ . They restrict attention to the case  $J = 2$  (three possible outcomes). “As is always the case with discrete response models, location and scale normalizations are required. As a location normalization, to identify the intercept term,  $[\alpha]$  we set  $[\mu_0] = 0$ . As a scale normalization, we set  $[\mu_1] = 1$ .” (Again, our notation.) The last assumption is, of course, crucial. Heretofore, we have achieved scale normalization by assuming  $\sigma_\varepsilon = 1$ . The implication of the new assumption in the three outcome is as follows:

$$\begin{aligned} \text{Prob}(y_i = 0 \mid \mathbf{x}_i) &= \Phi\left(\frac{-\alpha - \boldsymbol{\beta}'\mathbf{x}_i}{\sigma(\mathbf{x}_i)}\right) = P_{0i} \\ 1 - \text{Prob}(y_i = 2 \mid \mathbf{x}_i) &= \Phi\left(\frac{1 - \alpha - \boldsymbol{\beta}'\mathbf{x}_i}{\sigma(\mathbf{x}_i)}\right) = 1 - P_{2i} \end{aligned}$$

It follows that

$$\Phi^{-1}(1 - P_{2i}) - \Phi^{-1}(P_{0i}) = \Phi^{-1}(P_{0i} + P_{1i}) - \Phi^{-1}(P_{0i}) = \frac{1}{\sigma(\mathbf{x}_i)} > 0.$$

This suggests that the variance is estimable. The authors propose a kernel estimator,

$$\hat{P}_{ji} = \frac{\frac{1}{H_N} \sum_{l \neq i, l=1}^N m_{ij} K\left[\frac{x_l - x_i}{H_N}\right]}{\frac{1}{H_N} \sum_{l \neq i, l=1}^N \left[\frac{x_l - x_i}{H_N}\right]}, j = 0, 2.$$

(This is a multivariate kernel in any realistic case. For the application, the authors used a product of Epanechnikov kernel functions. Details on selection of the bandwidth may be found in their paper.) With these estimates of  $P_{10}$  and  $P_{12}$  in hand, the estimator of  $\sigma(\mathbf{x}_i)$  is

$$\hat{\sigma}(\mathbf{x}_i) = \frac{1}{\Phi^{-1}(1 - \hat{P}_{12}) - \Phi^{-1}(\hat{P}_{10})}.$$

The second step MLEs of  $\alpha$  and  $\beta$  are obtained by maximizing a log likelihood function,

$$\log \hat{L} = \sum_{i=1}^N \tau(\mathbf{x}_i) \log \left[ m_{i0} \Phi \left( \frac{-\alpha - \beta' \mathbf{x}_i}{\hat{\sigma}(\mathbf{x}_i)} \right) + m_{i1} \left( \Phi \left( \frac{1 - \alpha - \beta' \mathbf{x}_i}{\hat{\sigma}(\mathbf{x}_i)} \right) - \Phi \left( \frac{-\alpha - \beta' \mathbf{x}_i}{\hat{\sigma}(\mathbf{x}_i)} \right) \right) + m_{i2} \Phi \left( \frac{\alpha + \beta' \mathbf{x}_i - 1}{\hat{\sigma}(\mathbf{x}_i)} \right) \right].$$

Where  $\tau(\mathbf{x}_i)$  is a trimming function “often adopted in two-step estimators, whose support is assumed to be a compact subset of the support of  $\mathbf{x}_i$ . For the study done here,  $\tau(\mathbf{x}) = 0$  if either predicted probability is outside  $[0.005, 0.995]$  and 1 otherwise. The Monte Carlo study that follows agrees with expectations; when the ordered probit model is well specified, it performs well, and when it is not, it performs poorly. Likewise confirming expectations, the authors find that when there is pronounced heteroscedasticity, their estimator outperforms the MLE that assumes homoscedastic disturbances.

## 8.2 A Distribution Free Estimator with Unknown Heteroscedasticity

Lewbel’s (2000) formulation of an *ordered choice* model that allows heteroscedasticity of unknown form is

$$y_i^* = z_i + \beta' \mathbf{x}_i + \sigma_i \varepsilon_i \\ y_i = j \text{ if } \mu_{j-1} < y_i^* \leq \mu_j, j = 0, 1, \dots, J.$$

(We rely heavily on Stewart’s (2005) very concise exposition of this model.) In this instance, the normalization is transferred to one of the slope coefficients. Lewbel’s model is initially formulated in terms of a constant  $\sigma$ , but it is noted that the estimator is robust to heteroscedasticity of unknown form. It is convenient to carry the more general form above. Lewbel’s estimator is noniterative and requires only ordinary least squares regressions. The “special variable,”  $z_i$  whose coefficient is normalized, is required to satisfy certain requirements [see Lewbel (2000) and Stewart (2005).] Among other features, the sign of  $z_i$  must be observed. Then, define the indicator  $y_{.i} = y_i/J$ ; values range from 0 to 1. Then construct,

$$\tilde{y}_{.i} = \frac{y_{.i} - \mathbb{I}[z_i > 0]}{f(z_i | \mathbf{x}_i)} \\ \tilde{y}_{ji} = \frac{\mathbb{I}[y_i > j] - \mathbb{I}[z_i > 0]}{f(z_i | \mathbf{x}_i)}.$$

The numerators are trivial to compute, however, the density of  $z$  given  $x$  requires some additional computation. Stewart (2005, p. 559) navigates some of the developments in the literature for this computation. Assuming the estimator of  $f(z|x)$  is in hand and in the estimator of  $\tilde{y}_{.i}$ , the estimate of  $\beta$  is obtained by least squares regression of  $\tilde{y}_{.i}$  on  $\mathbf{x}_i$ . The estimates of the threshold parameters are the negatives of the constant terms in the  $J-1$  regressions of  $\tilde{y}_{ji}$  on  $\mathbf{x}_i$ .

Lewbel provides this approach for binary, ordered and unordered choice models, censored regressions, and a variety of other settings. Stewart notes, that he found no empirical applications of the ordered choice model, and only a few about binary responses. There have also been “few” studies that compare the estimator to other semiparametric approaches. Little is known about the behavior of this estimator beyond the asymptotic properties that Lewbel, himself, has established in a series of papers [e.g., Lewbel (1997, 2000), Lewbel and Schennach (2007), Honore and Lewbel (2002).]

### 8.3 A Semi-nonparametric Approach

Stewart (2003, 2005) proposes a model that nests the ordered probit model in a general estimator of an unknown density. The alternative density, proposed by Gallant and Nychka (1987) is

$$f_K(\varepsilon) = \frac{1}{\theta} \left( \sum_{k=0}^K \gamma_k \varepsilon^k \right)^2 \phi(\varepsilon).$$

The constant,  $\theta$ , normalizes the density so that it integrates to 1;

$$\theta = \int_{-\infty}^{\infty} \left( \sum_{k=0}^K \gamma_k \varepsilon^k \right)^2 \phi(\varepsilon) d\varepsilon.$$

With the normalization, the density is homogeneous of degree 0 in  $\gamma = (\gamma_0, \dots, \gamma_K)$ , so the normalization  $\gamma_1 = 0$  is imposed. If the remaining  $\gamma_k = 0$ , the normal distribution results. The class of distributions is defined by the order of the polynomial,  $K$ . The model shares a feature with the latent class model examined earlier (Section 5.2.6); the index,  $K$ , is not parametric, and must be located by a specification search. Surprisingly, it turns out that the normal model emerges with  $K = 1$  and  $K = 2$  as well as  $K = 0$ ; the first model in the series that extends the ordered probit model has  $K = 3$ . The model selection problem is a bit more straightforward here in that the order of the model is reduced by one if  $\gamma_K = 0$ , so a likelihood based approach can be used for the specification search.

Stewart notes that the implicit scaling is needed to interpret the coefficients in any ordered choice model. For the application he considers, he suggests that ratios of coefficients are likely to be useful for several reasons. Figure 15 is extracted from Table 1 in Stewart (2005). (An alternative model formulation has been omitted.) The OP and SNP estimates are broadly similar, but the least squares estimates show some pronounced differences from both of the others. The SNP model is a parametric extension of the ordered probit model – hence the name “semi-nonparametric.” It is not in the same class as the Lewbel or Chen and Khan specifications. The likelihood ratio test rejects the ordered probit model. The results in Figure 15 do not include the polynomial parameters or the threshold parameters from the ordered choice models. Figure 16 is Table 2 from Stewart’s earlier study using the same data and a much large model. Moving across the results, we see the changes from  $K=2$  (OP) to the 3 and 5 order polynomials. The hypothesis tests against the null model reject the ordered probit model in both cases. The third order model is also rejected in favor of the fifth order one.

Table 1  
Job satisfaction model—alternative estimators

	Mean	OP	LLS	SNP(3)
log(earnings)	6.66	1	1	1
log('comparison')	6.66	-2.73 (0.98)	-3.24 (1.49)	-2.86 (0.40)
log(hours)	4.95	-1.66 (0.78)	0.20 (3.38)	-1.22 (0.48)
Male	0.50	-1.56 (0.79)	0.56 (0.90)	-1.25 (0.37)
Age/10	3.72	-1.81 (1.12)	-0.51 (2.28)	-1.35 (0.75)
Age <sup>2</sup> /100	15.19	0.33 (0.18)	0.03 (0.28)	0.26 (0.10)
Health	0.18	-2.28 (1.11)	-1.27 (0.51)	-1.85 (0.40)
Second job	0.10	-0.91 (0.66)	1.33 (1.53)	-0.85 (0.44)
Temporary	0.06	-1.44 (0.90)	0.31 (1.00)	-1.16 (0.58)
Manager	0.38	1.68 (0.86)	1.63 (1.03)	1.29 (0.33)
Log-likelihood		-6210.14		-6204.58
L-R test of OP				11.13
[p-value]				[0.001]

Notes:

(1) Standard errors in parentheses.

