



Technical Report 113

A New Generalized Heterogeneous Data Model (GHDM) to Jointly Model Mixed Types of Dependent Variables

Chandra R. Bhat
Center for Transportation Research

September 2015

Data-Supported Transportation Operations & Planning Center (D-STOP)

A Tier 1 USDOT University Transportation Center at The University of Texas at Austin



**CENTER FOR
TRANSPORTATION
RESEARCH**



**Wireless Networking &
Communications Group**

D-STOP is a collaborative initiative by researchers at the Center for Transportation Research and the Wireless Networking and Communications Group at The University of Texas at Austin.

DISCLAIMER

The contents of this report reflect the views of the authors, who are responsible for the facts and the accuracy of the information presented herein. This document is disseminated under the sponsorship of the U.S. Department of Transportation's University Transportation Centers Program, in the interest of information exchange. The U.S. Government assumes no liability for the contents or use thereof.

Technical Report Documentation Page

1. Report No. D-STOP/2016/113		2. Government Accession No.		3. Recipient's Catalog No.	
4. Title and Subtitle A New Generalized Heterogeneous Data Model (GHDM) to Jointly Model Mixed Types of Dependent Variables				5. Report Date September 2015	
				6. Performing Organization Code	
7. Author(s) Chandra R. Bhat				8. Performing Organization Report No. Report 113	
9. Performing Organization Name and Address Data-Supported Transportation Operations & Planning Center (D-STOP) The University of Texas at Austin 1616 Guadalupe Street, Suite 4.202 Austin, Texas 78701				10. Work Unit No. (TRAVIS)	
				11. Contract or Grant No. DTRT13-G-UTC58	
12. Sponsoring Agency Name and Address Data-Supported Transportation Operations & Planning Center (D-STOP) The University of Texas at Austin 1616 Guadalupe Street, Suite 4.202 Austin, Texas 78701				13. Type of Report and Period Covered	
				14. Sponsoring Agency Code	
15. Supplementary Notes Supported by a grant from the U.S. Department of Transportation, University Transportation Centers Program.					
16. Abstract This paper formulates a generalized heterogeneous data model (GHDM) that jointly handles mixed types of dependent variables—including multiple nominal outcomes, multiple ordinal variables, and multiple count variables, as well as multiple continuous variables—by representing the covariance relationships among them through a reduced number of latent factors. Sufficiency conditions for identification of the GHDM parameters are presented. The maximum approximate composite marginal likelihood (MACML) method is proposed to estimate this jointly mixed model system. This estimation method provides computational time advantages since the dimensionality of integration in the likelihood function is independent of the number of latent factors. The study undertakes a simulation experiment within the virtual context of integrating residential location choice and travel behavior to evaluate the ability of the MACML approach to recover parameters. The simulation results show that the MACML approach effectively recovers underlying parameters, and also that ignoring the multi-dimensional nature of the relationship among mixed types of dependent variables can lead not only to inconsistent parameter estimation, but also have important implications for policy analysis.					
17. Key Words Latent factors, big data analytics, high dimensional data, MACML estimation approach, mixed dependent variables, structural equations models, integrated land use-transportation modeling, factor analysis			18. Distribution Statement No restrictions. This document is available to the public through NTIS (http://www.ntis.gov): National Technical Information Service 5285 Port Royal Road Springfield, Virginia 22161		
19. Security Classif.(of this report) Unclassified		20. Security Classif.(of this page) Unclassified		21. No. of Pages 58	22. Price

Disclaimer

The contents of this report reflect the views of the authors, who are responsible for the facts and the accuracy of the information presented herein. Mention of trade names or commercial products does not constitute endorsement or recommendation for use.

Acknowledgements

This research was partially supported by the U.S. Department of Transportation through the Data-Supported Transportation Operations and Planning (D-STOP) Tier 1 University Transportation Center. The author would also like to acknowledge support from a Humboldt Research Award from the Alexander von Humboldt Foundation, Germany. Finally, the author is grateful to Lisa Macias for her help in formatting this document, to Subodh Dubey and Xuemei Fu for help with the simulation runs, and two anonymous referees who provided useful comments on an earlier version of the paper.

Table of Contents

Chapter 1. Introduction.....	1
Chapter 2. The GHDM Formulation	5
2.1 Latent Variable SEM.....	7
2.2 Latent Variable Measurement Equation Model Components	7
Chapter 3. The Model System Identification and Estimation	13
3.1 Model Identification	13
3.2 Model Estimation	17
3.3 The Joint Mixed Model System and the MACML Estimation Approach	18
3.4 Positive Definiteness	22
Chapter 4. Simulation Experiment.....	23
4.1 Experimental Design	24
4.2 The Structural Equation System.....	24
4.3 The Measurement Equation System.....	27
4.4 Data Generation Process	35
4.5 Performance Evaluation	36
4.6 Simulation Results.....	37
4.7 Procedure for Treatment Effects Based on Residential Choice	46
References	48

List of Illustrations

Figure 1: Diagrammatic representation of the structural equation.....	26
Figure 2a: Diagrammatic representation of the measurement equation for the non-nominal variables.....	28
Figure 2b: Diagrammatic representation of the measurement equation for the nominal variables.....	29
Figure 2c: Endogeneous effects.....	30

List of Tables

Table 1: Matrix Notation, Description, and Dimension.....	5
Table 2: Simulation Results for the 1000-Observations Case with 200 Datasets.....	40
Table 3: Simulation Results for the 2000-Observations Case with 200 Datasets.....	42
Table 4: Simulation Results for the 3000-Observations Case with 200 Datasets.....	44

Chapter 1. Introduction

The joint modeling of data with mixed types of dependent variables (including ordered-response or ordinal variables, unordered-response or nominal variables, count variables, and continuous variables) is of interest in several fields, including biology, developmental toxicology, finance, economics, epidemiology, social science, and transportation (see a good synthesis of applications in De Leon and Chough, 2013). For instance, in the clinical biology field, alternative treatments for a specific condition are assessed based on binary, ordered, and continuous indicators of the treatment's after-effects; this approach has been used to assess the effectiveness of depression medication in reducing the occurrence, frequency, and intensity of depression (such as in Gueorguieva and Sanacora, 2006). In the health field, in addition to binary, count, and continuous variables related to the occurrence, frequency, and intensity, respectively, of specific health problems, it is not uncommon to obtain ordinal information on quality of life outcomes/perceptions. In the toxicology field, the focus is on regulating the use of chemical and pharmaceutical drugs (Sutton et al., 2000). Typically, varying quantities of a drug are administered to mice; the effects on their offspring are studied in terms of combinations of discrete outcomes (such as the presence of congenital deformations) and continuous outcomes (such as birth weight). In the transportation field, households that are not auto-oriented are likely to locate in transit- and pedestrian-friendly neighborhoods that are characterized by mixed and high-density land use; pedestrian-oriented design in such communities may also further structurally reduce motorized vehicle miles of travel. If that is the case, then it is likely that the choices of residential location (nominal variable), vehicle ownership (count), and vehicle miles of travel (continuous) are being made jointly as a bundle (see, for example, Bhat et al., 2014a).

The interest in mixed model systems has been spurred particularly by the recent availability of high-dimensional heterogeneous data with complex dependence structures, thanks to technology that allows the collection and archival of voluminous amounts of data ("big data"). Unlike standard correlated linear data that can be analyzed using traditional multivariate linear regression models, the presence of non-commensurate outcomes creates difficulty because of the absence of a convenient multivariate distribution to jointly (and directly) represent the relationship between discrete and continuous outcomes. Several approaches have been developed to handle such situations. The first and simplest is, of course, to simply ignore the dependence and estimate separate models. However, such an approach is inefficient in estimating covariate effects for each outcome because it fails to borrow information on other outcomes, and is limiting in its ability to answer intrinsically multivariate questions such as the effect of a covariate on a multidimensional outcome (Teixeira-Pinto and Harezlak, 2013). Besides, joint analysis of mixed outcomes obviates the need for multiple tests and facilitates global tests, offering superior power in testing and better control of type I error rates (De Leon and Zhu, 2008). But, more importantly, if some endogenous outcomes are used to explain other endogenous outcomes (such as examining the effect of density of residence on auto-ownership model), and if the outcomes are not modeled jointly to recognize the presence of unobserved exogenous variable effects, the result is inconsistent estimation of the effects of one endogenous outcome on another (see Bhat and Guo, 2007, and

Mokhtarian and Cao, 2008). A second common approach to joint mixed outcome modeling originates in the general location model (GLOM), which assumes an arbitrary marginal distribution for the discrete outcomes and a conditional (on the discrete component) normality assumption for the continuous outcomes (De Leon and Chough, 2013). However, the GLOM is not suitable for ordinal outcome variables and does not accommodate dependence between nominal and ordinal outcomes. A third “reverse-factorization” approach is to employ a latent variable representation for binary/ordinal outcomes, and assume a multivariate normal (MVN) distribution for the continuous outcomes and the latent variables underlying the binary/ordinal outcomes. Then, the joint distribution is derived using a marginal distribution of the continuous outcomes and the conditional distribution of the latent variables (given the continuous variables) underlying the binary/ordinal outcomes. This approach is referred to as the conditional grouped continuous model (CGCM) by De Leon and Chough (2013). However, this approach cannot be directly extended to the case of nominal outcomes, since nominal outcomes do not arise from the partitioning of a single latent variable using thresholds (as is the case for binary/ordinal outcomes). So, De Leon and Carriere (2007) and De Leon et al. (2011) proposed an extended factorization approach, which they label as the general mixed data model (GMDM), to accommodate nominal outcomes. They use a GLOM for the joint distribution of the nominal and continuous outcomes, and a CGCM for the joint distribution of the ordinal and continuous outcomes. Specifically, the GMDM uses a multinomial distribution for the marginal distribution of the possible multidimensional discrete states obtained from the combinatorics of a set of nominal outcomes, followed by a conditional MVN distribution for the latent variables (underlying the ordinal outcomes) and the continuous outcomes. The mean vector for this latter conditional MVN distribution is specified to be a function of the multidimensional discrete state, engendering an association between the nominal discrete outcomes and the ordinal/continuous outcomes. However, the covariance matrix of the conditional MVN distribution is constant across the nominal discrete states. A further problem with the GMDM is that the number of multidimensional discrete states explodes as the number of nominal discrete outcomes increases, and as the number of elemental categories within each nominal discrete outcome increases. Besides, the GMDM (like the GLOM) resorts to a factorization approach in which an artificial hierarchy is implicitly assumed. In this hierarchy, the multidimensional discrete outcomes are intermediate responses and the ordinal/continuous outcomes are the ultimate responses (see Wu et al., 2013).

Independent from the work discussed above, a fourth approach originates in the economics and transportation fields, wherein mixed models with nominal outcomes are based on latent variable representations of nominal outcomes. Surprisingly, such studies are rarely mentioned in papers in the statistical field that deal with mixed outcomes. The studies in this strand may be viewed as extensions of the CGCM approach to the case of nominal outcomes, except that each nominal outcome is represented by a series of latent variables. An early example of such a multivariate model may be found in Keane (1992), who considered one nominal variable and one continuous variable. However, only relatively recently has this methodology been extended to include mixed nominal, binary, ordinal, count, and continuous variables (for example, see Paleti et. al., 2013 and Bhat et al., 2014a). The resulting mixed models may be viewed as an alternative to the GMDM, and have the advantage that all outcomes are tied based on their latent or observed

continuous variable representations (rather than using different types of linkages for different types of outcomes, as in the GMDM). Further, these models treat the mixed outcomes symmetrically rather than imposing any form of hierarchy. The models typically assume an MVN distribution over the entire set of latent and observed continuous variables characterizing the many types of outcomes. A variant of this methodology uses a Gaussian copula function to tie the latent and observed continuous variables if the variables have different marginal distributions, though this approach has been confined to scenarios without a nominal outcome (see, for example, Wu et al., 2013). Another variant introduces random error terms linearly in the latent and observed continuous variable equations associated with the discrete outcomes and continuous outcomes, respectively. The underlying continuous variables are considered to be independent, conditional on these random error terms. Then, if these random error terms are common or correlated, the result is an association structure among the mixed outcomes. Such a specification falls under the label of a multivariate generalized linear latent and mixed model (GLLMM), and is particularly helpful when considering clustering effects (due to multiple observations from the same person or due to spatial dependency) in addition to correlation across mixed outcomes (see, for example, Faes et al., 2009 and Bhat et al., 2014a). An extension of this approach that accommodates clustering as well as an association structure among mixed outcomes (that is, mixed outcomes are independent, conditional on appropriately specified latent variables) is referred to as the item response theory (IRT) model in the literature (see Bartholomew et al., 2011 and Feddag, 2013). However, again, these GLLMM and IRT models have been predominantly used for cases with no nominal variables, though similar approaches can be used to generate dependence between a nominal variable and other kinds of variables too (see, for example, Bhat and Guo, 2007 and Pinjari et al., 2008).

A fifth approach, originating from the social sciences, implicitly generates dependence among mixed outcomes by writing the latent and observed continuous variables as a function of unobserved psychological constructs. These relationships are characterized as measurement equations, in that the psychological constructs are manifested in the larger combination of mixed outcomes. The constructs themselves are related to exogenous variables and may be correlated with one another in a structural relationship. In this approach, the unobserved psychological constructs serve as latent factors that provide a structure to the dependence among the many mixed indicator variables. Seen from this perspective, the approach can also be viewed as a parsimonious attempt to explain the covariance relationship among a large set of mixed outcomes through a much smaller number of unobservable latent factors. Sometimes referred to as factor analysis, the approach represents a powerful dimension-reduction technique to analyze high-dimensional heterogeneous outcome data by representing the covariance relationship among the data through a smaller number of unobservable latent factors. An entire field of structural equations modeling (SEM) has been developed around this psychological construct-based dependence modeling, originating in some of the early works of Jöreskog (1977). However, the SEM field has focused almost exclusively on non-nominal outcome analysis (see Gates et al., 2011 and Hoshino and Bentler, 2013). Indeed, traditional SEM software (such as LISREL, MPLUS, and EQS) is either not capable of handling nominal indicators or at least are not readily suited to handle nominal indicators (see Temme et al., 2008). But when this approach is extended to include a nominal indicator, it

essentially takes the form of an integrated choice and latent variable (ICLV) model (Ben-Akiva et al., 2002, and Bolduc et al., 2005). Also, while traditional SEM techniques typically adopt normally distributed latent factors along with normally distributed measurement error terms (leading to probit models in the presence of binary/ordered outcomes), ICLV models tend to use normally distributed latent factors mixed with logistically distributed errors in the measurement equations for ordinal variables and type-1 extreme value errors in the nominal outcome utility functions (leading to a probability expression that involves a multivariate integral over the product of logit-type probabilities for the outcomes). In both the SEM and ICLV cases, the standard estimation methodology is the method of maximum likelihood estimation. When there are many binary/ordered-response outcomes (indicators) and/or a nominal variable, the integrals in the overall probability expression are computed using simulation techniques. As indicated by Hoshino and Bentler (2013), this can “be difficult to impossible when the model is complex or the number of variables is large.” This is particularly the case with the traditional mixture formulation of ICLV models in general, and particularly when there are several latent factors (see Daziano and Bolduc, 2013).