(2) Estimators: OP = Ordered Probit estimator, LLS = Lewbel least squares estimator, SNP = Semi- nonparametric estimator,

Figure 15 Table 1 From Stewart (2005)

	OP	SNP(3)	SNP(5)
	coef (s.e.)	coef (s.e.)	coef (s.e.)
log(earnings)	0.134 (.054)	0.096 (.050)	0.087 (.056)
log(comp. earn.)	-0.283 (.064)	-0.254 (.060)	-0.323 (.068)
Male	-0.156 (.044)	-0.147 (.044)	-0.130 (.046)
Age/10	-0.210 (.100)	-0.154 (.092)	-0.141 (.104)
Age <sup>2</sup> /100	0.038 (.013)	0.031 (.011)	0.030 (.013)
<u>Thresholds:</u>			
1	-4.125 (.377)	-4.125	-4.125
2	-3.917 (.376)	-3.879 (.043)	-3.847 (.047)
3	-3.562 (.375)	-3.476 (.093)	-3.390 (.075)
4	-2.995 (.375)	-2.877 (.161)	-2.603 (.125)
5	-2.416 (.374)	-2.319 (.211)	-1.889 (.164)
6	-1.664 (.374)	-1.657 (.254)	-1.151 (.177)
<u>Polynomial:</u>			
1		-0.050 (.193)	0.415 (.062)
2		-0.097 (.088)	0.366 (.252)
3		-0.051 (.021)	-0.073 (.171)
4			-0.087 (.034)
5			-0.002 (.016)
Log-likelihood	-6174.25	-6169.87	-6165.48
Standard deviation	1	0.979	1.369
Skewness	0	0.034	0.064
Kurtosis	3	4.600	4.665
Test sum=0 [ $\chi^2(1)$ ]	7.40	9.26	15.83
p-value	[0.007]	[0.002]	[0.000]

Notes: (1) Sample size = 3895. (2) Models also contain 34 other variables.

<sup>3</sup>The other variables included in the model are as in Table 3 of Clark and Oswald (1996). The data used here are from the current release of Wave 1, while they used the original release. The ordered probit results are very close to theirs, but not identical.

**Figure 16 Job Satisfaction Application, Extended**

## 8.4 A Partially Linear Model

Bellemare, Melenberg and van Soest (2002) propose the following ordered choice model based on a partially linear (semiparametric) latent regression, ordered probit model:

$$y_i^* = g(\mathbf{z}_i) + \boldsymbol{\beta}'\mathbf{x}_i + \varepsilon_i$$

$$y_i = j \text{ if } \mu_{j-1} < y_i^* < \mu_j.$$

Their model specifies  $\varepsilon_i \sim N[0, \sigma^2]$ , however,  $\sigma$  remains unidentified. The usual normalizations for the threshold parameters are also required. There is an interesting intersection of the different aspects of “semiparametric” at this point. It seems that concern about the distribution of  $\varepsilon_i$  would be a moot point here; if  $g(\cdot)$  is unspecified, then it seems unlikely that an observed sample could support estimation of a model that is also built around an unspecified density for  $\varepsilon$ . The implied probabilities for the model are

$$\text{Prob}(y_i = j | \mathbf{x}_i, \mathbf{z}_i) = \Phi(\mu_j - g(\mathbf{z}_i) - \boldsymbol{\beta}'\mathbf{x}_i) - \Phi(\mu_{j-1} - g(\mathbf{z}_i) - \boldsymbol{\beta}'\mathbf{x}_i).$$

Estimation of the model is suggested using a technique by Hardle, Huet, Mammen and Sperlich (2004) and Severini and Staniswalis (1994). This involves iterating back and forth between maximum likelihood estimation of  $\boldsymbol{\theta} = (\boldsymbol{\beta}, \boldsymbol{\mu}, \sigma)$  conditioned on estimates of  $g(\mathbf{z}_i)$  and estimates of  $g(\mathbf{z}_i)$  given the other parameters. The former uses the conventional MLE carrying the current estimates of  $g(\mathbf{z}_i)$  as known constants. The latter is accomplished by maximizing a separated weighted likelihood function for each  $i$  to obtain the current estimate of  $g(\mathbf{z}_i)$ .

## 8.5 Semiparametric Analysis

We have examined most of the received developments in the area of semiparametric and nonparametric analyses of the ordered choice model. The central focus of the developments is consistent estimation of the regression slope parameters,  $\boldsymbol{\beta}$  in the absence of an assumption about the distribution or the variance of the disturbance. As we have observed repeatedly in the preceding analyses, however, these elements of the model are crucial for translating the coefficient estimates into meaningful characterizations of the underlying data generating process, and these features are absent by design from the semiparametric estimators. Perhaps the signature feature of the ordered choice model is the vexing result that neither the sign nor the magnitude of  $\boldsymbol{\beta}$  is informative about the impact of interesting right hand variables on the process that generates the outcome variable. For example, Copejans (2007) comments at length on the difference in magnitude of a particular coefficient (a fee elasticity) estimated by the ordered probit MLE compared to that obtained by a distribution-free sieve estimator. But, the difference in magnitude observed there is comparable to the difference that would emerge in the same context if he had used an ordered logit model compared to an ordered probit model. The fact that the scaling induced by the distributional model has been obscured in the estimation process is crucial to the finding. That is, the comparison of the estimates of -0.20 for an ordered probit model to a -0.063 for the semiparametric estimator is meaningless without information on the scaling induced by the underlying distributions. No evidence is available to eliminate the possibility that the partial effect in the ordered probit model is actually larger, not smaller, than that in the semiparametric model.

The presence of unaccounted for heteroscedasticity makes this worse. In the Chen and Khan (2003) model, the heteroscedasticity involves the same  $\mathbf{x}$  as the mean of the regression. The upshot is that in neither model is  $\boldsymbol{\beta}$  the partial effect of interest – indeed, the sign of that

partial effect could be different from that of  $\beta$  in all cells of the outcome, since the mean effect,  $\beta$ , and the variance effect,  $\sigma(\mathbf{x})$  would typically have opposite signs. In their formulation of the model, any partial effect will have to include  $\partial\sigma(\mathbf{x})/\partial\mathbf{x}$ , however,  $\sigma(\mathbf{x})$  is not estimated parametrically; we have no idea what this derivative looks like.

Of all the papers that we examined in this section of the literature (perhaps 10 in total), only one, Stewart (2005) dwells on this issue at any length. As he notes, “Estimated coefficients in the standard parameterization of the Ordered Probit model cannot be interpreted directly and are only identified up to a scale normalization ... However, ratios of coefficients can be usefully interpreted.” Strictly, this claim is correct when the partial effects in the true model obeys the “parallel regressions” feature and it is somewhat misleading as it only applies to a particular outcome – the partial effects change sign and magnitude as one moves through the set of outcomes. That is, when the partial effects are of the form  $\partial\text{Prob}(y = j|\mathbf{x})/\partial x_k = K\beta$  for some  $K$  that is independent of  $k$ . Stewart notes that this feature is useful for examining “indifference curves,” that is, for examining what trades of two variables will leave the outcome (or, underlying preference) unchanged. [Boes and Winkelmann (2006a) pursue this same point at great length.] A second motivation for examining the ratios of coefficients is to see the ratios of specific partial effects, relative to a particular variable. He notes, in the Lewbel formulation, one of the coefficients (that on the “special  $z$ ”) is normalized to 1. As such, each coefficient on another variable is interpretable as relative to this variable. Of course, the normalization could be on any other variable to secure identification of the model, but that would leave Stewart’s observation intact. The ratios of coefficients on other variables to the  $z$  in question would survive renormalization of the model. However, even with all this in place, the analysis hangs on the assumption that the ratios of partial effects in the model equal the ratios of the parameters. In some of the model extensions we have examined, this is not the case.

The upshot of all this is that there is a loose end remaining to be tied up in the development of the semiparametric estimators. In the parametric formulations, the otherwise annoying scale difference between, say, probit and logit estimates is reconciled by the scaling of the model, itself. That reconciliation remains to be developed for the semiparametric approaches. This is needed in order to make the “robust” parameter estimates meaningful.

## 9 CONCLUSIONS

This review began as a short note to propose the new estimator in Section 5.2.7. In researching the recent developments in ordered choice modeling, it appeared that it might be useful to include some pedagogical material about uses and interpretation of the model at the most basic level. We continue to believe that practitioners (and theorists) focus too sharply on coefficient estimation and do not place enough attention on the meaning of the model or its components. As that effort proceeded, it struck us that a more thorough survey of the model, including its historical development might be useful and (we hope) interesting for readers. The preceding is (we hope as well) a survey of the entire literature on the model of ordered choice. (We have, of course, omitted mention of many – perhaps most – of the huge number of applications.)

The development of the ordered choice regression model has emerged in two surprisingly disjoint strands of literature, in its earliest forms in the bioassay literature and in its modern social science counterpart with the pioneering paper by McElvey and Zavoina (1975) and its successors, such as Terza (1985). There are a few prominent links between these two literatures, notably Walker and Duncan (1967). However, even up to the contemporary literature, biological scientists and social scientists have largely successfully avoided bumping into each other.

The earliest applications of modeling ordered outcomes involved grouped data assembled in table format, and with moderate numbers of levels of usually a single stimulus. The fundamental ordered logistic (“cumulative odds”) model in its various forms serves well as an appropriate modeling framework for such data. Walker and Duncan (1967) focused on a major limitation of the approach. When data are obtained with large numbers of inputs – the models in Brewer et al. (2008) involve over 40 covariates – and many levels of those inputs, then crosstabulations are no longer feasible or adequate. Two requirements become obvious, the use of the individual data and the heavy reliance on what amount to multiple regression techniques. McElvey and Zavoina (1975) added to the model a reliance on a formal underlying “data generating process,” the latent regression, a mechanism that makes an occasional appearance in the bioassay treatment, but is never absent from the social science application.