Recently, Bhat and Dubey (2014) proposed a different way of formulating ICLV models, in which they use a SEM-like probit approach while also accommodating a single nominal variable. Essentially, this approach combines the power and parsimony of the dimension-reduction factor analysis structure of SEMs (as just discussed above) with the extended CGCM approach that uses a symmetric, latent continuous variable representation for all non-continuous outcomes (as in Paleti et al., 2013 and Bhat et al., 2014a). In this paper, we generalize Bhat and Dubey’s approach to the case of multiple nominal outcomes, multiple ordinal variables, multiple count variables, and multiple continuous variables. The resulting model, which we label simply as the generalized heterogeneous data model (GHDM), is general enough to accommodate other models in the literature as special cases. Straightforward extensions of the model are available to accommodate longitudinal and spatial clustering, though we focus on the non-clustered mixed outcome model in the current paper. We propose the estimation of the GHDM using Bhat’s maximum approximate composite marginal likelihood (MACML) inference approach. In particular, in our approach, the dimensionality of integration in the composite marginal likelihood (CML) function that needs to be maximized to obtain a consistent estimator (under standard regularity conditions) for the GHDM parameters is independent of the number of latent factors and easily accommodates general covariance structures for the structural equation and for the utilities of the discrete alternatives for each nominal outcome. Further, the use of the analytic approximation in the MACML approach to evaluate the multivariate cumulative normal distribution (MVNCD) function in the CML function simplifies the estimation procedure even further so that the proposed MACML procedure requires the maximization of a function that has no more than bivariate normal cumulative distribution functions to be evaluated.

Chapter 2. The GHDM Formulation

There are two components to the model: (1) the latent variable SEM, and (2) the latent variable measurement equation model. These components are discussed in turn below. In the following presentation, for ease in notation, we will consider a cross-sectional model. As appropriate and convenient, we will suppress the index q for decision-makers ($q=1,2,\dots,Q$) in parts of the presentation, and assume that all error terms are independent and identically distributed across decision-makers. Table 1 summarizes all matrix notations and corresponding matrix dimensions used below in the GHDM formulation.

Table 1: Matrix Notation, Description, and Dimension

Symbol		Represents...	
L		Number of latent variables	
\tilde{D}		Total number of exogenous variables in the structural equation system	
H		Number of continuous outcomes in the measurement equation system	
N		Number of ordinal outcomes in the measurement equation system	
C		Number of count outcomes in the measurement equation system	
A		Total number of exogenous and endogenous variables in the measurement equation system	
\tilde{G}		Total number of alternatives across all nominal variables in the choice model component of the measurement equation system	
Equation	Notation	Represents...	Dimension
Structural Equation (Equation 12 in text)	\mathbf{z}^*	Vector of latent variables	$L \times 1$
	\mathbf{a}	Matrix of exogenous variable loadings on \mathbf{z}^*	$L \times \tilde{D}$
	\mathbf{w}	Vector of exogenous variables affecting \mathbf{z}^*	$\tilde{D} \times 1$
	$\boldsymbol{\eta}$	Vector of errors in structural equation	$L \times 1$
	$\boldsymbol{\Gamma}$	Correlation matrix of error vector $\boldsymbol{\eta}$ in latent variable structural equation	$L \times L$

Equation	Notation	Represents...	Dimension
Measurement Equation (Equation 13 in text; \tilde{y} originates from Equation 7)	\tilde{y}	Vector of observed latent measurement equation dependent variables	$(H + N + C) \times 1$
	$\tilde{\gamma}$	Matrix of coefficients representing the effect of exogenous and possible endogenous variables	$(H + N + C) \times A$
	\tilde{d}	Matrix of coefficients representing the effect of latent variables on measurement equation dependent variables	$(H + N + C) \times L$
Measurement Equation	$\tilde{\varepsilon}$	Vector of errors in measurement equation	$(H + N + C) \times 1$
	$\tilde{\Sigma}$	Covariance matrix of $\tilde{\varepsilon}$ (assumed diagonal for identification)	$(H + N + C) \times (H + N + C)$
	$\tilde{\gamma}$	Matrix of coefficients representing the effect of exogenous and possible endogenous variables on the count outcome	$C \times A$
Choice Model (Equation 14 in text; see text above Equation 10 for β and ϑ)	U	Vector of alternative utilities	$\vec{G} \times 1$
	b	Matrix of exogenous and possible endogenous variable effects on U	$\vec{G} \times A$
	x	Vector of exogenous variables in choice model	$A \times 1$
	β	Matrix of coefficients capturing effects of latent variables and their interactions with exogenous variables	$\left[\sum_{i_g=1}^{I_g} N^{gi_g} \right] \times L$ (Please see text for construction)
	ϑ	Matrix of variables interacting with latent variables	$\vec{G} \times \left[\sum_{i_g=1}^{I_g} N^{gi_g} \right]$ (Please see text for construction)
	ζ	Utility error vector	$\vec{G} \times 1$
	Λ	Covariance matrix of ζ	$\vec{G} \times \vec{G}$

2.1 Latent Variable SEM

Let l be an index for latent variables ($l=1,2,\dots,L$). Consider the latent variable z_l^* and write it as a linear function of covariates:

$$z_l^* = \boldsymbol{\alpha}_l' \boldsymbol{w} + \eta_l, \quad (1)$$

where \boldsymbol{w} is a $(\tilde{D} \times 1)$ vector of observed covariates (excluding a constant), $\boldsymbol{\alpha}_l$ is a corresponding $(\tilde{D} \times 1)$ vector of coefficients, and η_l is a random error term assumed to be standard normally distributed for identification purposes (see Stapleton, 1978).¹ Next, define the $(L \times \tilde{D})$ matrix $\boldsymbol{\alpha} = (\boldsymbol{\alpha}_1, \boldsymbol{\alpha}_2, \dots, \boldsymbol{\alpha}_L)'$, and the $(L \times 1)$ vectors $\boldsymbol{z}^* = (z_1^*, z_2^*, \dots, z_L^*)'$ and $\boldsymbol{\eta} = (\eta_1, \eta_2, \eta_3, \dots, \eta_L)'$. Unlike much of the earlier research in ICLV modeling, we allow an MVN correlation structure for $\boldsymbol{\eta}$ to accommodate interactions among the unobserved latent variables: $\boldsymbol{\eta} \sim MVN_L[\mathbf{0}_L, \boldsymbol{\Gamma}]$, where $\mathbf{0}_L$ is an $(L \times 1)$ column vector of zeros, and $\boldsymbol{\Gamma}$ is $(L \times L)$ correlation matrix. In matrix form, we may write Equation (1) as:

$$\boldsymbol{z}^* = \boldsymbol{\alpha} \boldsymbol{w} + \boldsymbol{\eta}. \quad (2)$$

It is not uncommon in the SEM literature to have latent variables affecting each other in the SEM. However, it may also not be easy to justify *a priori* inter-relationships between unobserved variables, and so we prefer a general covariance structure for the latent variables as in Equation (2). However, in some cases, it may indeed be appropriate to allow inter-relationships between the latent variables. Section 3.1 discusses the identification considerations in this case. Note also that our model formulation and estimation technique are readily applicable to this case of inter-related latent constructs too as long as the identification considerations in Section 3.1 are met.

2.2 Latent Variable Measurement Equation Model Components

We will consider a combination of continuous, ordinal, count, and nominal outcomes (indicators) of the underlying latent variable vector \boldsymbol{z}^* . However, these outcomes may be a function of a set of exogenous variables too.

Let there be H continuous outcomes (y_1, y_2, \dots, y_H) with an associated index h ($h = 1, 2, \dots, H$). Let $y_h = \boldsymbol{\gamma}_h' \boldsymbol{x} + \boldsymbol{d}_h' \boldsymbol{z}^* + \varepsilon_h$ in the usual linear regression fashion, where \boldsymbol{x} is an $(A \times 1)$ vector of exogenous variables (including a constant) as well as possibly the observed values of other endogenous continuous variables, other endogenous ordinal variables, other endogenous count variables, and other endogenous nominal variables (introduced as dummy variables). $\boldsymbol{\gamma}_h$ is a corresponding compatible coefficient vector.²

¹ The reason for excluding the constant in the covariate vector \boldsymbol{w} will become clear in Section 3.

² In joint limited-dependent variable systems in which one or more dependent variables are not observed on a continuous scale, such as the joint system considered in the current paper that has discrete dependent and count variables (which we will more generally refer to as limited-dependent variables), the structural effects of one limited-dependent variable on another can only be in a single direction. That is, it is not possible to have correlated

\mathbf{d}_h is an $(L \times 1)$ vector of latent variable loadings on the h th continuous outcome, and ε_h is a normally distributed measurement error term. Stack the H continuous outcomes into an $(H \times 1)$ vector \mathbf{y} , and the H error terms into another $(H \times 1)$ vector $\boldsymbol{\varepsilon} = (\varepsilon_1, \varepsilon_2, \dots, \varepsilon_H)'$. Also, let $\boldsymbol{\Sigma}$ be the covariance matrix of $\boldsymbol{\varepsilon}$, which is restricted to be diagonal. This helps identification because there is already an unobserved latent variable vector \mathbf{z}^* that serves as a vehicle to generate covariance between the outcome variables (as we discuss in the next section). Define the $(H \times A)$ matrix $\boldsymbol{\gamma} = (\gamma_1, \gamma_2, \dots, \gamma_H)'$ and the $(H \times L)$ matrix of latent variable loadings $\mathbf{d} = (\mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_H)'$. Then, one may write, in matrix form, the following measurement equation for the continuous outcomes:

$$\mathbf{y} = \boldsymbol{\gamma}\mathbf{x} + \mathbf{d}\mathbf{z}^* + \boldsymbol{\varepsilon}. \quad (3)$$

Next, consider N ordinal outcomes (indicator variables) for the individual, and let n be the index for the ordinal outcomes ($n = 1, 2, \dots, N$). Also, let J_n be the number of categories for the n^{th} ordinal outcome ($J_n \geq 2$) and let the corresponding index be j_n ($j_n = 1, 2, \dots, J_n$). Let \tilde{y}_n^* be the latent underlying variable whose horizontal partitioning leads to the observed outcome for the n^{th} ordinal variable. Assume that the individual under consideration chooses the a_n^{th} ordinal category. Then, in the usual ordered response formulation, for the individual, we may write:

$$\tilde{y}_n^* = \tilde{\boldsymbol{\gamma}}_n' \mathbf{x} + \tilde{\mathbf{d}}_n' \mathbf{z}^* + \tilde{\varepsilon}_n, \text{ and } \tilde{\boldsymbol{\psi}}_{n, a_n-1} < \tilde{y}_n^* < \tilde{\boldsymbol{\psi}}_{n, a_n}, \quad (4)$$

where \mathbf{x} is a vector of exogenous and possibly endogenous variables as defined earlier, $\tilde{\boldsymbol{\gamma}}_n$ is a corresponding vector of coefficients to be estimated, $\tilde{\mathbf{d}}_n$ is an $(L \times 1)$ vector of latent variable loadings on the n^{th} continuous outcome, the $\tilde{\boldsymbol{\psi}}$ terms represent thresholds, and $\tilde{\varepsilon}_n$ is the standard normal random error for the n^{th} ordinal outcome. For each ordinal outcome, $\tilde{\boldsymbol{\psi}}_{n,0} < \tilde{\boldsymbol{\psi}}_{n,1} < \tilde{\boldsymbol{\psi}}_{n,2} \dots < \tilde{\boldsymbol{\psi}}_{n, J_n-1} < \tilde{\boldsymbol{\psi}}_{n, J_n}$; $\tilde{\boldsymbol{\psi}}_{n,0} = -\infty$, $\tilde{\boldsymbol{\psi}}_{n,1} = 0$, and $\tilde{\boldsymbol{\psi}}_{n, J_n} = +\infty$. For later use, let $\tilde{\boldsymbol{\psi}}_n = (\tilde{\boldsymbol{\psi}}_{n,2}, \tilde{\boldsymbol{\psi}}_{n,3}, \dots, \tilde{\boldsymbol{\psi}}_{n, J_n-1})'$ and $\tilde{\boldsymbol{\psi}} = (\tilde{\boldsymbol{\psi}}_1', \tilde{\boldsymbol{\psi}}_2', \dots, \tilde{\boldsymbol{\psi}}_N)'$. Stack the N underlying continuous variables \tilde{y}_n^* into an $(N \times 1)$ vector $\tilde{\mathbf{y}}^*$, and the N error terms $\tilde{\varepsilon}_n$ into another $(N \times 1)$ vector $\tilde{\boldsymbol{\varepsilon}}$. Define $\tilde{\boldsymbol{\gamma}} = (\tilde{\boldsymbol{\gamma}}_1, \tilde{\boldsymbol{\gamma}}_2, \dots, \tilde{\boldsymbol{\gamma}}_N)'$ [$(N \times A)$ matrix] and $\tilde{\mathbf{d}} = (\tilde{\mathbf{d}}_1, \tilde{\mathbf{d}}_2, \dots, \tilde{\mathbf{d}}_N)$ [$(N \times L)$ matrix], and let \mathbf{IDEN}_N be the identity matrix of dimension N representing the correlation matrix of $\tilde{\boldsymbol{\varepsilon}}$ (so, $\tilde{\boldsymbol{\varepsilon}} \sim MVN_N(\mathbf{0}_N, \mathbf{IDEN}_N)$); again, this is for identification purposes, given the presence of the unobserved \mathbf{z}^* vector to generate covariance. Finally, stack the lower thresholds for the decision-maker $\tilde{\boldsymbol{\psi}}_{n, a_n-1}$ ($n = 1, 2, \dots, N$) into an $(N \times 1)$

unobserved effects underlying the propensities determining two limited-dependent variables, as well as have the observed limited-dependent variables themselves structurally affect each other in a bi-directional fashion. This creates a logical inconsistency problem (see Maddala, 1983, page 119 for a good discussion). It is critical to note that, regardless of which directionality of structural effects among the endogenous variables is specified (or even if no relationships are specified), the system is a joint bundled system because of the correlation in unobserved factors impacting the underlying propensities.