The cumulative odds model for contingency tables and the fundamental ordered probit model for individual data are now standard tools. The recent advances in ordered choice modeling have involved modeling heterogeneity, in cross sections and in panel data sets. These include a variety of threshold models and models of parameter variation such as latent class and mixed and hierarchical models.

## References

- Abramovitz, M. and I. Stegun, 1971. *Handbook of Mathematical Functions*, New York, Dover Press.
- Abrevaya, J., 1997. "The Equivalence of Two Estimators of the Fixed Effects Logit Model," *Economics Letters*, 55, pp. 41-43.
- Acharya, S., 1988. "A Generalized Econometric Model and Tests of a Signalling Hypothesis with Two Discrete Signals," *Journal of Finance*, 43, pp. 413-429.
- Adams, J., 2006. "Learning, Internal Research and Spillovers," *Economics of Innovation and New Technology*, 15, pp. 5-36.
- Agresti, A., 1984. *Analysis of Ordinal Categorical Data*, New York, Wiley.
- Agresti, A., 1990. *Categorical Data Analysis*, New York, John Wiley and Sons.
- Agresti, A., 1999. "Modelling Ordered Categorical Data: Recent Advances and Future Challenges," *Statistics in Medicine*, 18, pp. 2191-2207.
- Aitchison, J. and J. Bennett, 1970. "Polychotomous Quantal Response by Maximum Indicant," *Biometrika*, 57, pp. 253-262.
- Aitchison, J. and S. and Silvey, 1957. "The Generalization of Probit Analysis to the Case of Multiple Responses," *Biometrika*, 44, pp. 131-140.
- Albert, J. and S. Chib, 1993. "Bayesian Analysis of Binary and Polytomous Response Data," *Journal of the American Statistical Association*, 88, pp. 669-679.
- Allison, P., 1999. "Comparing Logit and Probit Coefficients Across Groups," *Sociological Methods and Research*, 28, pp. 186-208.
- Amel, D. and J. Liang, 1994. "A Dynamic Model of Entry and Performance in the U.S. Banking Industry," Board of Governors of the Federal Reserve System, Finance and Economics Discussion Series, 210.
- Amel, D. and J. Liang, 1997. "Determinants of Entry and Profits in Local Banking Markets," *Review of Industrial Organization*, 12, pp. 59-78.
- Amemiya, T., 1975. "Qualitative Response Models," *Annals of Economic and Social Measurement*, 4, pp. 363-372.
- Amemiya, T. 1981. "Qualitative Response Models: A Survey," *Journal of Economic Literature*, 19, 4, pp. 481-536.
- Amemiya, T., 1980. "The  $n^2$  - order Mean Squared Errors of the Maximum Likelihood and the Minimum Logit Chi Squared Estimator," *Annals of Statistics*, 8, pp. 488-505.
- Amemiya, T., 1985a. "Tobit Modeling: A Survey," *Journal of Econometrics*, 24, 1/2, pp. 3-61.
- Amemiya, T., 1985b. *Advanced Econometrics*, Cambridge, Harvard University Press
- Ananth, C. and D. Kleinbaum, 1997. "Regression Models for Ordinal Responses: A Review of Methods and Applications," *International Journal of Epidemiology*, 26, pp. 1232-1333.
- Andersen, D., 1970. "Asymptotic Properties of Conditional Maximum Likelihood Estimators," *Journal of the royal Statistical Society, Series B*, 32, pp. 283-301.
- Anderson, J., 1972. "Separate Sample Logistic Discrimination," *Biometrika*, 59, pp. 19-35.
- Anderson, J., 1984. "Regression and Ordered Categorical Variables," *Journal of the Royal Statistical Society, Series B (Methodological)*, 46, pp. 1-30.
- Anderson, J. and P. Philips, 1981. "Regression, Discrimination and Measurement Models for Ordered Categorical Variables," *Applied Statistics*, 30, pp. 22-31.
- Ando, T., 2006. "Bayesian credit rating analysis based on ordered probit regression model with functional predictor," *Proceeding of The Third IASTED International Conference on Financial Engineering and Applications*, 69-76.
- Andrews, D. and W. Ploberger, 1994. "Optimal Tests when a Nuisance Parameter is Present Only Under the Alternative," *Econometrica*, 62, pp. 1383-1414.

- Andrich, D., 1979, "A Model for Contingency Tables Having an Ordered Response Classification," *Biometrics*, 35, 403-415.
- Armstrong, D. and J. McVicar, 2000. "Value Added in Further Education and Vocational Training in Northern Ireland," *Applied Economics*, 32, pp. 1727-1736.
- Barnhart, H. and A. Sampson, 1994. "Overview of Multinomial Models for Ordered Data," *Communications in Statistics – A. Theory and Methods*, 23, pp. 3395-3416.
- Baltagi, B., 2007. *Econometric Analysis of Panel Data*, New York, John Wiley and Sons.
- Basu, D. and R. de Jong, 2006. "Dynamic Multinomial Ordered Choice With An Application to the Estimation of Monetary Policy Rules," Department of Economics, Ohio State University, Manuscript.
- Becker, W. and P. Kennedy, 1992. "A Graphical Exposition of the Ordered Probit Model," *Econometric Theory*, 8, pp. 127-131.
- Bedi, A. and I. Tunalı, 2004. "Testing for Market Imperfections: Participation in Land and Labor Contracts in Turkish Agriculture," Working Paper, Institute of Social Studies, The Hague.
- Bellemare, C., B. Melenbert and A. van Soest, 2002. "Semi-parametric Models for Satisfaction with Income," *Portuguese Economic Journal*, 1, pp. 181-203.
- Bera, A., C. Jarque and L. Lee, 1984. "Model Specification Tests: A Simultaneous Approach," *Journal of Econometrics*, 20, pp. 59-82.
- Berndt, E., B. Hall, R. Hall and J. Hausman, 1974. "Estimation and Inference in Nonlinear Structural Models," *Annals of Economic and Social Measurement*, 3/4, pp. 653-665.
- Berkson, J., 1944. "Application of the Logistic Function to Bioassay," *Journal of the American Statistical Association*, 39, pp. 357-365.
- Berkson, J., 1951. "Why I Prefer Logits to Probits," *Biometrics*, 7, 4, pp. 327-339.
- Berkson, J., 1953. "A Statistically Precise and Relatively Simple Method of Estimating the Bioassay with Quantal Response, Based on the Logistic Function," *Journal of the American Statistical Association*, 48, pp. 565-599.
- Berkson, J., 1955a. "Maximum Likelihood and Minimum  $\chi^2$  Estimates of the Logistic Function," *Journal of the American Statistical Association*, 50, pp. 130-162.
- Berkson, J., 1955b. "Estimate of the Integrated Normal Curve by Minimum Normit Chi-Square with Particular Reference to Bioassay," *Journal of the American Statistical Association*, 50, pp. 529-550.
- Berkson, J., 1957. "Tables for Use In Estimating the Normal Distribution Function by Normit Analysis," *Biometrika*, 44, pp. 411-435.
- Berkson, J., 1980. "Minimum Chi-Square, Not Maximum Likelihood," *Annals of Statistics*, 8, pp. 457-487.
- Bhat, C., 1994. "Imputing a Continuous Income Variable from Grouped and Missing Income Observations," *Economics Letters*, 46, 4, pp. 311-320.
- Biswas and Das, 2002. "A Bayesian Analysis of Bivariate Ordinal Data: Wisconsin Epidemiologic Study of Diabetic Retinopathy Revisited," *Statistics in Medicine*, 21, 4, pp. 549-559
- Bliss, C., 1934a. "The Method of Probits," *Science*, 79, 2037, pp. 38-39.
- Bliss, C., 1934b. "The Method of Probits: A Correction," *Science*, 79, 2053, pp. 409-410.
- Boes, S., 2007. "Nonparametric Analysis of Treatment Effects in Ordered Response Models," University of Surich, Socioeconomic Institute, Working Paper 0709.
- Boes, S. and R. Winkelmann, 2004. "Income and Happiness: New Results from Generalized Threshold and Sequential Models," IZA Discussion Paper No. 1175, SOI Working Paper 0407, IZA
- Boes, S. and R. Winkelmann, 2006a. "Ordered Response Models," *Allgemeines Statistisches Archiv*, 90, 1, pp. 165-180.