vector $\tilde{\psi}_{low}$ and the upper thresholds $\tilde{\psi}_{n,a_n}$ ($n = 1, 2, \dots, N$) into another vector $\tilde{\psi}_{up}$. Then, in matrix form, the measurement equation for the ordinal outcomes (indicators) for the decision-maker may be written as:

$$\tilde{y}^* = \tilde{\gamma}x + \tilde{d}z^* + \tilde{\varepsilon}, \quad \tilde{\psi}_{low} < \tilde{y}^* < \tilde{\psi}_{up}. \quad (5)$$

Let there be C count variables for a household, and let c be the index for the count variables ($c = 1, 2, \dots, C$). Let the count index be k_c ($k_c = 0, 1, 2, \dots, \infty$) and let r_c be the actual observed count value for the household. Then, following the recasting of a count model in a generalized ordered-response probit formulation (see Castro, Paleti, and Bhat, or CPB, 2012 and Bhat *et al.*, 2014b), a generalized version of the negative binomial count model may be written as:

$$\tilde{y}_c^* = \tilde{d}_c' z^* + \tilde{\varepsilon}_c, \quad \tilde{\psi}_{c,r_c-1} < \tilde{y}_c^* < \tilde{\psi}_{c,r_c}, \quad (6)$$

$$\tilde{\psi}_{c,r_c} = \Phi^{-1} \left[\frac{(1-v_c)^{\theta_c}}{\Gamma(\theta_c)} \sum_{t=0}^{r_c} \left(\frac{\Gamma(\theta_c + t)}{t!} (v_c)^t \right) \right] + \varphi_{c,r_c}, \quad v_c = \frac{\lambda_c}{\lambda_c + \theta_c}, \quad \text{and } \lambda_c = e^{\tilde{\gamma}_c' x}. \quad (7)$$

In the above equation, \tilde{y}_c^* is a latent continuous stochastic propensity variable associated with the count variable c that maps into the observed count r_c through the $\tilde{\psi}_c$ vector (which is a vertically stacked column vector of thresholds $(\tilde{\psi}_{c,-1}, \tilde{\psi}_{c,0}, \tilde{\psi}_{c,1}, \tilde{\psi}_{c,2}, \dots)'$). \tilde{d}_c is an $(L \times 1)$ vector of latent variable loadings on the c^{th} count outcome, and $\tilde{\varepsilon}_c$ is a standard normal random error term. $\tilde{\gamma}_c$ is a column vector corresponding to the vector x . Φ^{-1} in the threshold function of Equation (7) is the inverse function of the univariate cumulative standard normal. θ_c is a parameter that provides flexibility to the count formulation, and is related to the dispersion parameter in a traditional negative binomial model ($\theta_c > 0 \forall c$). $\Gamma(\theta_c)$ is the traditional gamma function; $\Gamma(\theta_c) = \int_{\tilde{t}=0}^{\infty} \tilde{t}^{\theta_c-1} e^{-\tilde{t}} d\tilde{t}$. The

threshold terms in the $\tilde{\psi}_c$ vector satisfy the ordering condition (*i.e.*, $\tilde{\psi}_{c,-1} < \tilde{\psi}_{c,0} < \tilde{\psi}_{c,1} < \tilde{\psi}_{c,2} \dots < \infty \forall c$) as long as $\varphi_{c,-1} < \varphi_{c,0} < \varphi_{c,1} < \varphi_{c,2} \dots < \infty$. The presence of the φ_c terms in the thresholds provides substantial flexibility to accommodate high or low probability masses for specific count outcomes without the need for cumbersome traditional treatments using zero-inflated or related mechanisms in multi-dimensional model systems (see Castro *et al.*, 2011 for a detailed discussion). For identification, we set $\varphi_{c,-1} = -\infty$ and $\varphi_{c,0} = 0$ for all count variables c . In addition, we identify a count value e_c^* ($e_c^* \in \{0, 1, 2, \dots\}$) above which φ_{c,k_c} ($k_c \in \{1, 2, \dots\}$) is held fixed at φ_{c,e_c^*} ; that is, $\varphi_{c,k_c} = \varphi_{c,e_c^*}$ if $k_c > e_c^*$, where the value of e_c^* can be based on empirical testing. Doing so is the key to allowing the count model to predict beyond the range available in the estimation sample. For later use, let $\varphi_c = (\varphi_{c,1}, \varphi_{c,2}, \dots, \varphi_{c,e_c^*})'$ ($e_c^* \times 1$

vector) (assuming $e_c^* > 0$), $\varphi = (\varphi'_1, \varphi'_2, \dots, \varphi'_C)' \left[\left(\sum_c e_c^* \right) \times 1 \text{ vector} \right]$, and

$\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_C)'$ [$C \times 1$ vector]. Also, stack the C latent variables \tilde{y}_c^* into a $(C \times 1)$ vector $\tilde{\mathbf{y}}^*$, and the C error terms $\tilde{\boldsymbol{\varepsilon}}_c$ into another $(C \times 1)$ vector $\tilde{\boldsymbol{\varepsilon}}$. Let $\tilde{\boldsymbol{\varepsilon}} \sim MVN_C(\mathbf{0}_C, \mathbf{IDEN}_C)$ from identification considerations, and stack the lower thresholds of the individual $\tilde{\psi}_{c,r_{c-1}} (c = 1, 2, \dots, C)$ into a $(C \times 1)$ vector $\tilde{\boldsymbol{\psi}}_{low}$, and the upper thresholds $\tilde{\psi}_{c,r_c} (c = 1, 2, \dots, C)$ into another $(C \times 1)$ vector $\tilde{\boldsymbol{\psi}}_{up}$. Define $\tilde{\boldsymbol{\gamma}} = (\tilde{\gamma}_1, \tilde{\gamma}_2, \dots, \tilde{\gamma}_C)'$ [$(C \times A)$ matrix] and $\tilde{\mathbf{d}} = (\tilde{d}_1, \tilde{d}_2, \dots, \tilde{d}_C)'$ [$(C \times L)$ matrix]. With these definitions, the latent propensity underlying the count outcomes may be written in matrix form as:

$$\tilde{\mathbf{y}}^* = \tilde{\mathbf{d}}\mathbf{z}^* + \tilde{\boldsymbol{\varepsilon}}, \quad \tilde{\boldsymbol{\psi}}_{low} < \tilde{\mathbf{y}}^* < \tilde{\boldsymbol{\psi}}_{up}. \quad (8)$$

Note also that the interpretation of the generalized ordered-response recasting is that consumers have a latent “long-term” propensity \tilde{y}_c^* associated with the demand for each product/service represented by the count c , which is a linear function of the latent variable vector \mathbf{z}^* (see CPB for a discussion of the interpretation of the generalized ordered-response recasting of count models). Such a specification enables covariance across the count outcomes (through the propensity variables \tilde{y}_c^*) and between the count outcomes and other mixed outcomes. On the other hand, there may be some specific consumer contexts and characteristics (embedded in \mathbf{x}) that may dictate how the long-term propensity is manifested in a count demand at any given *instant of time*. Our implicit assumption is that the latent variable vector \mathbf{z}^* affects the “long-term” latent demand propensity \tilde{y}_c^* , but does not play a role in the instantaneous translation of propensity to actual manifested count demand. This allows us to easily incorporate count outcomes within a mixed outcome model, and estimate the resulting model using Bhat (2011) MACML approach. Similarly, an implicit assumption in Equation (8) is that the factors/constraints that are responsible for the instantaneous translation of propensity to manifested count demand (that is, the elements of the \mathbf{x} vector) do not affect the “long-term” demand propensity, though this is being imposed purely for parsimony purposes. Relaxing this assumption does not complicate the model system or the estimation process in any way.

Finally, let there be G nominal (unordered-response) variables for an individual, and let g be the index for the nominal variables ($g = 1, 2, 3, \dots, G$). Also, let I_g be the number of alternatives corresponding to the g^{th} nominal variable ($I_g \geq 3$) and let i_g be the corresponding index ($i_g = 1, 2, 3, \dots, I_g$). Consider the g^{th} nominal variable and assume that the individual under consideration chooses the alternative m_g . Also, assume the usual random utility structure for each alternative i_g .

$$U_{gi_g} = \mathbf{b}'_{gi_g} \mathbf{x} + \boldsymbol{\vartheta}'_{gi_g} (\boldsymbol{\beta}_{gi_g} \mathbf{z}^*) + \zeta_{gi_g}, \quad (9)$$

where \mathbf{x} is as defined earlier, \mathbf{b}_{gi_g} is an $(A \times 1)$ column vector of corresponding coefficients, and ζ_{gi_g} is a normal error term. $\boldsymbol{\beta}_{gi_g}$ is an $(N_{gi_g} \times L)$ -matrix of variables interacting with latent variables to influence the utility of alternative i_g , and $\boldsymbol{\vartheta}_{gi_g}$ is an

$(N_{g^{i_g}} \times 1)$ -column vector of coefficients capturing the effects of latent variables and their interaction effects with other exogenous variables. If each of the latent variables impacts the utility of the alternatives for each nominal variable purely through a constant shift in the utility function, $\beta_{g^{i_g}}$ will be an identity matrix of size L , and each element of $\vartheta_{g^{i_g}}$ will capture the effect of a latent variable on the constant specific to alternative i_g of nominal variable g . Let $\zeta_g = (\zeta_{g1}, \zeta_{g2}, \dots, \zeta_{g^{I_g}})'$ ($I_g \times 1$ vector), and $\zeta_g \sim MVN_{I_g}(\mathbf{0}, \Lambda_g)$. Taking the difference with respect to the first alternative, the only estimable elements are found in the covariance matrix $\tilde{\Lambda}_g$ of the error differences, $\tilde{\zeta}_g = (\tilde{\zeta}_{g2}, \tilde{\zeta}_{g3}, \dots, \tilde{\zeta}_{g^{I_g}})$ (where $\tilde{\zeta}_{gi} = \zeta_{gi} - \zeta_{g1}$, $i \neq 1$).³ Further, the variance term at the top left diagonal of $\tilde{\Lambda}_g$ ($g = 1, 2, \dots, G$) is set to 1 to account for scale invariance. Λ_g is constructed from $\tilde{\Lambda}_g$ by adding a row on top and a column to the left. All elements of this additional row and column are filled with values of zero. In addition, the usual identification restriction is imposed such that one of the alternatives serves as the base when introducing alternative-specific constants and variables that do not vary across alternatives (that is, whenever an element of \mathbf{x} is individual-specific and not alternative-specific, the corresponding element in $\mathbf{b}_{g^{i_g}}$ is set to zero for at least one alternative i_g). To proceed, define $U_g = (U_{g1}, U_{g2}, \dots, U_{g^{I_g}})'$ ($I_g \times 1$ vector), $\mathbf{b}_g = (\mathbf{b}_{g1}, \mathbf{b}_{g2}, \mathbf{b}_{g3}, \dots, \mathbf{b}_{g^{I_g}})'$ ($I_g \times A$ matrix), and $\beta_g = (\beta'_{g1}, \beta'_{g2}, \dots, \beta'_{g^{I_g}})'$ $\left(\sum_{i_g=1}^{I_g} N_{g^{i_g}} \times L \right)$ matrix. Also, define the $\left(I_g \times \sum_{i_g=1}^{I_g} N_{g^{i_g}} \right)$ matrix ϑ_g , which is initially filled with all zero values. Then, position the $(1 \times N_{g1})$ row vector ϑ'_{g1} in the first row to occupy columns 1 to N_{g1} , position the $(1 \times N_{g2})$ row vector ϑ'_{g2} in the second row to occupy columns $N_{g1} + 1$ to $N_{g1} + N_{g2}$, and so on until the $(1 \times N_{g^{I_g}})$ row vector $\vartheta'_{g^{I_g}}$ is appropriately positioned. Further, define $\varpi_g = (\vartheta_g \beta_g)$ ($I_g \times L$ matrix), $\tilde{G} = \sum_{g=1}^G I_g$, $\tilde{G} = \sum_{g=1}^G (I_g - 1)$, $\mathbf{U} = (\mathbf{U}'_1, \mathbf{U}'_2, \dots, \mathbf{U}'_G)'$ ($\tilde{G} \times 1$ vector), $\boldsymbol{\zeta} = (\boldsymbol{\zeta}'_1, \boldsymbol{\zeta}'_2, \dots, \boldsymbol{\zeta}'_G)'$ ($\tilde{G} \times 1$ vector), $\mathbf{b} = (\mathbf{b}'_1, \mathbf{b}'_2, \dots, \mathbf{b}'_G)'$ ($\tilde{G} \times A$ matrix), $\boldsymbol{\varpi} = (\boldsymbol{\varpi}'_1, \boldsymbol{\varpi}'_2, \dots, \boldsymbol{\varpi}'_G)'$ ($\tilde{G} \times L$ matrix), and $\boldsymbol{\vartheta} = \text{Vech}(\boldsymbol{\vartheta}_1, \boldsymbol{\vartheta}_2, \dots, \boldsymbol{\vartheta}_G)$ (that is, $\boldsymbol{\vartheta}$ is a column vector that includes all elements of the matrices $\boldsymbol{\vartheta}_1, \boldsymbol{\vartheta}_2, \dots, \boldsymbol{\vartheta}_G$). Then, in matrix form, we may write Equation (9) as:

$$U = \mathbf{b}\mathbf{x} + \boldsymbol{\varpi}\mathbf{z}^* + \boldsymbol{\zeta}, \quad (10)$$

where $\boldsymbol{\zeta} \sim MVN_{\tilde{G}}(\mathbf{0}_{\tilde{G}}, \Lambda)$. As earlier, to ensure identification, we specify Λ as follows:

³ Also, in multinomial probit models, identification is tenuous when only individual-specific covariates are used in the vector \mathbf{x} (see Keane, 1992 and Munkin and Trivedi, 2008). In particular, exclusion restrictions are needed in the form of at least one individual characteristic being excluded from each alternative's utility in addition to being excluded from a base alternative (but appearing in some other utilities). But these exclusion restrictions are not needed when there are alternative-specific variables.

$$\Lambda = \begin{bmatrix} \Lambda_1 & \mathbf{0} & \mathbf{0} & \mathbf{0} \dots \mathbf{0} \\ \mathbf{0} & \Lambda_2 & \mathbf{0} & \mathbf{0} \dots \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \Lambda_3 & \mathbf{0} \dots \mathbf{0} \\ \vdots & \vdots & \vdots & \vdots \dots \vdots \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \dots \Lambda_G \end{bmatrix} (\vec{G} \times \vec{G} \text{ matrix}). \quad (11)$$

In the general case, this allows the estimation of $\sum_{g=1}^G \left(\frac{I_g^* (I_g - 1)}{2} - 1 \right)$ terms across all the G nominal variables, as originating from $\left(\frac{I_g^* (I_g - 1)}{2} - 1 \right)$ terms embedded in each $\check{\Lambda}_g$ matrix; ($g=1,2,\dots,G$).

Chapter 3. The Model System Identification and Estimation

Let $E = (H + N + C)$. Define $\tilde{\mathbf{y}} = \left(\mathbf{y}', [\tilde{\mathbf{y}}^*]', [\tilde{\mathbf{y}}^*]' \right)' [E \times 1 \text{ vector}]$, $\tilde{\boldsymbol{\gamma}} = (\boldsymbol{\gamma}', \tilde{\boldsymbol{\gamma}}', \mathbf{0}_{AC})' [E \times A \text{ matrix}]$, $\tilde{\mathbf{d}} = (\mathbf{d}', \tilde{\mathbf{d}}', \tilde{\mathbf{d}}')' [E \times L \text{ matrix}]$, and $\tilde{\boldsymbol{\varepsilon}} = (\boldsymbol{\varepsilon}', \tilde{\boldsymbol{\varepsilon}}', \tilde{\boldsymbol{\varepsilon}}')' (E \times 1 \text{ vector})$, where $\mathbf{0}_{AC}$ is a matrix of zeros of dimension $A \times C$. Let $\boldsymbol{\delta}$ be the collection of parameters to be estimated:

$\boldsymbol{\delta} = [\text{Vech}(\boldsymbol{\alpha}), \text{Vech}(\boldsymbol{\Sigma}), \text{Vech}(\tilde{\boldsymbol{\gamma}}), \text{Vech}(\tilde{\mathbf{d}}), \text{Vech}(\tilde{\boldsymbol{\gamma}}), \boldsymbol{\varphi}, \boldsymbol{\theta}, \text{Vech}(\mathbf{b}), \boldsymbol{\vartheta}, \text{Vech}(\boldsymbol{\Lambda})]$, where the operator "Vech(.)" vectorizes all the non-zero elements of the matrix/vector on which it operates. We will assume that the error vectors $\boldsymbol{\tau}$, $\boldsymbol{\varepsilon}$, $\boldsymbol{\zeta}$, and $\boldsymbol{\varsigma}$ are independent of each other. While this assumption is not strictly necessary (and can be relaxed in a very straightforward manner within the estimation framework of our model system as long as the resulting model is identified), the assumption aids in developing general sufficiency conditions for identification of parameters in a mixed model when the latent variable vector \mathbf{z}^* already provides a mechanism to generate covariance among the mixed outcomes.

With the matrix definitions above, the continuous components of the model system may be written compactly as:

$$\mathbf{z}^* = \boldsymbol{\alpha}\mathbf{w} + \boldsymbol{\eta}, \quad (12)$$

$$\tilde{\mathbf{y}} = \tilde{\boldsymbol{\gamma}}\mathbf{x} + \tilde{\mathbf{d}}\mathbf{z}^* + \tilde{\boldsymbol{\varepsilon}}, \text{ with } \text{Var}(\tilde{\boldsymbol{\varepsilon}}) = \tilde{\boldsymbol{\Sigma}} = \begin{bmatrix} \boldsymbol{\Sigma} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \text{IDEN}_N & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \text{IDEN}_C \end{bmatrix} (E \times E \text{ matrix}), \quad (13)$$

$$\mathbf{U} = \mathbf{b}\mathbf{x} + \boldsymbol{\vartheta}\mathbf{z}^* + \boldsymbol{\varsigma}. \quad (14)$$

To develop the reduced form equations, replace the right side of Equation (12) for \mathbf{z}^* in Equations (13) and (14) to obtain the following system:

$$\tilde{\mathbf{y}} = \tilde{\boldsymbol{\gamma}}\mathbf{x} + \tilde{\mathbf{d}}\mathbf{z}^* + \tilde{\boldsymbol{\varepsilon}} = \tilde{\boldsymbol{\gamma}}\mathbf{x} + \tilde{\mathbf{d}}(\boldsymbol{\alpha}\mathbf{w} + \boldsymbol{\eta}) + \tilde{\boldsymbol{\varepsilon}} = \tilde{\boldsymbol{\gamma}}\mathbf{x} + \tilde{\mathbf{d}}\boldsymbol{\alpha}\mathbf{w} + \tilde{\mathbf{d}}\boldsymbol{\eta} + \tilde{\boldsymbol{\varepsilon}}, \quad (15)$$

$$\mathbf{U} = \mathbf{b}\mathbf{x} + \boldsymbol{\vartheta}\mathbf{z}^* + \boldsymbol{\varsigma} = \mathbf{b}\mathbf{x} + \boldsymbol{\vartheta}(\boldsymbol{\alpha}\mathbf{w} + \boldsymbol{\eta}) + \boldsymbol{\varsigma} = \mathbf{b}\mathbf{x} + \boldsymbol{\vartheta}\boldsymbol{\alpha}\mathbf{w} + \boldsymbol{\vartheta}\boldsymbol{\eta} + \boldsymbol{\varsigma}.$$

Now, consider the $[(E + \bar{G}) \times 1]$ vector $\mathbf{y}\mathbf{U} = [\tilde{\mathbf{y}}', \mathbf{U}']'$. Define

$$\mathbf{B} = \begin{bmatrix} \mathbf{B}_1 \\ \mathbf{B}_2 \end{bmatrix} = \begin{bmatrix} \tilde{\boldsymbol{\gamma}}\mathbf{x} + \tilde{\mathbf{d}}\boldsymbol{\alpha}\mathbf{w} \\ \mathbf{b}\mathbf{x} + \boldsymbol{\vartheta}\boldsymbol{\alpha}\mathbf{w} \end{bmatrix} \text{ and } \boldsymbol{\Omega} = \begin{bmatrix} \boldsymbol{\Omega}_1 & \boldsymbol{\Omega}'_{12} \\ \boldsymbol{\Omega}_{12} & \boldsymbol{\Omega}_2 \end{bmatrix} = \begin{bmatrix} \tilde{\mathbf{d}}\boldsymbol{\Gamma}\tilde{\mathbf{d}}' + \tilde{\boldsymbol{\Sigma}} & \tilde{\mathbf{d}}\boldsymbol{\Gamma}\boldsymbol{\vartheta}' \\ \boldsymbol{\vartheta}\boldsymbol{\Gamma}\tilde{\mathbf{d}}' & \boldsymbol{\vartheta}\boldsymbol{\Gamma}\boldsymbol{\vartheta}' + \boldsymbol{\Lambda} \end{bmatrix}. \quad (16)$$

Then $\mathbf{y}\mathbf{U} \sim \text{MVN}_{E+\bar{G}}(\mathbf{B}, \boldsymbol{\Omega})$.

3.1 Model Identification

The question of identification relates to whether all the elements of $\boldsymbol{\delta}$ are estimable from the elements of \mathbf{B} and $\boldsymbol{\Omega}$ (that is, from $\mathbf{B}_1, \mathbf{B}_2, \boldsymbol{\Omega}_1, \boldsymbol{\Omega}_2, \boldsymbol{\Omega}_{12}$). A simple approach would be to develop easy-to-apply sufficiency conditions for identification (even if they may lead to over-identification and may be more restrictive than needed). A starting point for

this is O'Brien (1994) and Reilly and O'Brien (1996), who develop sufficiency conditions for multiple-indicator multiple-cause (MIMIC) models, and whose discussion is applicable to SEM-based models with no nominal variables. Conforming with the setup of earlier MIMIC models, we will assume in our mixed model that the number of measurement equations with non-nominal variables exceeds the number of latent factors (this will typically be the case, and indeed forms the backbone of modeling a high-dimensional mixed data model through a lower dimensional factor analytic structure). That is, we will assume that $E > L$. We will also assume the presence of more than one latent variable, as is quite common in MIMIC models ($L > 1$). However, in contrast to earlier MIMIC studies, we allow nominal dependent variables, allow the variable vector \mathbf{x} to appear in the measurement equations, and allow the observed endogenous variables to be inter-related. In this situation, we can develop sufficiency conditions in five steps as follows.

- (1) First, if the exogenous covariates do not appear in the measurement equations, one can use O'Brien's (1994) exposition for MIMIC models with no nominal variables (that is, for the sub-model given by Equations (12) and (13) with $\vec{\gamma} = \mathbf{0}$) to show that the elements of this sub-model (*i.e.*, $\boldsymbol{\alpha}$, $\boldsymbol{\Gamma}$, $\vec{\mathbf{d}}$, and $\vec{\boldsymbol{\Sigma}}$) are all identifiable as long as:
 - (a) $\boldsymbol{\Gamma}$ in the structural equation is specified to be a correlation matrix, with each latent variable correlated with at least one other latent variable,
 - (b) diagonality is maintained across the elements of the error term vector $\vec{\boldsymbol{\varepsilon}}$ (that is, $\vec{\boldsymbol{\Sigma}}$ is diagonal),
 - (c) for each latent variable, there are at least two non-nominal outcome variables that load only on that latent variable and no other latent variable (that is, there are at least two factor complexity one outcome variables for each latent variable) (see Reilly and O'Brien, 1996).

The first two of these conditions have already been imposed in the development of our mixed model formulation (the specification that the covariance matrices of $\vec{\boldsymbol{\varepsilon}}$ and $\vec{\boldsymbol{\varepsilon}}$ are identity matrices is a result of imposing diagonality combined with a scaling restriction for ordinal and count outcomes). Intuitively speaking, the reason for the first condition is that only the entire diagonal terms of the covariance matrix elements of the non-nominal outcomes in the reduced form Equation (16) are identified: that is, only the diagonal terms of $\vec{\mathbf{d}}\boldsymbol{\Gamma}\vec{\mathbf{d}}' + \vec{\boldsymbol{\Sigma}}$ as a whole are identified. Thus, as long as there are diagonal variance terms to be estimated in $\vec{\boldsymbol{\Sigma}}$ (subject to identification considerations as discussed in the previous section), it immediately implies that diagonal terms in $\boldsymbol{\Gamma}$ cannot be identified solely from the estimated diagonal entries of $\vec{\mathbf{d}}\boldsymbol{\Gamma}\vec{\mathbf{d}}' + \vec{\boldsymbol{\Sigma}}$ (and so the diagonal terms of $\boldsymbol{\Gamma}$ are normalized to one, leading to the correlation matrix for $\boldsymbol{\Gamma}$). The second sufficiency condition is related to the off-diagonal terms in $\vec{\mathbf{d}}\boldsymbol{\Gamma}\vec{\mathbf{d}}' + \vec{\boldsymbol{\Sigma}}$. If we allow $\vec{\boldsymbol{\Sigma}}$ to have a full set of off-diagonal elements, it immediately implies that the off-diagonal elements of $\boldsymbol{\Gamma}$ are not identified. That is, one can ignore the correlations (the off-diagonals) in $\boldsymbol{\Gamma}$ (set these to zero), and estimate all the off-diagonal elements of $\vec{\boldsymbol{\Sigma}}$. The problem with this though is that it leads to an explosion in the number of covariance parameters to be estimated. Thus, if there are a total of six ordinal/count/ continuous dependent variables, the number of off-diagonal parameters in a fully specified $\vec{\boldsymbol{\Sigma}}$ matrix is 15. With 10 ordinal/count/continuous dependent variables, the number of off-diagonal parameters in a fully specified $\vec{\boldsymbol{\Sigma}}$ matrix is 45. On the other hand, the value of the latent factor approach arises through the effective dimensionality reduction that accrues from having all off-diagonal

elements in a full covariance matrix for Γ , but no off-diagonal elements in $\tilde{\Sigma}$. Doing so essentially places a factor-analytic structure for the covariances among the ordinal/count/continuous dependent variables, with this structure being represented by the off-diagonal elements of $\vec{d}\Gamma\vec{d}'$. Thus, if there are three latent variables that underlie the 10 ordinal/count/continuous variables, there are effectively only three off-diagonal elements in Γ to be estimated to characterize the 45 off-diagonal entries for the covariance elements among the ordinal/count/continuous dependent variables. Of course, one can keep all the off-diagonal elements of Γ and introduce additional off-diagonal elements very selectively in $\tilde{\Sigma}$ to still achieve theoretical identification, but this can become *ad hoc* and will require examination for each specific case to ensure identification. Overall, keeping $\tilde{\Sigma}$ diagonal and allowing Γ to have all off-diagonal elements ensures identification, while also being the vehicle to reduce high-dimensional problems through a factor-analytic structure. This increases econometric efficiency, and allows the estimation of high-dimensional models with the order of sample sizes typically available for model estimation. Note, however, that our estimation procedure itself is agnostic to the number of parameters to be estimated in terms of computational ability. The third condition can be imposed through the empirical specification based on theoretical/intuitive considerations. This condition, referred to as the two indicator rule (see, Bollen, 1989, page 244), essentially allows identification of the matrices \vec{d} , α , and covariance matrix Γ of the structural matrix errors.

- (2) Next, we consider the result from the first step, but now relax the constraint that $\vec{\gamma} = \mathbf{0}$, and allow some exogenous variables to influence the non-nominal variables. In this situation, there is an identification problem in Equation (13) if the same exogenous variable is allowed to have a direct impact through the \mathbf{x} vector as well as an indirect impact through a latent variable. That is, in general, it is not possible to disentangle the separate effects of the same variable through the direct $\vec{\gamma}$ effect and through the indirect \vec{d} effect. A sufficient identification condition is then to ensure that the element corresponding to the effect of each exogenous variable is zero in either the $\vec{\gamma}$ vector or the α vector (this is also the reason that we include a constant in the \mathbf{x} vector, but not in the \mathbf{w} vector). In other words, a sufficient condition for identification of the parameters in the structural equation and the measurement equations for non-nominal outcomes (that is, α , Γ , $\vec{\gamma}$, \vec{d} , and $\tilde{\Sigma}$) is:

- (a) the three conditions from the first step hold, plus
- (b) the condition holds that each element of $\vec{\gamma}$ in Equation (13) is either
 - (i) directly related to an exogenous variable without being a function of any latent variable that itself has the exogenous variable as a covariate in the structural equation, or
 - (ii) loaded onto latent variables, but then not directly related to any exogenous variable that itself impacts any of the latent variables on which the outcome variable loads.