- Boes, S. and R. Winkelmann, 2006b. "The Effect of Income on Positive and Negative Subjective Well-Being," University of Zurich, Socioeconomic Institute, Manuscript, IZA Discussion Paper Number 1175.
- Boyes, W., D. Hoffman and S. Low, 1989. "An Econometric Analysis of the Bank Credit Scoring Problem," *Journal of Econometrics*, 40, pp. 3-14.
- Brant, R., 1990. "Assessing Proportionality in the Proportional Odds Model for Ordered Logistic Regression," *Biometrics*, 46, pp. 1171-1178.
- Bresnahan, T. F., 1987. "Competition and Collusion in the American Automobile Industry: The 1955 Price War," *Journal of Industrial Economics*, 35, pp. 457-482.
- Breusch, T. and A. Pagan, 1979. "A Simple Test for Heteroscedasticity and Random Parameter Variation," *Econometrica*, 47, pp. 1287-1294.
- Brewer, C., C. Kovner, W. Greene, Y. Cheng, 2008. "Predictors of RNs' Intent to Work and Work Decisions One Year Later in a U.S. National Sample," *The International Journal of Nursing Studies*, forthcoming.
- Britt, C., 2000. "Comment on Paternoster and Brame," *Criminology*, 38, 3, pp. 965-970.
- Brooks, R., M. Harris and C. Spencer, 2007. "A Inflated Ordered Probit Model of Monetary Policy: Evidence from MPC Voting Data," Department of Econometrics and Business Statistics, Monash University, Manuscript.
- Buckle, R. and J. Carlson 2000, 2000. "Menu Costs, Firm Size and Price Rigidity," *Economics Letters*, 66, pp. 59-63.
- Buliung, R., 2005. "Activity/Travel Behaviour Research: Approaches and Findings with Identification of Researching Themes and Emerging Methods," Center for Spatial Analysis, McMaster Univ, Working paper 008.  
(<http://www.science.mcmaster.ca/cspa/papers/CSpA%20WP%20008.pdf>).
- Butler, J., T. Finegan and J. Siegfried, 1994. "Does More Calculus Improve Student Learning in Intermediate Micro and Macro Economic Theory?" *American Economic Review*, 84, 2, pp. 206-210.
- Butler, J., T. Finegan and J. Siegfried, 1998. "Does More Calculus Improve Student Learning in Intermediate Micro- and Macroeconomic Theory?" *Journal of Applied Econometrics*, 13, pp. 185-202.
- Butler, J. and P. Chatterjee, 1995. "Pet Econometrics: Ownership of Cats and Dogs," Department of Economics, Vanderbilt University, Working Paper Number 95-WP1.
- Butler, J. and P. Chatterjee, 1997. "Tests of the Specification of Univariate and Bivariate Ordered Probit," *Review of Economics and Statistics*, 79, pp. 343-347.
- Butler, J. and R. Moffitt, 1982. "A Computationally Efficient Quadrature Procedure for the One Factor Multinomial Probit Model," *Econometrica*, 50, pp. 761-764.
- Calhoun, C. A. 1986. "BIVOPROB: Maximum Likelihood Program for Bivariate Ordered-Probit Regression." Washington, D.C.: The Urban Institute.
- Calhoun, C., 1989. "Estimating the Distribution of Desired Family Size and Excess Fertility," *Journal of Human Resources*, 24, 4, pp. 709-24.
- Calhoun, C., 1991. "Desired and Excess Fertility in Europe and the United States: Indirect Estimates from World Fertility Survey Data," *European Journal of Population*, 7, pp. 29-57.
- Calhoun, C., 1994. "The Impact of Children on the Labor Supply of Married Women: Comparative Estimates from European and U.S. Data," *European Journal of Population*, 10, pp. 293-318.
- Calhoun, C., 1995. "BIVOPROB: Computer Program for Maximum Likelihood Estimation of Bivariate Ordered Probit Models for Censored Data: Version 11.92," *Economic Journal*, 105, pp. 786-787.

- Cameron, S. and J. Heckman, 1998. "Life Cycle Schooling and Dynamic Selection Bias: Models and Evidence for Five Cohorts of American Males," *Journal of Political Economy*, 106, pp. 262–333.
- Carneiro, P., K. Hansen and J. Heckman, 2001. "Removing the Veil of Ignorance in Assessing the Distributional Impacts of Social Policies," *Swedish Economic Policy Review*, 8, pp. 273-301.
- Carneiro, P., K. Hansen and J. Heckman, 2003. "Estimating Distributions of Treatment Effects with an Application to Schooling and Measurement of the Effects of Uncertainty on College Choice," *International Economic Review*, 44, pp. 361-422.
- Carro, J., 2007. "Estimating Dynamic Panel Data Discrete Choice Models with Fixed Effects," *Journal of Econometrics*, 140, pp. 503-528.
- Chamberlain, G., 1980. "Analysis of Covariance with Qualitative Data," *Review of Economic Studies*, 47, pp. 225-238.
- Chen, S. and S. Khan, 2003. "Rates of Convergence for Estimating Regression Coefficients in Heteroskedastic Discrete Response Models," *Journal of Econometrics*, 117, pp. 245-278.
- Chesher, A. and M. Irish, 1987. "Residual Analysis in the Grouped Data and Censored Normal Linear Model," *Journal of Econometrics*, 34, pp. 33-62.
- Cheung, S., 1996. "Provincial Credit Rating in Canada: An Ordered Probit Analysis," Bank of Canada, Working Paper 96-6. (<http://www.bankofcanada.ca/en/res/wp/1996/wp96-6.pdf>)
- Christensen, K., H. Kohler, O. Basso, Olga, J. Olsen, J. Vaupel and J. Rodgers, 2003. "The Correlation of Fecundability Among Twins: Evidence of a Genetic Effect on Fertility?" *Epidemiology*, 14, 1, pp. 60-64.
- Clark, A., Y. Georgellis and P. Sanfey, 2001. "Scarring: The Psychological Impact of Past Unemployment," *Economica*, 68, pp. 221-241. et al. 2001
- Clogg, C. and E. Shihadeh, 1994. *Statistical Models for Ordered Variables*, Thousand Oaks, CA, Sage Publications.
- Contoyannis, A., A. Jones and N. Rice, 2004. "The Dynamics of Health in the British Household Panel Survey," *Journal of Applied Econometrics*, 19, 4, pp. 473-503.
- Coppejans, M., 2007. "On Efficient Estimation of the Ordered Response Model," *Journal of Econometrics*, 137, pp. 577-614.
- Cox, C., 1995. "Location-Scale Cumulative Odds Models for Ordered Data: A Generalized Nonlinear Model Approach," *Statistics in Medicine*, 14, pp. 1191-1203.
- Cox, D., 1970. *Analysis of Binary Data*, Methuen, London.
- Crawford, D., R. Pollak and F. Vella, 1988. "Simple Inference in Multinomial and Ordered Logit," *Econometric Reviews*, 17, pp. 289–299.
- Crouchley, R., 1995. "A Random Effects model for Ordered Categorical Data," *Journal of the American Statistical Association*, 90, pp. 489-498.
- Crouchley, B., 2005. "E-Science, The GRID and Statistical Modelling in Social Research," (<http://www.ccsr.ac.uk/methods/festival/programme/gss/crouchley.ppt>).
- Cunha, F., J. Heckman and S. Navarro, 2007. "The Identification & Economic Content of Ordered Choice Models with Stochastic Thresholds," University College Dublin, Gery Institute, Discussion Paper WP/26/2007.
- Czado, C., A. Heyn and G. Müller, 2005. "Modeling Migraine Severity with Autoregressive Ordered Probit Models," Technische Universität München, Working paper number 463.
- D'Addio, A., T. Eriksson and P. Frijters, 2007. "An Analysis of the Determinants of Job Satisfaction When Individuals' Baseline Satisfaction Levels May Differ," *Applied Economics*, 39, 19, pp. 2413-2423.
- Dardanomi, V. and A. Forcina, 2004. "Multivariate Ordered Logit Regressions," University of Palermo, Manuscript. ([http://www.cide.info/conf\\_old/papers/11128.pdf](http://www.cide.info/conf_old/papers/11128.pdf))
- Das, M. and A. van Soest, 2000. "A Panel Data Model for Subjective Information on Household Income Growth," *Journal of Economic Behavior and Organization*, 15, pp. 401-416.

- Davidson, R. and J. MacKinnon, 1983. "Small Sample Properties of Alternative Forms of the Lagrange Multiplier Test," *Economics Letters*, 12, pp. 269-275.
- Davidson, R. and J. MacKinnon, 1984. "Model Specification Tests Based on Artificial Linear Regressions," *International Economic Review*, 25, pp. 485-502.
- Davidson, R. and J. MacKinnon, 1993. *Estimation and Inference in Econometrics*, Oxford, Oxford University Press.
- Daykin, A. and P. Moffatt, 2002. "Analyzing Ordered Responses: A Review of the Ordered Probit Model," *Understanding Statistics*, 1, 3, pp. 157-166.
- DeMaris, A., 2004. *Regression with Social Data: Modeling Continuous and Limited Response Variables*, Hoboken, New Jersey, John Wiley and Sons.
- Dempster, A., N. Laird and D. Rubin 1977. "Maximum Likelihood Estimation from Incomplete Data via the EM Algorithm," *Journal of the Royal Statistical Society, Series B*, 39, pp. 1-38.
- Diggle, R., P. Liang and S. Zeger, 1994. *Analysis of Longitudinal Data*, Oxford: Oxford University Press.
- Drezner, Z., 1978. "Computation of the Bivariate Normal Integral," *Mathematics of Computation*, 32, 141, pp. 277-279.
- Dueker, M., S. Spuirr, A. Jacox and D. Kalist, 2005 "The Practice Boundaries of Advanced Practice Nurses: An Economic and Legal Analysis," Federal Reserve Bank of St. Louis, Working Paper 2005-071A
- Dupor, B., T. Mirzoev, T. Conley, T., 2004. "Does the Federal Reserve Do What It Says It Expects to Do?" Working Paper, Department of Economics, Ohio State University, Manuscript.
- Econometric Software, 2007. *NLOGIT: Version 4.0*, Plainview, New York.
- Eichengreen, B., M. Watson and R. Grossman, 1985. "Bank Rate Policy Under the Interwar Gold Standard: A Dynamic Probit Approach," *Economic Journal*, 95, pp. 725-745.
- Ekholm, A. and J. Palmgren, 1989. "Regression Models for an Ordinal Response Are Best Handled As Nonlinear Models," *GLIM Newsletter*, 18, pp. 31-35.
- Eluru, N., C. Bhat and D. Hensher, 2007. "A Mixed Generalized Ordered Response Model for Examining Pedestrian and Bicyclist Injury Severity Levels in Traffic Crashes," *Accident Analysis and Prevention*, 40, 3, pp. 1033-1054..
- Everitt, B., 1988. A Finite Mixture Model for the Clustering of Mixed-Mode Data," *Statistics and Probability Letters*, 6, pp. 305-309.
- EViews, 2008. *Eviews Version 6.0*, QMS, Irvine, CA.
- Fahrmeir, L. and G. Tutz, 1994. *Multivariate Statistical Modeling Based on Generalized Linear Models*, Berlin, Springer Verlag.
- Farewell, V., 1982. "A Note on Regression Analysis of Ordinal Data with Variability of Classification," *Biometrika*, 69, pp. 533-538.
- Feinberg, S., 1980. *The Analysis of Cross-Classified Categorical Data*, Cambridge, MIT Press.
- Ferrer-i-Carbonell, A. and P. Frijters, 2004. "How Important Is Methodology for the Estimates of the Determinants of Happiness," *Economic Journal*, 114, pp. 641-659.
- Fernandez-Val I., 2008. "Fixed Effects Estimation of Structural Parameters and Marginal Effects in Panel Probit Models," *Journal of Econometrics*, Forthcoming
- Fernández-Val, I. and F. Vella, 2007. "Bias Corrections for Two-Step Fixed Effects Panel Data Estimators," IZA Working Papers Number 2690.
- Filer, R. and M. Honig, 2005. "Endogenous Pensions and Retirement Behavior, Department of Economics, Hunter College, Manuscript.
- Finney, D. 1944a. "The Application of the Probit Method to Toxicity Test Data Adjusted for Mortality in the Control", *Annals of Applied Biology*, 31, pp.68-74.
- Finney, D., 1944b. "The Application of Probit Analysis to the Results of Mental Tests", *Psychometrika*, 9, pp. 31-39.