That is, an exogenous variable, as a sufficiency condition for identification, should not impact an element of $\vec{\gamma}$ both directly and indirectly.

- (3) Third, we proceed to the choice model components. Following Bhat and Dubey (2014), we ignore the information available from the covariance matrix $\Omega_{12} = \sigma\Gamma\vec{d}'$. While one can effectively use this covariance matrix to identify parameters in specific situations, we develop a simpler (albeit more restrictive than needed) and general sufficiency condition for identification

of the measurement equation parameters corresponding to the nominal outcomes based only on the mean element of the utilities $\mathbf{B}_2 = \mathbf{b}\mathbf{x} + \boldsymbol{\omega}\mathbf{a}\mathbf{w}$ (but we retain a general covariance matrix $\boldsymbol{\Lambda}_g$ across alternative utilities for each nominal outcome g). Specifically, all the parameters in the nominal measurement equation part in Equation (14) (that is, elements of \mathbf{b} , the elements of $\boldsymbol{\vartheta}_g$ ($g=1,2,\dots,G$) embedded in $\boldsymbol{\omega}$, and $\boldsymbol{\Lambda}$) are estimable if all latent variables appear only as interactions and not as direct shifters of utility. In this case, there are effectively no common exogenous variables in the \mathbf{x} effect and the \mathbf{w} effect, and so identification of the elements of \mathbf{b}_g and $\boldsymbol{\vartheta}_g$ is immediate for each nominal variable g through estimation of the mean \mathbf{B}_2 . But identification becomes more challenging in the case when the latent variables appear by themselves in the choice models (with or without additional interaction effects of the latent variables). In this case, if an element of \mathbf{b}_{g_i} corresponding to a specific variable in the vector \mathbf{x} is non-zero, a sufficient condition for identification is that the utility of alternative i_g not depend on any latent variable that contains that specific variable as a covariate in the structural equation system. This is the most common way that identification has been achieved in most earlier ICLV studies. In fact, most ICLV studies do not even seem to discuss this identification issue. Alternatively, one may include common elements (including alternative-specific attributes in the utilities of the alternatives of nominal variables and those same variables in the structural model for latent variables that impact the utilities), but appropriate restrictions have to be imposed (for example, a latent variable may affect the utility of one of three alternatives for a nominal variable, and a covariate affecting that latent variable may also impact the utility of the same alternative but the coefficient on the covariate may be constrained to be the same as a covariate appearing in the utility of one of the other two alternatives). However, given the sheer number of such specific situations, we leave an in-depth study of identification issues in the context of the overlapping explanatory variables in the structural equation and in the utilities of nominal variables for a later date.

- (4) Fourth, as indicated in footnote 2, endogenous variable effects can be specified only in a single direction. In addition, when a continuous observed endogenous variable (say variable A) appears as a right side variable in the regression for another continuous observed endogenous variable, or as a right side variable in the latent regression underlying another count or ordinal endogenous variable, each latent variable appearing in the regression/latent regression for the other endogenous continuous/count/ordinal variable (say variable B) should have two factor complexity one outcome variables after excluding the equation for variable B. Essentially, this sufficiency condition ensures that part c of the first step continues to hold. This latter condition is not needed when a non-continuous observed endogenous variable appears as a right side variable in the regression of any other observed endogenous variable because of the non-linear nature of the relationship between the latent regressions and the observed non-continuous endogenous variables.
- (5) Finally, moving to the structural equation system, in this paper we use a reduced form system as shown in Equation (2). In this case, only the above four sufficiency conditions are needed for identification. However, as discussed under Equation (2), there may be instances when the analyst wants to allow direct inter-relationships between the latent constructs or variables. In this situation, identification is still possible if a recursive relationship is used so that some latent variables appear as right side variables in the equations for other latent variables in a recursive fashion. But one of two conditions for identification should hold even in this recursive case. The first is that the error terms of the latent variables in the structural form are uncorrelated (though, in reduced form each latent variable should be correlated with at least another latent variable; that is, one must ensure that each latent variable, excepting the first one in the

recursive structure, is directly related to at least one other upstream latent variable in this uncorrelated case for the sufficiency conditions discussed in the first four steps above to hold). Alternatively, a second condition that also allows identification is that there should be at least one exogenous variable in each upstream latent variable equation that does not appear in each downstream latent variable equation that has the upstream latent variable as an explanatory variable (please see the online supplement to this paper at http://www.caee.utexas.edu/prof/bhat/ABSTRACTS/GHDM/Online_supp_GHDM.pdf for a discussion of these identification conditions).

3.2 Model Estimation

To estimate the model, note that, under the utility maximization paradigm, $U_{g i_g} - U_{g m_g}$ must be less than zero for all $i_g \neq m_g$ corresponding to the g th nominal variable, since the individual chose alternative m_g . Let $u_{g i_g m_g} = U_{g i_g} - U_{g m_g}$ ($i_g \neq m_g$), and stack the latent utility differentials into a vector $\mathbf{u}_g = \left[\left(u_{g 1 m_g}, u_{g 2 m_g}, \dots, u_{g I_g m_g} \right)'; i_g \neq m_g \right]$. Also, define

$\mathbf{u} = \left(\left[\mathbf{u}_1 \right]', \left[\mathbf{u}_2 \right]', \dots, \left[\mathbf{u}_G \right]' \right)'$. We now need to develop the distribution of the vector

$\mathbf{y}\mathbf{u} = (\mathbf{y}', \mathbf{u}')'$ from that of $\mathbf{y}\mathbf{U} = [\mathbf{y}', \mathbf{U}']'$. To do so, define a matrix \mathbf{M} of size $[E + \tilde{G}] \times [E + \tilde{G}]$. Fill this matrix with values of zero. Then, insert an identity matrix of size E into the first E rows and E columns of the matrix \mathbf{M} . Next, consider the rows from $E + 1$ to $E + I_1 - 1$, and columns from $E + 1$ to $E + I_1$. These rows and columns correspond to the first nominal variable. Insert an identity matrix of size $(I_1 - 1)$ after supplementing with a column of '-1' values in the column corresponding to the chosen alternative. Next, rows $E + I_1$ through $E + I_1 + I_2 - 2$ and columns $E + I_1 + 1$ through $E + I_1 + I_2$ correspond to the second nominal variable. Again position an identity matrix of size $(I_2 - 1)$ after supplementing with a column of '-1' values in the column corresponding to the chosen alternative for the second nominal variable. Continue this procedure for all G nominal variables. With the matrix \mathbf{M} as defined, we can write $\mathbf{y}\mathbf{u} \sim MVN_{E+\tilde{G}}(\tilde{\mathbf{B}}, \tilde{\mathbf{\Omega}})$, where $\tilde{\mathbf{B}} = \mathbf{M}\mathbf{B}$ and $\tilde{\mathbf{\Omega}} = \mathbf{M}\mathbf{\Omega}\mathbf{M}'$. Next,

partition the vector $\tilde{\mathbf{B}}$ into components that correspond to the mean of the vectors \mathbf{y} (for the continuous variables), $\tilde{\mathbf{y}} = \left(\left[\tilde{\mathbf{y}}^* \right]', \left[\tilde{\mathbf{y}}^* \right]' \right)'$ $[(N + C) \times 1 \text{ vector}]$, (for the ordinal and count outcomes), and \mathbf{u} (for the nominal outcomes), and the matrix $\tilde{\mathbf{\Omega}}$ into the corresponding variances and covariances:

$$\tilde{\mathbf{B}} = \begin{bmatrix} \tilde{\mathbf{B}}_y \\ \tilde{\mathbf{B}}_{\tilde{\mathbf{y}}} \\ \tilde{\mathbf{B}}_u \end{bmatrix} (E + \tilde{G}) \times 1 \text{ vector and } \tilde{\mathbf{\Omega}} = \begin{bmatrix} \tilde{\mathbf{\Omega}}_{yy} & \tilde{\mathbf{\Omega}}_{y\tilde{\mathbf{y}}} & \tilde{\mathbf{\Omega}}_{yu} \\ \tilde{\mathbf{\Omega}}'_{y\tilde{\mathbf{y}}} & \tilde{\mathbf{\Omega}}_{\tilde{\mathbf{y}}\tilde{\mathbf{y}}} & \tilde{\mathbf{\Omega}}'_{y\tilde{\mathbf{y}}} \\ \tilde{\mathbf{\Omega}}'_{yu} & \tilde{\mathbf{\Omega}}'_{\tilde{\mathbf{y}}\tilde{\mathbf{y}}} & \tilde{\mathbf{\Omega}}_{uu} \end{bmatrix} (E + \tilde{G}) \times (E + \tilde{G}) \quad (17)$$

matrix.

Define $\tilde{\mathbf{u}} = (\tilde{\mathbf{y}}', \mathbf{u}')'$, so that $\mathbf{y}\mathbf{u} = (\mathbf{y}', \tilde{\mathbf{u}})'$. Re-partition $\tilde{\mathbf{B}}$ and $\tilde{\mathbf{\Omega}}$ in a different way such that:

$$\tilde{\mathbf{B}} = \begin{bmatrix} \tilde{\mathbf{B}}_y \\ \tilde{\mathbf{B}}_{\tilde{\mathbf{u}}} \end{bmatrix}, \text{ where } \tilde{\mathbf{B}}_{\tilde{\mathbf{u}}} = \begin{bmatrix} \tilde{\mathbf{B}}_{\tilde{\mathbf{y}}} \\ \tilde{\mathbf{B}}_{\mathbf{u}} \end{bmatrix} (N + C + \tilde{G}) \times 1 \text{ vector, and} \quad (18)$$

$$\tilde{\mathbf{\Omega}} = \begin{bmatrix} \tilde{\mathbf{\Omega}}_y & \tilde{\mathbf{\Omega}}_{y\tilde{\mathbf{u}}} \\ \tilde{\mathbf{\Omega}}'_{y\tilde{\mathbf{u}}} & \tilde{\mathbf{\Omega}}_{\tilde{\mathbf{u}}} \end{bmatrix}, \tilde{\mathbf{\Omega}}_{\tilde{\mathbf{u}}} = \begin{bmatrix} \tilde{\mathbf{\Omega}}_{\tilde{\mathbf{y}}} & \tilde{\mathbf{\Omega}}_{\tilde{\mathbf{y}}\mathbf{u}} \\ \tilde{\mathbf{\Omega}}'_{\tilde{\mathbf{y}}\mathbf{u}} & \tilde{\mathbf{\Omega}}_{\mathbf{u}} \end{bmatrix} (N + C + \tilde{G}) \times (N + C + \tilde{G}), \text{ and } \tilde{\mathbf{\Omega}}'_{y\tilde{\mathbf{u}}} = \begin{bmatrix} \tilde{\mathbf{\Omega}}_{y\tilde{\mathbf{y}}} \\ \tilde{\mathbf{\Omega}}_{y\mathbf{u}} \end{bmatrix}.$$

The conditional distribution of $\tilde{\mathbf{u}}$, given \mathbf{y} , is MVN with mean $\tilde{\mathbf{B}}_{\tilde{\mathbf{u}}} = \tilde{\mathbf{B}}_{\tilde{\mathbf{u}}} + \tilde{\mathbf{\Omega}}'_{y\tilde{\mathbf{u}}} \tilde{\mathbf{\Omega}}_y^{-1} (\mathbf{y} - \tilde{\mathbf{B}}_y)$ and variance $\tilde{\mathbf{\Omega}}_{\tilde{\mathbf{u}}} = \tilde{\mathbf{\Omega}}_{\tilde{\mathbf{u}}} - \tilde{\mathbf{\Omega}}'_{y\tilde{\mathbf{u}}} \tilde{\mathbf{\Omega}}_y^{-1} \tilde{\mathbf{\Omega}}_{y\tilde{\mathbf{u}}}$. Next, define threshold vectors as follows:

$\tilde{\boldsymbol{\psi}}_{low} = \left[\tilde{\boldsymbol{\psi}}'_{low}, \tilde{\boldsymbol{\psi}}'_{low}, (-\infty_{\tilde{G}})' \right]'$ ($[(N + C + \tilde{G}) \times 1]$ vector) and $\tilde{\boldsymbol{\psi}}_{up} = \left[\tilde{\boldsymbol{\psi}}'_{up}, \tilde{\boldsymbol{\psi}}'_{up}, (\mathbf{0}_{\tilde{G}})' \right]'$ ($[(N + C + \tilde{G}) \times 1]$ vector), where $-\infty_{\tilde{G}}$ is a $\tilde{G} \times 1$ -column vector of negative infinities, and $\mathbf{0}_{\tilde{G}}$ is another $\tilde{G} \times 1$ -column vector of zeros. Then the likelihood function may be written as:

$$\begin{aligned} L(\boldsymbol{\delta}) &= f_H(\mathbf{y} | \tilde{\mathbf{B}}_y, \tilde{\mathbf{\Omega}}_y) \times \Pr \left[\tilde{\boldsymbol{\psi}}_{low} \leq \tilde{\mathbf{u}} \leq \tilde{\boldsymbol{\psi}}_{up} \right], \\ &= f_H(\mathbf{y} | \tilde{\mathbf{B}}_y, \tilde{\mathbf{\Omega}}_y) \times \int_{D_r} f_{N+C+\tilde{G}}(\mathbf{r} | \tilde{\mathbf{B}}_{\tilde{\mathbf{u}}}, \tilde{\mathbf{\Omega}}_{\tilde{\mathbf{u}}}) d\mathbf{r}, \end{aligned} \quad (19)$$

where the integration domain $D_r = \{\mathbf{r} : \tilde{\boldsymbol{\psi}}_{low} \leq \mathbf{r} \leq \tilde{\boldsymbol{\psi}}_{up}\}$ is simply the multivariate region of the elements of the $\tilde{\mathbf{u}}$ vector determined by the observed ordinal indicator outcomes, and the range $(-\infty_{\tilde{G}}, \mathbf{0}_{\tilde{G}})$ for the utility differences is taken with respect to the utility of the observed choice alternative for the nominal outcome. $f_H(\mathbf{y} | \tilde{\mathbf{B}}_y, \tilde{\mathbf{\Omega}}_y)$ is the MVN density function of dimension H with a mean of $\tilde{\mathbf{B}}_y$ and a covariance of $\tilde{\mathbf{\Omega}}_y$, and evaluated at \mathbf{y} . The likelihood function for a sample of Q decision-makers is obtained as the product of the individual-level likelihood functions.