- Finney, D. 1947, "The Principles of Biological Assays", *Journal of the Royal Statistical Association B*, 9, pp. 46-91.
- Finney, D., 1947. *Probit analysis: A Statistical Treatment of the Sigmoid Response Curve*, Cambridge: Cambridge University Press.
- Finney, D., 1952. *Probit Analysis: A Statistical Treatment of the Sigmoid Response Curve, 2<sup>nd</sup> Edition*, Cambridge: Cambridge University Press.
- Finney, D., 1971. *Probit Analysis: A Statistical Treatment of the Sigmoid Response Curve, 3<sup>rd</sup> Edition*, Cambridge: Cambridge University Press.
- Formisano, J., K. Still, W. Alexander and M. Lippmann, 2001. "Application of Statistical Models for Secondary Data Usage of the U.S. Navy's Occupational Exposure Database (NOED)," *Applied Occupational and Environmental Hygiene*, 16, pp. 201-209.
- Frazis, H., 1993. "Selection Bias and the Degree Effect," *Journal of Human Resources*, 28, pp. 538-554.
- Freedman, D., 2006. "On the So-Called 'Huber Sandwich Estimator' and Robust Standard Errors," *The American Statistician*, 60, 4, pp. 299-302.
- Frijters, P., J. Haisken-DeNew and M. Shields, 2004. "The Value of Reunification in Germany: An Analysis of Changes in Life Satisfaction," *Journal of Human Resources*, 39, 3, pp. 649-674.
- Fu, V., 1998. "Estimating Generalized Ordered Logit Models," *Stata Technical Bulletin*, 44, pp. 27-30.
- Fu, A., M. Gordon, G. Liu B. Dale and R. Christensen, 2004. "Inappropriate Medication Use and Health Outcomes in the Elderly" *Journal of the American Geriatrics Society*, 52, 11, pp. 1934-1939.
- Gallant, R. and D. Nychka, 1987. "Semi-nonparametric Maximum Likelihood Estimation," *Econometrica*, 55, pp. 363-390.
- Garen, J., 1984. "The Returns to Schooling: A Selectivity Bias Approach with a Continuous Choice Variable," *Econometrica*, 52, 5, pp. 1199-1218.
- Genius, M., C. Pantzios, V. Tzouvelakis, 2005. "Information Acquisition and Adoption of Organic Farming Practices: Evidence from Farm Operations in Crete, Greece" Department of Economics, University of Crete, Manuscript.
- Gaddum, J., 1933. "Reports on Biological Standards, III. Methods of Biological Assay Depending on a Quantal Response" Special Report Series 183. Medical Research Council, HM Statistical Office, London.
- Gallant, R. and D. Nychka, 1987. "Semi-nonparametric Maximum Likelihood Estimation," *Econometrica*, 55, pp. 363-390.
- Genberg, H. and S. Gerlach, 2004. "Estimating Central Bank Reaction Functions with Ordered Probit: A Note" Graduate Institute of International Studies, Geneva, Manuscript.
- Geweke, J., 1991. Efficient Simulation From the Multivariate Normal and Student t-Distributions Subject to Linear Constraints," in *Computer Sciences and Statistics Proceedings of the 23<sup>d</sup> Symposium on the Interface*, pp. 571-578.
- Girard, P. and E. Parent, 2001. "Bayesian Analysis of Autocorrelated Ordered Categorical Data for Industrial Quality Monitoring," *Technometrics*, 43, 2, pp. 180-191.
- Glewwe, P., 1997. "A Test of the Normality Assumption in the Ordered Probit Model," *Econometric Reviews*, 16, 1, pp. 1-19.
- Glewwe, P. and H. Jacoby, 1994. "Student Achievement and Schooling Choice in Low-Income Countries: Evidence from Ghana," *Journal of Human Resources*, 29, 3, pp. 843-864.
- Glewwe, P. and H. Jacoby, 1995. "An Economic Analysis of Delayed Primary School Enrollment in a Low Income Country: The Role of Early Childhood Malnutrition," *Review of Economics and Statistics*, 77, 1, pp. 156-169.
- Godfrey, L., 1988. *Misspecification Tests in Econometrics*, Cambridge, Cambridge University Press.

- Gourieroux, C., A. Monfort and E. Renault, 1987. "Generalized Residuals," *Journal of Econometrics*, 34, pp. 5-32.
- Greene, W., 1981. "Sample Selection Bias As a Specification Error: Comment," *Econometrica*, 49, pp. 795-798.
- Greene, W., 1990. *Econometric Analysis*, New York, Macmillan.
- Greene, W., 1994. "Accounting for Excess Zeros and Sample Selection in Poisson and Negative Binomial Regression Models," Working Paper 94-10, Department of Economics, Stern School of Business, New York University.
- Greene, W., 1995. "Sample Selection in the Poisson Regression Model," Department of Economics, Stern School of Business, New York University, Working paper #95-06, 1995.
- Greene, W., 2002. *LIMDEP Version 8.0, Reference Guide*, Plainview, NY, Econometric Software.
- Greene, W., 2003. *Econometric Analysis, 5<sup>th</sup> Edition*, Englewood Cliffs, Prentice Hall.
- Greene, W., 2004a. "Fixed Effects and Bias Due To The Incidental Parameters Problem in the Tobit Model," *Econometric Reviews*, 23, 2, pp. 125-147.
- Greene, W., 2004b. "The Behavior of the Fixed Effects Estimator in Nonlinear Models," *The Econometrics Journal*, 7, 1, pp. 98-119.
- Greene, W., 2005. "Functional form and Heterogeneity in Models for Count Data," *Foundations and Trends in Econometrics*, 1, 2, pp. 113-218.
- Greene, W., 2006. "A General Approach to Incorporating Selectivity in a Model," Department of Economics, Stern School of Business, New York University, Working Paper 06-10.
- Greene, W., 2007a. *LIMDEP Version 9.0: Reference Guide*, Plainview, New York, Econometric Software, Inc.
- Greene, W., 2007b. *NLOGIT Version 4.0: Reference Guide*, Plainview, NY, Econometric Software.
- Greene, W., 2008a. *Econometric Analysis, 6<sup>th</sup> Edition*, Englewood Cliffs, Prentice Hall.
- Greene, W., 2008b. "A Stochastic Frontier Model with Correction for Selection," Department of Economics, Stern School of Business, New York University, Working Paper EC-08-09.
- Greene, W., M. Harris, B. Hollingworth, P. Maitra, 2008. "A Bivariate Latent Class Correlated Generalized Ordered Probit Model with an Application to Modeling Observed Obesity Levels," Department of Economics, Stern School of Business, New York University, Working Paper 08-18.
- Greene, W. and D. Hensher, 2008. "Ordered Choices and Heterogeneity in Attribute Processing," IITLS, Sydney University, Manuscript.
- Greene, W., L. Knapp and T. Seaks, 1993. "Estimating the Functional Form of the Independent Variables in Probit Models," *Applied Economics*, pp. 193-196.
- Greenland, S., 1994. "Alternative Models for Ordinal Logistic Regression," *Statistics in Medicine*, 13, 1665-1677.
- Greenwood, C. and V. Farewell, 1988. "A Comparison of Regression Models for Ordinal Data in Analysis of Transplanted Kidney Function," *Canadian Journal of Statistics*, 16, pp. 325-336.
- Grizzle, J., C. Starmet and G. Koch, 1969. "Analysis of Categorical Data By Linear Models," *Biometrics*, 25, pp. 489-504.
- Gurland, J., J. Lee and P. Dahm, 1960. "Polychotomous Quantal Response in Biological Assay," *Biometrics*, 16, pp. 382-398.
- Gustaffson, S. and Stafford, F., 1992. "Child Care Subsidies and Labor Supply in Sweden," *Journal of Human Resources*, 27, pp. 204-230.
- Hamermesch, D., 2004. "Subjective Outcomes in Economics," *Southern Economic Journal*, 71, 1, pp. 2-11.
- Han, A. and J. Hausman, 1988. "Semiparametric Estimation of Duration and Competing Risk Models," Department of Economics, MIT, Manuscript.

- Hahn, J. and G. Kuersteiner, 2003. "Bias Reduction for Dynamic Nonlinear Panel Data Models with Fixed Effects," Department of Economics, UCLA, Manuscript.
- Hahn, J. and W. Newey, 2004. "Jackknife and Analytical Bias Reduction for Nonlinear Panel Models," *Econometrica*, 72, 4, pp. 1295-1319.
- Han, A. and J. Hausman, 1986. "Semiparametric Estimation of Duration and Competing Risk Models," Department of Economics, MIT, Working Paper 450.
- Harvey, A., 1976. "Estimating Regression Models with Multiplicative Heteroscedasticity," *Econometrica*, 44, pp. 461-465.
- Hardle, W., S. Huet, E. Mammen and E. Sperlich, 2004. "Bootstrap Inference in Semiparametric Generalized Additive Models," *Econometric Theory*, 20, pp. 265-300.
- Hausman, J., 1978. "Specification Tests in Econometrics," *Econometrica*, 46, pp. 1251-1271.
- Hausman, J., A. Lo and C. MacKinlay, 1992. "An Ordered Probit Analysis of Transaction Stock Prices," *Journal of Financial Economics*, 31, pp. 319-379.
- Harris, M. and X. Zhao, 2007. "Modeling Tobacco Consumption with a Zero Inflated Ordered Probit Model," *Journal of Econometrics*, 141, pp. 1073-1099.
- Heckman, J., 1979. "Sample Selection Bias as a Specification Error," *Econometrica*, 47, pp. 153-161.
- Heckman, J., 1981. "Heterogeneity and State Dependence," in S. Rosen, ed., *Studies in Labor Markets*, Chicago, University of Chicago Press.
- Heckman, J. J. and T. E. MaCurdy, 1981. "New Methods for Estimating Labor Supply Functions," in R. Ehrenberg, ed., *Research in Labor Economics*, Greenwich, CT., JAI Press, pp. 65-102.
- Heckman, J. and B. Singer, 1984. "A Method for Minimizing the Impact of Distributional Assumptions in Econometric Models," *Econometrica*, 52, pp. 271-320.
- Heilbron, D., 1994. "Zero-Altered and Other Regression Models for Count Data with Added Zeros," *Biometrical Journal*, 36, pp. 531-547.
- Hemmingsen, A., 1933. "The Accuracy of Insulin Assay on White Mice" *Quarterly Journal of Pharmacy and Pharmacology*, 6, 39-80 and 187-283.
- Hensher, D., 2006. "How Do Respondents Process Stated Choice Experiments? – Attribute Consideration Under Varying Information Load," *Journal of Applied Econometrics*, 21, pp. 861-878.
- Herbert, A., N. Gerry and N. McQueen, 2006. "A Common Genetic Variant is Associated with Adult and Childhood Obesity," *Science*, 312, pp. 279-283.
- Hinde, J., G. Clarice and Demetrio, 1998. "Overdispersion in Models and Estimation," *Computational Statistics and Data Analysis*, 27, 2, pp. 151-170.
- Holmes, M. and R. Williams, 1954. "The Distribution of Carriers of Streptococcus Pyogenes Among 2413 Healthy Children," *Journal of Hygiene*, 52, pp. 165-179.
- Honore, B. and A. Lewbel, 2002. "Semiparametric Binary Choice Panel Data Models without Strictly Exogenous Regressors," *Econometrica*, 70, pp. 2053-2063.
- Horowitz, J., 1992. "A Smoothed Maximum Score Estimator for the Binary Response Model," *Econometrica*, 60, pp. 505-531.
- Hsiao, C., 1986. *Analysis of Panel Data*, Cambridge, Cambridge University Press.
- Hsiao, C., 2003. *Analysis of Panel Data, 2<sup>nd</sup> Ed.*, Cambridge, Cambridge University Press.
- Hutchison, V., 1985. "Ordinal Variable Regression Using the McCullagh (Proportional Odds) Model," *Canadian Journal of Statistics*, 16, pp. 325-336.
- Imai, K., G. King and O. Lau, 2008. "Zelig: Everyone's Statistical Software," Department of Government, Harvard University. (<http://gking.harvard.edu/zelig/>)
- ISI, 1982. "Citation Classic: Finney D. J. Probit Analysis," ISI, 31, August.
- Jansen, J., 1990. "On the Statistical Analysis of Ordinal Data when Extravariation is Present," *Applied Statistics*, 39, pp. 75-84.

- Jiminez , E. and B. Kugler, 1987. "The Earnings Impact of Training Duration in a Developing County: An Ordered Probit Selection Model of Colombia's Servicio Nacional de Aprendizaje," *Journal of Human Resources*, 22, 2, pp. 228-247.
- Johnson, N., S. Kotz and A. Balakrishnan, 1994. *Continuous Univariate Distributions, Vol. I*, New York, John Wiley and Sons.
- Johnson, P., 1996. "A Test of the Normality Assumption in the Ordered Probit Model," *Metron*, 54, pp. 213-221.
- Jones, S. and D. Hensher, 2004. "Predicting Firm Financial Distress: A Mixed Logit Model," *The Accounting Review (American Accounting Association)*, 79, pp. 1011-1038.
- Kadam, A and P. Lenk, 2008. "Bayesian Inference for Issuer Heterogeneity in Credit Ratings Migration" . *Journal of Banking and Finance*, Forthcoming.  
(SSRN:<http://ssrn.com/abstract=1084006>)
- Kalman, R., 1960. "A New Approach to Linear Filtering and Prediction Problems," *Journal of Basic Engineering (ASME Transactions)*, 82D, pp. 35-45.
- Kao, C. and C. Wu, 1990. "Two Step Estimation of Linear Models with Ordinal Unobserved Variables: The Case of Corporate Bonds," *Journal of Business and Economic Statistics*, 8, pp. 317-325.
- Kasteridis, P., M Munkin, and S. Yen., 2008. "A Binary-Ordered Probit Model of Cigarette Demand." *Applied Economics*, 41, forthcoming.
- Keele, L. and D. Park, 2005. "Difficult Choices: An Evaluation of Heterogeneous Choice Models," Presented at the 2004 Meeting of the American Political Science Association, Department of Politics and International Relations, Oxford University, Manuscript.
- Kenny, L., L. Lee, G. Maddala and R. rost, 1979. "Returns to College Education: An Investigation of Self-Selection Bias and the Project Talent Data," *International Economic Review*, 20, 3, pp. 775-789.
- Kerkhofs, M., Lindeboom, M., 1995. "Subjective Health measures and State Dependent Reporting Errors" *Health Economics* 4, pp. 221-235.
- Kim, K., 1995. "A Bivariate Cumulative Probit Regression Model for Ordered Categorical Data," *Statistics in Medicine*, 14, pp. 1341-1352.
- Klein, R. and R. Sherman, 2002. "Shift Restrictions and Semiparametric estimation in Ordered Response Models," *Econometrica*, 70, pp. 663-692.
- Klein, R. and R. Spady, 1993. "An Efficient Semiparametric Estimator for Discrete Choice Models," *Econometrica*, 61, pp. 387-421.
- Kohler, H. and J. Rodgers, 1999. "DF-Like Analyses of Binary, Ordered and Censored Variables using Probit and Tobit Approaches," *Behavior Genetics*, 20, 4, pp. 221-232.
- Koop, G. and J. Tobias, 2006. "Semiparametric Bayesian Inference in Smooth Coefficient Models," *Journal of Econometrics*, 134, 1, pp. 283-315.
- Krailo, M. and M. Pike, 1984. "Conditional Multivariate Logistic Analysis of Stratified Case-Control Studies," *Applied Statistics*, 44, 1, pp. 95-103.
- Kuriama, K., Y. Kitibatake, Y. Oshima, 1998. "The Downward Bias Due to "No-Vote" Option in Contingent Valuation Study, World Congress of Environmental and Resource Economists, Venice. (<http://www.f.waseda.jp/kkuri/research/workingpaper/WP9805.PDF>) (1998)
- Lancaster, T., 2000. "The Incidental Parameters Problem Since 1948," *Journal of Econometrics*, 95, pp. 391-413.
- Lambert, D., 1992. "Zero-inflated Poisson Regression With An Application To Defects In Manufacturing," *Technometrics*, 34, 1, pp. 1-14.
- Lawrence, C. and H. Palmer, 2002. Heuristics, Hillary Clinton and Health care Reform, Annual Meeting of the Midwest Political Science Association, Chicago.
- Lechner, M., 1991. "Testing Logit Models in Practice," *Empirical Economics*, 16, pp. 77-108.

- Lee, L. and R. Trost, 1978. Estimation of Some Limited Dependent Variable Models with Application to Housing Demand," *Journal of Econometrics*, 8, pp. 357-382.
- Lewbel, A., 1997. "Semiparametric Estimation of Location and Other Discrete Choice Moments," *Econometric Theory*, 13, pp. 32-51.
- Lewbel, A. and S. Schennach, 2007. "A Simple Ordered Data Estimator for Inverse Density Weighted Expectations," *Journal of Econometrics*, 136, pp.189-211..
- Lewbel, A., 2000. Semiparametric Qualitative Response Model Estimation with Unknown Heteroscedasticity or Instrumental Variables," *Journal of Econometrics*, 97, pp. 145-177.
- Li, M. and J. Tobias, 2006a. "Calculus Attainment and Grades Received in Intermediate Economic Theory," *Journal of Applied Econometrics*, 21,6, pp. 893-896.
- Li, M. and J. Tobias 2006b. "Bayesian Analysis of Structural Effects in an Ordered Equation System," Department of Economics, Iowa State University, Working Paper.
- Li, M. and J. Tobias 2006c. "Bayesian Analysis of Structural Effects in An Ordered Equation System," *Studies in Nonlinear Dynamics & Econometrics*, 10, 4, Article 7, pp. 1-24.
- Lillard, L. and E. King, 1987. "Education Policy and Schooling Attainment in Malaysia and the Philippines," *Economics of Education Review*," 6, pp. 167-181.
- Lindeboom, M. and E. van Doorslayer, 2003. "Cut Point Shift and Index Shift in Self Reported Health," *Ecuity III Project Working Paper #2*.
- Long, S. 1997. *Regression Models for Categorical and Limited Dependent Variables*, Thousand Oaks, CA, Sage Publications.
- Long, S. and J. Freese, 2006. *Regression Models for Categorical Dependent Variables Using Stata*, College Station, Stata Press.
- Machin, S. and A. Vignoles, 2005. *What's the Good of Education? The Economics of Education in the UK*, Princeton, N.J., Princeton University Press.
- MacKinnon, J., 1992. "Model Specification Tests and Artificial Regressions," *Journal of Economic Literature*, 30, 1, pp. 102-146.
- Maddala, J., 1983. *Limited Dependent and Qualitative Variables in Econometrics*, Cambridge, Cambridge University Press.
- Magee, L., J. Burbidge and A. Robb, 2000. "The Correlation Between Husband's and Wife's Education: Canada , 1971-1996," SEDAP Research Paper #24, McMaster.
- Manski, C., 1975. "The Maximum Score Estimator of the Stochastic Utility Model fo Choice," *Journal of Econometrics*, 3, pp. 205-228.
- Manski, C., 1985. "Semiparametric Analysis of Discrete Response: Asymptotic Properties of the Maximum Score Estimator," *Journal of Econometrics*, 27, pp. 313-333.
- Manski, C., 1986. "Operational Characteristics of the Maximum Score Estimator," *Journal of Econometrics*, 32, pp. 85-100.
- Manski, C., 1988. "Identification of Binary Response Models," *Journal of the American Statistical Association*, 83, pp. 729-738.
- Manski, C. and S. and Thompson, 1985. "MSCORE: A Program for Maximum Score Estimation of Linear Quantile Regressions from Binary Response Data, Mimeo, Department of Economics, University of Wisconsin, Madison.
- Marcus, A. and W. Greene, 1983. "The Determinants of Rating Assignment and Performance," Working Paper CRC528, Alexandria, VA, Center for Naval Analyses.
- McCullagh, P., 1977. "Analysis of Ordered Categorical Data," Ph.D. Thesis, University of London.
- McCullagh, P., 1979, "The Use of the Logistic Function in the Analysis of Ordinal Data," *Bulletin of the International Statistical Institute*, 48, pp. 21-33.
- McCullagh, P., 1980. "Regression Models for Ordered Data," *Journal of the Royal Statistical Society, Series B (Methodological)*, 42, pp. 109-142.
- McCullagh, P. and J. Nelder, 1983. *Generalized Linear Models*, London, Chapman and Hall.