The above likelihood function involves the evaluation of an $N + C + \tilde{G}$ -dimensional rectangular integral for each decision-maker, which can be computationally expensive. Thus, the MACML approach of Bhat (2011) is used.

3.3 The Joint Mixed Model System and the MACML Estimation Approach

Consider the following (pairwise) composite marginal likelihood (CML) function formed by taking the products (across the N ordinal variables, the C count variables, and G nominal variables) of the joint pairwise probability of the chosen alternatives for a

decision-maker, and computed using the analytic approximation of the multivariate normal cumulative distribution (MVNCD) function.

$$\begin{aligned}
L_{CML}(\boldsymbol{\delta}) = & f_H(\mathbf{y} | \tilde{\mathbf{B}}_y, \tilde{\boldsymbol{\Omega}}_y) \times \left(\prod_{n=1}^{N-1} \prod_{n'=n+1}^N \Pr(j_n = a_n, j_{n'} = a_{n'}) \right) \times \left(\prod_{c=1}^{C-1} \prod_{c'=c+1}^C \Pr(k_c = r_c, k_{c'} = r_{c'}) \right) \times \\
& \left(\prod_{n=1}^N \prod_{c=1}^C \Pr(j_n = a_n, k_c = r_c) \right) \times \left(\prod_{n=1}^N \prod_{g=1}^G \Pr(j_n = a_n, i_g = m_g) \right) \times \\
& \left(\prod_{c=1}^C \prod_{g=1}^G \Pr(k_c = r_c, i_g = m_g) \right) \times \left(\prod_{g=1}^{G-1} \prod_{g'=g+1}^G \Pr(i_g = m_g, i_{g'} = m_{g'}) \right).
\end{aligned} \tag{20}$$

In the above CML approach, the MVNCD function appearing in the CML function is of dimension equal to (1) two for the second component (corresponding to each pair of observed ordinal outcomes), (2) two for the third component (corresponding to each pair of count outcomes), (3) two for the fourth component (corresponding to each pair of an ordinal outcome and a count outcome), (4) I_g for the fifth component (corresponding to each pair of a nominal variable and an ordinal variable), (5) I_g for the sixth component (corresponding to a nominal variable and a count variable), and (6) $I_g + I_{g'} - 2$ for the seventh component (corresponding to a pair of nominal outcomes g and g'). The net result is that the pairwise likelihood function now only needs the evaluation of a cumulative normal distribution function of dimension that is utmost equal to the sum of the alternatives associated with the pair of nominal variables with the two highest number of alternatives.

To explicitly write out the CML function in terms of the standard and bivariate standard normal density and cumulative distribution function, define $\boldsymbol{\omega}_\Delta$ as the diagonal matrix of standard deviations of matrix Δ , using $\phi_R(\cdot; \Delta^*)$ for the multivariate standard normal density function of dimension R and correlation matrix Δ^* ($\Delta^* = \boldsymbol{\omega}_\Delta^{-1} \Delta \boldsymbol{\omega}_\Delta^{-1}$), and $\Phi_E(\cdot; \Delta^*)$ for the multivariate standard normal cumulative distribution function of dimension E and correlation matrix Δ^* . Define a set of two selection matrices as follows:

(1) \mathbf{D}_{vg} is an $I_g \times (N + C + \tilde{G})$ selection matrix with an entry of ‘1’ in the first row and the v^{th} column, an identity matrix of size $I_g - 1$ occupying the last $I_g - 1$ rows and the

$N + C + \left[\sum_{j=1}^{g-1} (I_j - 1) + 1 \right]^{\text{th}}$ through $N + C + \left[\sum_{j=1}^g (I_j - 1) \right]^{\text{th}}$ columns (with the convention

that $\sum_{j=1}^0 (I_j - 1) = 0$), and entries of ‘0’ everywhere else, (2) $\mathbf{R}_{gg'}$ is a

$(I_g + I_{g'} - 2) \times (N + C + \tilde{G})$ selection matrix with an identity matrix of size $(I_g - 1)$

occupying the first $(I_g - 1)$ rows and the $N + C + \left[\sum_{j=1}^{g-1} (I_j - 1) + 1 \right]^{\text{th}}$ through

$N + C + \left[\sum_{j=1}^g (I_j - 1) \right]^{th}$ columns (with the convention that $\sum_{j=1}^0 (I_j - 1) = 0$), and another identity matrix of size $(I_{g'} - 1)$ occupying the last $(I_{g'} - 1)$ rows and the $N + C + \left[\sum_{j=1}^{g'-1} (I_j - 1) + 1 \right]^{th}$ through $N + C + \left[\sum_{j=1}^{g'} (I_j - 1) \right]^{th}$ columns; all other elements of $\mathbf{R}_{gg'}$ take a value of zero. Also, let $\widehat{\Omega}_{vg} = \mathbf{D}_{vg} \widetilde{\Omega}_{\tilde{u}} \mathbf{D}'_{vg}$, $\widehat{\Omega}_{gg'} = \mathbf{R}_{gg'} \widetilde{\Omega}_{\tilde{u}} \mathbf{R}'_{gg'}$, $\mu_{v,up} = \frac{[\widetilde{\psi}_{up}]_v - [\widetilde{\mathbf{B}}_{\tilde{u}}]_v}{\sqrt{[\widetilde{\Omega}_{\tilde{u}}]_{vv}}}$, $\mu_{v,low} = \frac{[\widetilde{\psi}_{low}]_v - [\widetilde{\mathbf{B}}_{\tilde{u}}]_v}{\sqrt{[\widetilde{\Omega}_{\tilde{u}}]_{vv}}}$, $\rho_{vv'} = \frac{[\widetilde{\Omega}_{\tilde{u}}]_{vv'}}{\sqrt{[\widetilde{\Omega}_{\tilde{u}}]_{vv} [\widetilde{\Omega}_{\tilde{u}}]_{vv'}}$, where $[\widetilde{\psi}_{up}]_v$ represents the v^{th} element of $\widetilde{\psi}_{up}$ (and similarly for other vectors), and $[\widetilde{\Omega}_{\tilde{u}}]_{vv'}$ represents the vv'^{th} element of the matrix $\widetilde{\Omega}_{\tilde{u}}$. Then,

$$\begin{aligned}
L_{CML}(\delta) = & \left(\prod_{h=1}^H \omega_{\widehat{\Omega}_y} \right)^{-1} \phi_H \left(\left[\omega_{\widehat{\Omega}_y} \right]^{-1} [\mathbf{y} - \widetilde{\mathbf{B}}_y], \widehat{\Omega}_y^* \right) \times \\
& \left(\prod_{v=1}^{N+C-1} \prod_{v'=v+1}^{N+C} \left[\Phi_2(\mu_{v,up}, \mu_{v',up}, \rho_{vv'}) - \Phi_2(\mu_{v,up}, \mu_{v',low}, \rho_{vv'}) \right] \right) \times \\
& \left(\prod_{v=1}^{N+C} \prod_{g=1}^G \Phi_{I_g} \left[\omega_{\widehat{\Omega}_{vg}}^{-1} \mathbf{D}_{vg} \{ \widetilde{\psi}_{up} - \widetilde{\mathbf{B}}_{\tilde{u}} \}, \widehat{\Omega}_{vg}^* \right] - \Phi_{I_g} \left[\omega_{\widehat{\Omega}_{vg}}^{-1} \mathbf{D}_{vg} \{ \widetilde{\psi}_{low} - \widetilde{\mathbf{B}}_{\tilde{u}} \}, \widehat{\Omega}_{vg}^* \right] \right) \times \\
& \left(\prod_{g=1}^{G-1} \prod_{g'=1}^G \Phi_{I_g + I_{g'} - 2} \left[\omega_{\widehat{\Omega}_{gg'}}^{-1} \mathbf{R}_{gg'} \{ -\widetilde{\mathbf{B}}_{\tilde{u}} \}, \widehat{\Omega}_{gg'}^* \right] \right),
\end{aligned} \tag{21}$$

where $\widetilde{\psi}_{low} = \left[\widetilde{\psi}'_{low}, \widetilde{\psi}'_{low}, (\mathbf{0}_{\widehat{G}})' \right]'$.

In Equation (21), the first component corresponds to the marginal likelihood of the continuous outcomes, the second component corresponds to the likelihood of pairs of outcomes across all ordinal and count outcomes (essentially this combines the second, third, and fourth components of Equation (20)), the third component corresponds to the pairwise likelihood for ordinal/count outcomes and nominal outcomes (this combines the fifth and sixth components of Equation (20)), and the last component corresponds to the pairwise likelihood for the nominal outcomes (this is also the last component of expression (20)). In the MACML approach, all MVNVD function evaluations greater than two dimensions are evaluated using an *analytic approximation* method rather than a simulation method. This combination of the CML with an analytic approximation for the MVNCD function is effective because the analytic approximation involves only univariate and bivariate cumulative normal distribution function evaluations. The MVNCD analytic approximation method used here is based on linearization with binary variables (see Bhat, 2011). As has been demonstrated by Bhat and Sidharthan (2011), the MACML method has the virtue of computational robustness in that the approximate CML surface is smoother and easier to maximize than are traditional simulation-based

likelihood surfaces. We can write the resulting equivalent of Equation (21) computed using the analytic approximation for the MVNCD function as $L_{MACML,q}(\boldsymbol{\delta})$, after introducing the index q for individuals. The MACML estimator is then obtained by maximizing the following function:

$$\log L_{MACML}(\boldsymbol{\delta}) = \sum_{q=1}^Q \log L_{MACML,q}(\boldsymbol{\delta}). \quad (22)$$

The covariance matrix of the parameters $\boldsymbol{\delta}$ may be estimated by the inverse of Godambe's (1960) sandwich information matrix (see Zhao and Joe, 2005; Bhat, 2014).

$$V_{MACML}(\boldsymbol{\delta}) = \frac{[\hat{\mathbf{G}}(\boldsymbol{\delta})]^{-1}}{Q} = \frac{[\hat{\mathbf{H}}^{-1}][\hat{\mathbf{J}}][\hat{\mathbf{H}}^{-1}]}{Q}, \quad (23)$$

$$\text{with } \hat{\mathbf{H}} = -\frac{1}{Q} \left[\sum_{q=1}^Q \frac{\partial^2 \log L_{MACML,q}(\boldsymbol{\delta})}{\partial \boldsymbol{\delta} \partial \boldsymbol{\delta}'} \right]_{\hat{\boldsymbol{\delta}}_{MACML}}$$

$$\hat{\mathbf{J}} = \frac{1}{Q} \sum_{q=1}^Q \left[\left(\frac{\partial \log L_{MACML,q}(\boldsymbol{\delta})}{\partial \boldsymbol{\delta}} \right) \left(\frac{\partial \log L_{MACML,q}(\boldsymbol{\delta})}{\partial \boldsymbol{\delta}'} \right) \right]_{\hat{\boldsymbol{\delta}}_{MACML}}. \quad (24)$$

An alternative estimator for $\hat{\mathbf{H}}$ may be obtained by computing the quantity below for each decision-maker, and averaging across decision-makers:

$$\hat{\mathbf{H}} \text{ for each } q = \left(\begin{array}{l} \left[\frac{\partial \log[f_H(\mathbf{y} | \tilde{\mathbf{B}}_y, \tilde{\boldsymbol{\Omega}}_y)]}{\partial \boldsymbol{\delta}} \right] \left[\frac{\partial \log[f_H(\mathbf{y} | \tilde{\mathbf{B}}_y, \tilde{\boldsymbol{\Omega}}_y)]}{\partial \boldsymbol{\delta}'} \right] + \\ \sum_{n=1}^{N-1} \sum_{n'=n+1}^N \left[\frac{\partial \log[\Pr(j_n = a_n, j_{n'} = a'_{n'})]}{\partial \boldsymbol{\delta}} \right] \left[\frac{\partial \log[\Pr(j_n = a_n, j_{n'} = a'_{n'})]}{\partial \boldsymbol{\delta}'} \right] + \\ \sum_{c=1}^{C-1} \sum_{c'=c+1}^C \left[\frac{\partial \log[\Pr(k_c = r_c, k_{c'} = r_{c'})]}{\partial \boldsymbol{\delta}} \right] \left[\frac{\partial \log[\Pr(k_c = r_c, k_{c'} = r_{c'})]}{\partial \boldsymbol{\delta}'} \right] + \\ \sum_{n=1}^N \sum_{c=1}^C \left[\frac{\partial \log[\Pr(j_n = a_n, k_c = r_c)]}{\partial \boldsymbol{\delta}} \right] \left[\frac{\partial \log[\Pr(j_n = a_n, k_c = r_c)]}{\partial \boldsymbol{\delta}'} \right] + \\ \sum_{n=1}^N \sum_{g=1}^G \left[\frac{\partial \log[\Pr(j_n = a_n, i_g = m_g)]}{\partial \boldsymbol{\delta}} \right] \left[\frac{\partial \log[\Pr(j_n = a_n, i_g = m_g)]}{\partial \boldsymbol{\delta}'} \right] + \\ \sum_{c=1}^C \sum_{g=1}^G \left[\frac{\partial \log[\Pr(k_c = r_c, i_g = m_g)]}{\partial \boldsymbol{\delta}} \right] \left[\frac{\partial \log[\Pr(k_c = r_c, i_g = m_g)]}{\partial \boldsymbol{\delta}'} \right] + \\ \sum_{g=1}^{G-1} \sum_{g'=g+1}^G \left[\frac{\partial \log[\Pr(i_g = m_g, i_{g'} = m_{g'})]}{\partial \boldsymbol{\delta}} \right] \left[\frac{\partial \log[\Pr(i_g = m_g, i_{g'} = m_{g'})]}{\partial \boldsymbol{\delta}'} \right] + \end{array} \right) \quad (25)$$

An important part of optimizing any such function is the generation of good start values. In our procedure, we came up with good start values in two steps as follows: (1) First, the

reduced form Equation (15) is estimated ignoring the latent variables; that is, setting all elements of \vec{d} and ω to zero, and setting the elements of α to zero and setting Γ to be a unit diagonal matrix, (2) Next, all the estimated parameters from step 1 are fixed, and the matrices/vectors α , \vec{d} , ω , and Γ are estimated. This produces initial estimates of all the relevant parameters, which is used to begin the iterations to maximize Equation (22). The optimization was undertaken using the GAUSS programming language, and we did not encounter any convergence issues during the optimization procedure.