- McCullagh, P. and J. Nelder, 1989. *Generalized Linear Models, 2<sup>nd</sup>. Ed.*, London, Chapman and Hall.
- McElvey, R. and W. Zavoina, 1975. "A Statistical Model for the Analysis of Ordered Level Dependent Variables," *Journal of Mathematical Sociology*, 4, pp. 103-120.
- McFadden, D., 1974. "Conditional Logit Analysis of Qualitative Choice Behavior," in P. Zarembka, ed., *Frontiers in Econometrics*, New York, Academic Press.
- McFadden, D. and K. Train, 2000. "Mixed MNL Models for Discrete Choice," *Journal of Applied Econometrics*, 15, pp. 447-470.
- McLachlan, G. and D. Peel, 2000. *Finite Mixture Models*, New York, John Wiley and Sons.
- McVicar, M. and J. McKee, 2002. "Part Time Work During Post-Compulsory Education and Examination Performance: Help or Hindrance" *Scottish Journal of Political Economy*, 49,4, pp. 393-406.
- Mehta, C. and N. Patel, 1995. "Exact Logistic Regression: Theory and Examples," *Statistics in Medicine*, 14, pp. 2143-2160.
- Mitchell, J. and M. Weale, 2007. "The Rationality and Reliability of Expectations Reported by British Households: Micro Evidence From the BHPS, National Institute of Economic and Social Research, Manuscript. (<http://www.niesr.ac.uk/pubs/dps/DP287.pdf>)
- Mora, J. and A. Moro-Egido, 2008. "On Specification Testing of Ordered Probit Models," *Journal of Econometrics*, 143, pp. 292-205.
- Mora, N., 2006., "Sovereign Credit Ratings: Guilty Beyond Reasonable Doubt?" *Journal of Banking and Finance*, 30, pp. 2041-2062.
- Mora, J. and A. Moro-Egido, 2008. "On Specification Testing of Ordered Discrete Choice Models," *Journal of Econometrics*, 143, pp. 191-205.
- Mullahy, J., 1986. "Specification and Testing of Some Modified Count Data Models," *Journal of Econometrics*, 33, pp. 341-365.
- Mullahy, J., 1997. "Heterogeneity, Excess Zeros and the Structure of Count Data Models," *Journal of Applied Econometrics*, 12, pp. 337-350.
- Müller, G. and C. Czado, 2005. "An Autoregressive Ordered Probit Model with Application to High Frequency Finance," *Journal of Computational and Graphical Statistics*, 14, pp. 320-338.
- Mundlak, Y., 1978. "On the Pooling of Time Series and Cross Section Data," *Econometrica*, 56, pp. 69-86.
- Munkin, M. and P. Trivedi, 2008. "Bayesian Analysis of the Ordered Probit Model with Endogenous Selection," *Journal of Econometrics*, 143, pp. 334-348.
- Murad, H., A. Fleischman, S. Sadetzki, O. Geyer and L. Freedman, 2003. "Small Samples and Ordered Logistic Regression: Does it Help to Collapse Categories of Outcome?" *The American Statistician*, 57, 3, pp. 155-160.
- Murphy, A., 1994. "Artificial Regression Based LM tests of Misspecification for Discrete Choice Models," *Economic and Social Review*, 26, pp. 69-74.
- Murphy, A., 1996. "Simple LM Tests of Misspecification for Ordered Logit Models," *Economics Letters*, 52, pp. 137-141.
- Murphy, K. and R. Topel, 2002. "Estimation and Inference in Two Stem Econometric Models," *Journal of Business and Economic Statistics*, 20, pp. 88-97 (reprinted from 2, pp. 370-379).
- Nagin, D. and J. Waldfogel, 1995. "The Effect of Criminality and Conviction on the Labor Market Status of Young British Offenders," *International Review of Law and Economics*, 15, 1, pp. 109-126.
- Nagler, J., 1994. "Scobit: An Alternative Estimator to Logit and Probit," *American Journal of Political Science*, 38, 1, pp. 230-255.
- Nakosteen, R. and M. Zimmer, 1980. "Migration and Income: The Question of Self Selection," *Southern Economic Journal*, 46, pp. 840-851.

- NDSHS, 2001. Computer Files for the Unit Record Data from the National Drug Strategy Household Surveys.
- Nelder, J. and R. Wedderburn, 1972. "Generalized Linear Models," *Journal of the Royal Statistical Society, A*, 135, pp. 370-384.
- Newey, W., 1985. "Maximum Likelihood Specification and Testing and Conditional Moment tests," *Econometrica*, 53, pp. 1047-1070.
- Nerlove, M. and J. Press, 1972. "Univariate and Multivariate Log-Linear and Logistic Models," Santa Monica, CA, RAND – R1306-EDA/HIS.
- Norton, E. and C. Ai, 2003. "Interaction Terms in Logit and Probit Models." *Economics Letters*, 80, 1, pp. 123–129.
- Olsen, R. 1978. "A Note on the Uniqueness of the Maximum Likelihood Estimator in the Tobit Model," *Econometrica*, 46, pp. 37-44.
- Olssen, U., 1979. "Maximum Likelihood Estimation of the Polychoric Correlation Coefficient," *Psychometrika*, 44, 4, pp. 443-460.
- Olsson, U., 1980. "Measuring of Correlation in Ordered Two Way Contingency Tables," *Journal of Marketing Research*, 17, 3, pp. 391-394
- Orme, C., 1990. "The small-sample performance of the information-matrix test," *Journal of Econometrics* 46, pp. 309–331
- Pagan, A. and F. Vella, 1989. "Diagnostic Tests for Models Based on Individual Data: A Survey," *Journal of Applied Econometrics*, 4, Supplement, pp. S29-S59.
- Paternoster, R. and R. Brame, 1998. "The Structural Similarity of Processes Generating Criminal and Analogous Behaviors." *Criminology* 36, pp. :633-670.
- Pearson, K., 1914. *Tables for Statisticians and Biometricians: Part I*, Cambridge, Cambridge University Press.
- Plackett, R., 1974. *The Analysis of Categorical Data*, London, Griffin.
- Plackett, R., 1981. *The Analysis of Categorical Data*, 2<sup>nd</sup>. Ed., London, Griffin
- Pratt, J., 1981. "Concavity of the Log Likelihood," *Journal of the American Statistical Association*, 76, pp. 103-116.
- Pregibon, D., 1984. "Book Review: P. McCullagh and J. A. Nelder, *Generalized Linear Models*," *The Annals of Statistics*, 12, 4, pp. 1589-1596.
- Prescott, E. and M. Visscher, 1977. "Sequential Location among Firms with Foresight," *Bell Journal of Economics*, 8, pp. 378–893.
- Pudney, S. and M. Shields, 2000. "Gender, Race, Pay and Promotion in the British Nursing Profession: Estimation of a Generalized Ordered Probit Model," *Journal of Applied Econometrics*, 15, pp. 367-399.
- Quednau, H., 1988. "An Extended Threshold Model for Analyzing Ordered Categorical Data," *Biometrical Journal*, 31, pp. 781-793.
- Rabe-Hesketh, S., A. Skrondal, A. and A. Pickles, 2005. "Maximum Likelihood Estimation of Limited and Discrete Dependent Variable Models with Nested Random Effects," *Journal of Econometrics*, 128, pp. 301-323.
- Rasch, G., 1960. "Probabilistic Models for Some Intelligence and Attainment Tests," Copenhagen, Danish Institute for Educational Research.
- Raudenbusch, S. and A. Bryk, 2002. *Hierarchical Linear Models: Applications and Data Analysis Methods*, Sage Publications, Thousand Oaks, CA.
- Ridder, G., 1990. "The Non-parametric Identification of Generalized Accelerated Failure-Time Models," *Review of Economic Studies*, 57 pp. 167–181.
- Riphahn, R., A. Wambach and A. Million, 2003. "Incentive Effects on the Demand for Health Care: A Bivariate Panel Count Data Estimation," *Journal of Applied Econometrics*, 18, 4, pp. 387-405.