3.4 Positive Definiteness

The matrix $\tilde{\Omega}$ for each household has to be positive definite (that is, all the eigenvalues of the matrix should be positive, or, equivalently, the determinant of the entire matrix and every principal submatrix of $\tilde{\Omega}$ should be positive). The simplest way to guarantee this in our mixed model system is to ensure that the $(L \times L)$ correlation matrix Γ is positive definite, and each matrix $\tilde{\Lambda}_g$ ($g=1,2,\dots,G$) is also positive definite. An easy way to ensure the positive-definiteness of these matrices is to use a Cholesky decomposition and parameterize the CML function in terms of the Cholesky parameters. Then, we use the Cholesky-decomposed parameters as the ones to be estimated. That is, the Cholesky of an initial positive-definite specification of the correlation matrix Γ and the covariance matrices $\tilde{\Lambda}_g$ ($g=1,2,\dots,G$) is taken before starting the optimization routine to maximize the CML function. Then, within the optimization procedure, one can construct the $\tilde{\Omega}$ matrix, and then pick off the appropriate elements of this matrix to obtain the CML function at each iteration. Further, because the matrix Γ is a correlation matrix, we write each diagonal element (say the aa^{th} element) of the lower triangular Cholesky matrix of Γ as $\sqrt{1 - \sum_{j=1}^{a-1} p_{aj}^2}$, where the p_{aj} elements are the Cholesky factors that are to be estimated. In addition, note that the top diagonal element of each $\tilde{\Lambda}_g$ matrix has to be normalized to one (as discussed in Section 2.2), which implies that the first element of the Cholesky matrix of each $\tilde{\Lambda}_g$ is fixed to the value of one.

Chapter 4. Simulation Experiment

In this section, we present the design of, and results from, a simulation experiment to evaluate the performance of the MACML approach to recover parameters in a GHDM system from different finite sample sizes. For ease in interpretation and understanding, the simulation design is motivated from an integrated land use-transportation context. Specifically, consider the situation where an analyst wants to examine residential location choices and travel choices of an individual using a cross-sectional data set, with a specific interest on whether (and how much) a neo-urbanist design (compact built environment design, high bicycle lane and roadway street density, good land-use mix, and good transit and non-motorized mode accessibility/facilities) would help in reducing motorized auto ownership of the household of which the individual is a part, and in influencing the individual's commute mode in a way that reduces solo auto mode use. In doing so, the analyst should consider what is commonly labeled as residential self-selection; that is, cross-sectional data reflect residential location preferences co-mingled with the travel preferences of individuals. For example, individuals who have an overall travel freedom and privacy orientation (typically associated with auto inclination) may locate themselves in suburban/rural neighborhoods (low population density, low bicycle lane and roadway street density, primarily single use residential land use, and auto-dependent urban design), own many motorized autos, and favor driving alone to work and other activities. On the other hand, a household whose members have a green and active lifestyle propensity may seek out urban neighborhoods so they can pursue their activities using non-motorized and transit modes of travel. If such self-selection effects in residence choices are ignored, when actually present, the result can be a “spurious” causal effect of neighborhood attributes on auto ownership and travel, and potentially misinformed BE design policies (see a detailed discussion in Bhat *et al.*, 2014a). But the self-selection may not be based solely on residential choice, and can also be based on auto ownership choice. Thus, individuals with a travel freedom and privacy orientation may both prefer more autos as well as be predisposed to traveling in motorized vehicles to work and other activities. As a consequence, any effect of the number of motorized vehicles on auto travel will be moderated by the travel freedom and privacy orientation of the individual.

The potential self-selection effects above can be acknowledged by considering workers' decisions associated with residential location, auto ownership, commute travel mode choice, and some quantification of non-commute travel as a multi-dimensional bundle. It is in this context that our simulation design is set. Residential location choice is represented as a nominal discrete choice among a multinomial set of three different types of BE designs as captured by designations as urban, suburban, and rural neighborhoods (these designations can be combinations of housing density and employment density; see Kim and Brownstone, 2013, Paleti *et al.*, 2013, Cao and Fan, 2012, and Bhat *et al.*, 2014a, who all use such a density-based classification scheme as a representation of residential location choice as this simplifies the representation of residential choice alternatives and also alleviates the problem of strong multi-collinearity of density with other built environment attributes). In addition, we also use a second continuous outcome, the (logarithm of) commute distance for the individual, to characterize residential location choice. This is because it has been well established in the literature

that commute distance is one of the most important determinants of residential location (see, for example, Clark *et al.*, 2003, Rashidi *et al.*, 2012).⁴ Auto ownership is a count outcome, while commute travel mode choice is represented as a second nominal choice in the system from among three different modes of transportation – non-motorized transportation (NM), public transportation (PT), and motorized (private) transportation or MT (either as a driver or a passenger). Non-commute travel is quantified as a multi-dimensional bundle of three ordinal variables that relate to intensities (occurrences) of weekly non-commute travel by NM, by PT, and by MT. However, since most household travel surveys capture only daily travel, we suppose that use of alternative modes over longer periods of time (as would be important particularly for NM and PT use) is obtained through an ordinal categorical indicator response from among three possibilities: (1) Never or about once a week, (2) about 2-3 times a week, and (3) four or more times in a week (see Sener *et al.*, 2009 for a survey that captures non-commute travel in such ordinal categories). In all, our system has seven endogenous outcomes/indicators, with one continuous outcome (commute distance), three ordinal indicators (non-commute travel occurrences by NM, PT, and MT), one count outcome (auto ownership), and two nominal outcomes (residential choice location based on density categorization and commute mode choice). While modeling all of these as a joint bundle, we also accommodate structural relationships among the endogenous outcomes/indicators. In particular, we specify that commute distance and auto ownership will affect commute mode choice, and the geographic area of residential location (urban, suburban, or rural) will affect auto ownership, commute distance, and non-commute travel occurrences by NM and PT.

4.1 Experimental Design

Consider a multi-dimensional choice bundle of residential location and activity-travel behavior, as discussed in the previous section. In previous studies on the integration of land-use patterns and activity-travel behavior, such as Pinjari *et al.* (2011) and Bhat *et al.* (2014a), correlated unobserved effects among multiple (but limited) choice dimensions were captured through the error terms of the many individual dimensions, resulting in a relatively large dimensional covariance matrix. The difference between these earlier studies and this simulation study is that, as discussed in Section 1, the covariance in a large number of choice dimensions is captured in a parsimonious manner through a factor-analytic structure where the choice dimensions are a function of a smaller dimension of correlated latent constructs. In addition, such a specification provides structure to the jointness among the choice dimensions by appealing to theoretical psychological constructs.

4.2 The Structural Equation System

Two latent variables associated with lifestyle and attitudes are employed as psychological constructs impacting the multi-dimensional choice bundle of residential location and

⁴ The implicit assumption here is that work location choices precede residential choice. While it is certainly possible that residential moves may motivate job moves, earlier research using panel data suggests that a vast majority (85% or more) of residential relocations follow a job move (see Rashidi *et al.*, 2012).

activity-travel behavior. The latent variables are shown in Figure 1, where the ovals represent the latent constructs, while rectangles represent observed explanatory variables. The first latent factor is *green lifestyle propensity* (z_1^*) or the individual's level of environmental consciousness, which is specified to be a function of whether the individual has a Bachelor's degree or higher (w_1 ; $w_1 = 1$ if individual has a Bachelor's degree or higher and 0 otherwise) and whether the individual is male or female (w_2 ; $w_2 = 1$ if individual is male and 0 otherwise). These reflect the finding from earlier studies that individuals with a Bachelor's degree or higher tend to be more active proponents and followers of ecologically friendly lifestyles (Paleti *et al.*, 2013), as do women compared to men (see, for example, Liu *et al.*, 2014 and Gifford and Nilsson, 2014). The specified values of these effects (embedded within the α_1 vector) are 0.8 (for the education effect) and -0.3 (for the male gender effect). The second factor is *travel freedom/privacy affinity* (z_2^*), generally associated with travel comfort/convenience and a sense of control over the travel experience. This latent variable is specified to be associated with men (w_2 ; $w_2 = 1$ if individual is male and 0 otherwise), and high income individuals (w_3 ; $w_3 = 1$ if individual earns a high income and zero otherwise). Earlier studies, including Schwanen and Mokhtarian (2007), Jansen, 2012, Shiftan *et al.*, 2008, and Day, 2000, have indicated that men and high income earners generally value travel freedom/privacy more than women and low income earners, respectively. The design values of these effects in the simulation (as embedded within the α_2 vector) are 0.2 and 0.5, respectively. In the vector notation of Equation (2), the effects in Figure 1 may be written as follows:

$$\begin{bmatrix} z_1^* = \text{GLP} \\ z_2^* = \text{TFA} \end{bmatrix} = \begin{bmatrix} 0.8 & -0.3 & 0.0 \\ 0.0 & 0.2 & 0.5 \end{bmatrix} \begin{bmatrix} w_1 = \text{Bachelor's degree or higher or not} \\ w_2 = \text{Male or not} \\ w_3 = \text{High income or not} \end{bmatrix} + \begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix},$$

where GLP is *green lifestyle propensity* and TFA is *travel freedom/privacy affinity*. The parameters in the matrix α to be estimated can be stacked into a vector $\text{Vech}(\alpha) = [\alpha_{11} = 0.8, \alpha_{12} = -0.3, \alpha_{22} = 0.2, \alpha_{33} = 0.5]$. The correlation matrix of the error vector η is specified as follows:

$$\text{Var}(\eta) = \Gamma = \begin{bmatrix} 1.0 & -0.6 \\ -0.6 & 1.0 \end{bmatrix} = \mathbf{L}_\Gamma \mathbf{L}'_\Gamma = \begin{bmatrix} 1.0 & 0.0 \\ -0.6 & 0.8 \end{bmatrix} \begin{bmatrix} 1.0 & -0.6 \\ 0.0 & 0.8 \end{bmatrix}.$$

In the matrix above, we allow a correlation (entry of -0.6) between the two latent propensity constructs of GLP and TFA to reflect the existence of the unobserved underlying value of individuality that affects both of these personality constructs. To ensure the positive definiteness of Γ , a Cholesky decomposition is conducted. In our specification, a single element is to be estimated in the matrix Γ : $l_\Gamma = -0.6$.

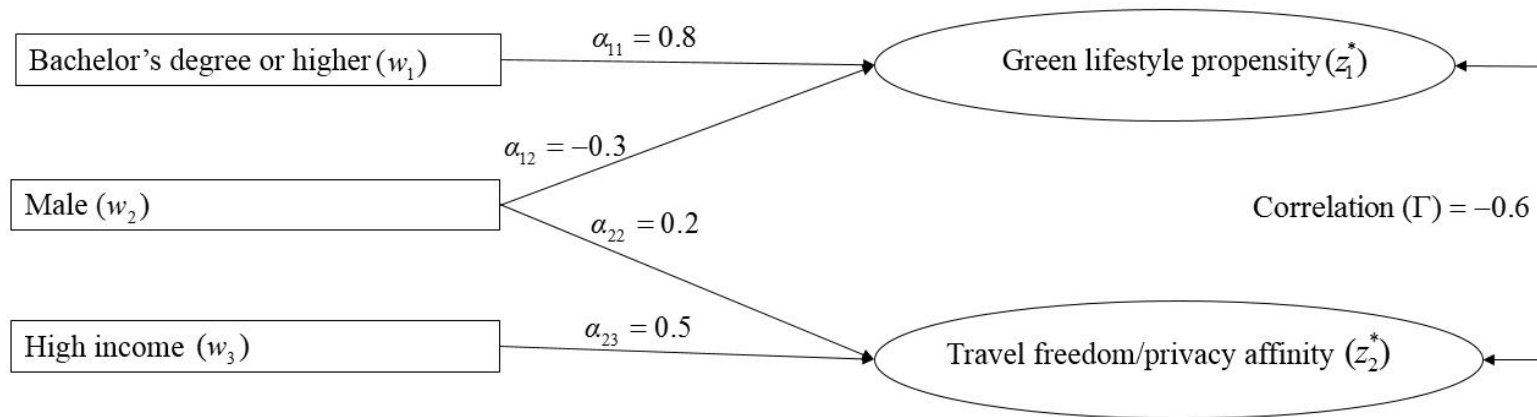


Figure 1: Diagrammatic representation of the structural equation

4.3 The Measurement Equation System

The measurement equation system includes the non-nominal equation system $\vec{y} = \vec{\gamma}x + \vec{d}z^* + \vec{\varepsilon}$ (Equation (13) earlier) as well as the nominal equation system $U = bx + \varpi z^* + \varsigma$ (Equation (14) earlier). Within each of these systems, there are exogenous and endogenous outcome effects (embedded in $\vec{\gamma}$ and $\vec{\gamma}$ for the non-nominal system and in b for the nominal system), as well as latent construct effects (embedded in \vec{d} and ϖ). The simulation design effects specified for the non-nominal equation system (including both the exogenous and latent construct effects) are presented in Figure 2a, while the corresponding effects for the nominal equation system are presented in Figure 2b. Finally, the endogenous variable effects (that is, the inter-relationships between the endogenous outcomes/indicators, which can only be recursive as discussed in Section 2.2), are presented in Figure 2c. Each of these effects is discussed in turn in the subsequent sections, while Section 4.3.4 brings all parameters to be estimated together in the measurement equation system. Note that the design considers four exogenous variables: (1) whether the individual is an immigrant or not (a dummy variable “immigrant” taking the value of 1 if the individual is born in the US and 0 otherwise), (2) whether the individual owns or rents her/his household (a dummy variable “owns hh” taking the value of 1 if the individual owns her/his household and 0 otherwise), (3) number of children less than 11 years of age, and (4) number of young active adults (to represent the presence of the so-called millenials born between 1981 and 1996).