- Ronning, G., 1990. "The Informational Content of Responses from Business Surveys," In J. Florens, M. Ivaldi, J. Laffont and F. Laisney, eds., *Microeconometrics: Surveys and Applications*, Oxford, Blackwell.
- Ronning, G. and M. Kukuk, 1996. "Efficient Estimation of Ordered Probit Models," *Journal of the American Statistical Association*, 91, 435, pp. 1120-1129.
- Sajaia, Z., 2008. "Maximum Likelihood Estimation of a Bivariate Ordered Probit Model: Implementation and Monte Carlo Simulations," *The Stata Journal*, 4, 2, pp. 1-18.
- Sanko, S., H. Maesoba, D. Dissanayake, T. Yamamoto, and T. Morikawa, 2004. "Inter-temporal and Inter-Regional Analysis of Household Cars and Motorcycles Ownership Behaviours in Asian Big Cities," Graduate School of Environmental Studies, Nagoya University, manuscript. ([http://cost355.inrets.fr/IMG/doc/SAKURA\\_Yamamoto\\_.doc](http://cost355.inrets.fr/IMG/doc/SAKURA_Yamamoto_.doc))
- SAS Institute, 2008. *SAS/STAT User's Guide. Version 9.2*, Cary NC, SAS Institute.
- SAS, 2008. *SAS User's Guide*, Cary NC, SAS.
- Scott, D. and K. Axhausen, 2006. "Household Mobility Tool Ownership, Modeling Interactions Between Cars and Season Tickets," *Transportation*, 33, 4, pp. 311-328.
- Scott, D. and Kanaroglou, P., 2001. "An Activity-Episode Generation Model that Captures Interactions Between Household Heads: Development and Empirical Analysis," *Transportation Research B: Methodological*, 36, 10, pp. 875-896
- Scotti, C., 2006. "A Bivariate Model of Fed and ECB Main Policy Rates," Board of Governors, International Finance Discussion paper 875
- Severini, T. and J. Staniswalis, 1994. "Quasi-Likelihood Estimation in Semiparametric Models," *Journal of the American Statistical Association*, 89, pp. 501-511.
- Shaked, A. and J. Sutton, 1982. "Relaxing Price Competition through Product Differentiation," *Review of Economic Studies*, 49, pp. 3-13.
- Simonoff, J., 2003. *Analyzing Categorical Data*, New York, Springer.
- Smith, R., 1989. "On the Use of Distributional Misspecification Checks in Limited Dependent Variables," *Economic Journal*, 99, pp. 178-192.
- Snell, E., 1964. "A Scaling Procedure for Ordered Categorical Data," *Biometrics*, 20, pp. 592-607.
- Stata, 2008. *Stata, Version 8.0*, College Station TX, Stata Corp.
- Stewart, M., 1983. "On Least Squares Estimation When the Dependent Variable Is Grouped," *Review of Economic Studies*, 50, pp. 141-149.
- Stewart, M., 2003. "Semi-nonparametric Estimation of Extended Ordered Probit Models," Department of Economics, University of Warwick, Manuscript.
- Stewart, M., 2005. "A Comparison of Semiparametric Estimators for the Ordered Response Model," *Computational Statistics and Data Analysis*, 49, pp. 555-573.
- Stutzer, A. and R. Lalive, 2001. "The Role of Social Work Norms in Job Searching and Subjective Well-Being," IZA Discussion Paper Number 300, Institute for the Study of Labor.
- Tanner, J. and W. Wong, 1987. "The Calculation of Posterior Distributions by Data Augmentation," *Journal of the American Statistical Association*, 82, pp. 528-549.
- Tattersfield, F., C. Gimingham and H. Morris, 1925. "Studies on Contact Insecticides. Part III. A Quantitative Examination of the Insecticidal Action of Chlor-, Nitro- and Hydroxyl Derivatives of Benzene and Naphthalene," *Annals of Applied Biology*, 12, p. 218.
- Tauchen, G., 1985. "Diagnostic Testing and Evaluation of Maximum Likelihood Models," *Journal of Econometrics*, 30, pp. 415-443.
- Terza, J., 1985. "Ordered Probit: A Generalization," *Communications in Statistics – A. Theory and Methods*, 14, pp. 1-11.
- Terza, J., 1987. "Estimating Linear Models with Ordinal Qualitative Regressors," *Journal of Econometrics*, 34, pp. 275-291.

- Terza, J., 1998. "Estimating Count Data Models with Endogenous Switching and Endogenous Treatment Effects," *Journal of Econometrics*, 84, pp. 129-154.
- Theil, H., 1969. "A Multinomial Extension of the Linear Logit Model," *International Economic Review*, 10, pp. 251-259.
- Theil, H., 1970. "On the Estimation of Relationships Involving Qualitative Variables," *American Journal of Sociology*, 76, pp. 103-154.
- Theil, H., 1971. *Principles of Econometrics*, New York, John Wiley and Sons.
- Tobin, J., 1958. "Estimation of Relationships for Limited Dependent Variables," *Econometrica*, 26, pp. 24-36.
- Tomoyuki, F. and F. Akira, 2006. "A Quantitative Analysis on Tourists' Consumer Satisfaction Via the Bayesian Ordered Probit Model," *Journal of the City Planning Institute of Japan*, 41, pp. 2-10. (In Japanese)
- Train, K., 2003. *Discrete Choice Methods with Simulation*, Cambridge, Cambridge University Press.
- Trost, R. and L. Lee, 1978. "Technical Training and Earnings: A Polychotomous Choice Model with Selectivity," *Review of Economics and Statistics*, 66, 1, pp. 151-156.
- Tsay, R., 2002. *Analysis of Financial Time Series*, New York, John Wiley and Sons.
- Tsay, R., 2005. *Analysis of Financial Time Series, 2<sup>nd</sup> Ed.*, New York, John Wiley and Sons.
- Tutz, G., 1989. "Compound Regression Models for Categorical Ordinal Data," *Biometrical Journal*, 31, pp. 259-272.
- Tutz, G., 1990. "Sequential Item Response Models with an Ordered Response," *British Journal of Mathematical and Statistical Psychology*, 43, pp. 39-55.
- Tutz, G., 1991. "Sequential Models in Ordered Regression," *Computational Statistics and Data Analysis*, 11, pp. 275-295.
- TSP, 2005. *TSP Version 5.0*, TSP International, Palo Alto, CA.
- UCLA/ATS, 2008. Academic Technology services  
([http://www.ats.ucla.edu/stat/mult\\_pkg/faq/general/Psuedo\\_RSquareds.htm](http://www.ats.ucla.edu/stat/mult_pkg/faq/general/Psuedo_RSquareds.htm))
- Uebersax, J., 1999. "Probit Latent Class Analysis with Dichotomous or Ordered Category Measures: Conditional Independence/Dependence Models," *Applied Psychological Measurement*, 23, pp. 283-297.
- Verbeek, M., 1990. "On the Estimation of a Fixed Effects Model with Selectivity Bias," *Economics Letters*, 34, pp. 267-270.
- Verbeek, M. and T. Nijman, 1992. "Testing for Selectivity Bias in Panel Data Models," *International Economic Review*, 33, pp. 681-703.
- Vuong, Q., 1989. "Likelihood Ratio Tests for Model Selection and Nonnested Hypotheses," *Econometrica*, 57, 2, pp. 307-333.
- Walker, S. and D. Duncan, 1967. "Estimation of the Probability of an Event As a Function of Several Independent Variables," *Biometrika*, 54, pp. 167-179.
- Wedderburn, A., 1974. "Quasi-Likelihood Functions, Generalized Linear Models, and the Gauss-Newton Method," *Biometrika*, 61, pp. 439-447.
- Weiss, A., 1993. "A Bivariate Ordered Probit Model with Truncation: Helmet Use and Motorcycle Injuries," *Applied Statistics*, 42, pp. 487-99.
- Weiss, A., 1997. "Specification tests in Ordered Logit and Probit Models," *Econometric Reviews*, 16, 4, pp. 361-391.
- White, H., 1980. "A Heteroscedasticity-Consistent Covariance Matrix Estimator and a Direct Test for Heteroscedasticity," *Econometrica*, 48, pp. 87-838.
- Williams, R., 2006. "Generalized Ordered Logit/ Partial Proportional Odds Models for Ordinal Dependent Variables," Department of Sociology University of Notre Dame, Notre Dame, IN, *The Stata Journal*, 6, 1, pp. 58-82.
- Winkelmann, R., 2005. "Subjective Well-being and the Family: Results from an Ordered Probit Model with Multiple Random Effects" *Empirical Economics*, 30, 3, pp 749-761,

- Winkelmann, L. and R. Winkelmann, 1998. "Why are the Unemployed So Unhappy? Evidence from Panel Data," *Economica*, 65, pp. 1-15.
- Winship, C. and R. Mare, 1984. "Regression Models with Ordered Variables," *American Sociological Review*, 49, 512-525.
- Wooldridge, J., 2002. "Inverse Probability Weighted M Estimators for Sample Stratification, Attrition and Stratification," *Portuguese Economic Journal*, 1, pp. 117-139.
- Wu, D., 1973. "Alternative Tests of Independence Between Stochastic Regressors and Disturbances," *Econometrica*, 41, pp. 733-750.
- Wynand, P. and B. van Praag, 1981. "The Demand for Deductibles in Private Health Insurance: A Probit Model with Sample Selection," *Journal of Econometrics*, 17, pp. 229-252.
- Zabel, J., 1992. "Estimating Fixed and Random Effects Models with Selectivity," *Economics Letters*, 40, pp. 269-272.
- Zhang, J., 2007. "Ordered Probit Modeling of User Perceptions of Protected Left-Turn Signals," *Journal of Transportation Engineering*, 133, 3, pp. 205-214.
- Zhang, Y, F. Liang and Y. Yuanchang, 2007. "Crash Injury Severity Analysis Using a Bayesian Ordered Probit Model," Transportation Research Board, Annual Meeting, Paper Number 07-2335.
- Zigante, V., 2007. "Ever Rising Expectations – The Determinants of Subjective Welfare in Croatia," School of Economics and Management, Lund University, Masters Thesis ([www.essays.se/about/Ordered+Probit+Model/](http://www.essays.se/about/Ordered+Probit+Model/)).