4.3.1 Non-Nominal Equation System with Exogenous and Latent Construct Effects

This system is shown diagrammatically in Figure 2a. Immigrant status positively influences (log) commute distance, as it has been observed that immigrants have longer commutes than do non-immigrants (see Paleti *et al.*, 2013). Further, individuals with young children are less likely to travel by non-motorized modes and more likely to travel by motorized vehicles (as they undertake pick up/drop off activities; see Sener *et al.*, 2009). Also, in the simulation design, we specify the number of young active adults in the individual’s household to negatively influence travel by motorized vehicles, as households with millenials tend to undertake their out-of-home activities less using private vehicles (see Bhat *et al.*, 2014a). A total of four exogenous variable effects are specified above. However, there are also constants to be specified in the (log) commute distance equation, and for the latent propensities for the ordinal indicators. The constant in the (log) commute distance equation as well as the constant effects for all the ordinal indicators are set to the value of 1.0.

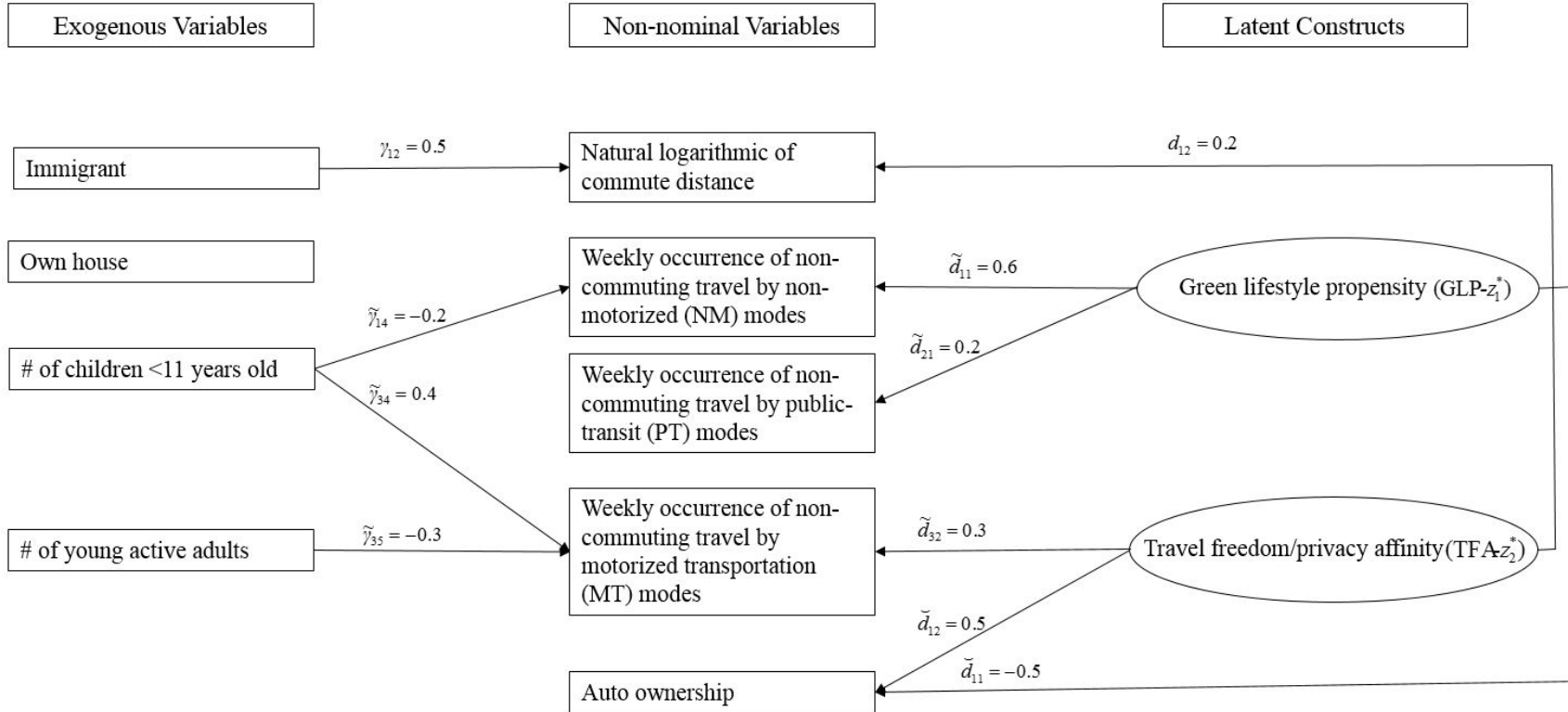


Figure 2a: Diagrammatic representation of the measurement equation for the non-nominal variables

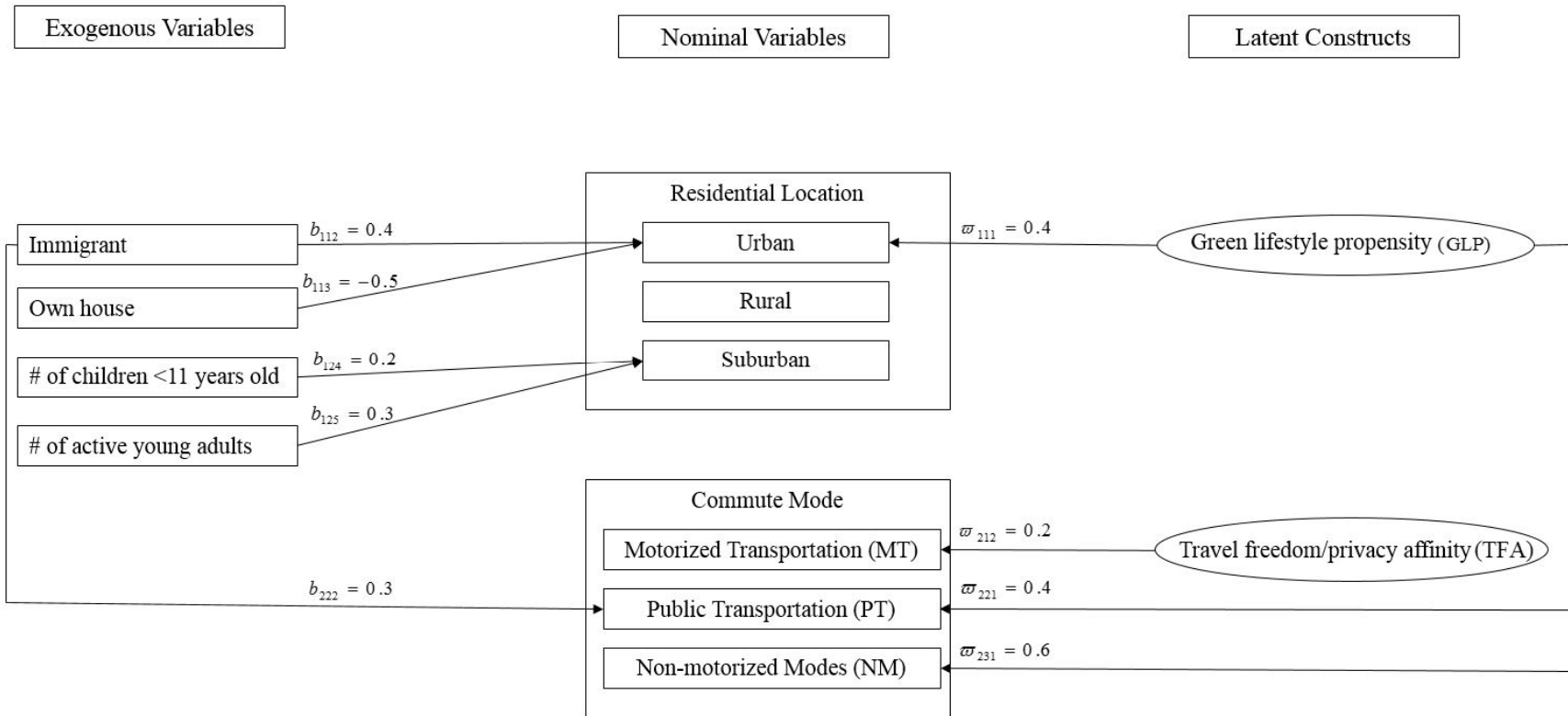


Figure 2b: Diagrammatic representation of the measurement equation for the nominal variables

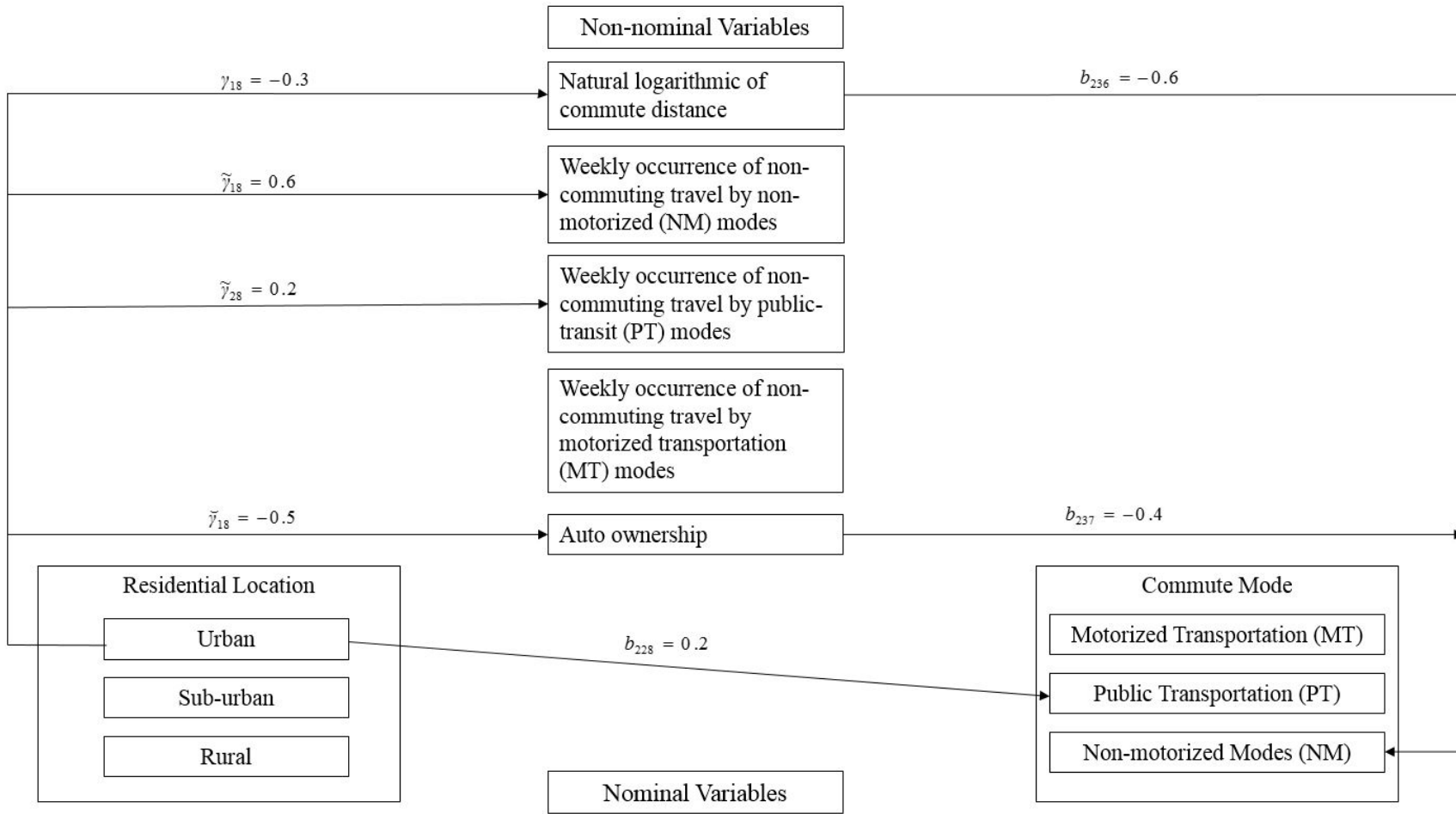


Figure 2c: Endogeneous effects

A total of five latent construct effects are also specified (see toward the right of Figure 2a). As expected, a green lifestyle propensity (GLP) increases non-commute travel occurrences by non-motorized (NM) modes as well as increases non-commute travel occurrences by public transit (PT) modes. These effects satisfy the two-indicator rule for the GLP latent construct. Similarly, we expect travel freedom/privacy affinity (TFA) to be positively related to commute distance (see, for example, Schwanen and Mokhtarian, 2007) and non-commute travel occurrences by motorized transport (MT) modes. These effects satisfy the two-indicator rule for the TFA latent construct. Finally, both GLP and TFA are specified to impact auto ownership, with the former having a negative effect and the latter a positive effect.

As presented in Equation (13), the covariance matrix $\bar{\Sigma}$ of random error $\bar{\epsilon}$ for non-nominal indicators is restricted to be diagonal, with elements corresponding to ordinal and count indicators being normalized to 1. This leaves the variance component for the continuous outcome (logarithm of commute distance), which is specified to be 1.25 in the simulation design. Thus, the one element to be estimated in the matrix $\bar{\Sigma}$ is 1.25, which we write as $l_{\bar{\Sigma}} = 1.25$.

There are three ordinal outcomes (non-commute travel occurrences by NM, PT, and MT), in the simulation design, which leads to the need to specify $\tilde{\psi}_{n,2}$ for each ordinal outcome n ($n = 1, 2, 3$) (see discussion in Section 2.2). All of these threshold values are set to 1.5. In addition, we need to specify the parameters in the threshold function for the count outcome (corresponding to auto ownership). This refers to the coefficient vector $\tilde{\gamma}$, the flexibility parameter vector $\boldsymbol{\varphi}_c = (\varphi_{c,1}, \varphi_{c,2}, \dots, \varphi_{c,e_c^*})'$, and the dispersion parameter vector $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_c)'$. For the $\tilde{\gamma}$ coefficient vector, we include only a constant effect and another endogenous effect (the latter is discussed in the next section). The coefficient on the constant is specified to be 1.0. For the flexibility vector, we will drop the index c since we have only one count outcome in the simulation design. We also specify a single flexibility parameter $\varphi_1 = 0.75$. For the dispersion parameter vector (which collapses to a scalar because there is only a single count outcome), we specify $\theta = 2.0$.

4.3.2 Nominal Equation System with Exogenous and Latent Construct Effects

Five exogenous effects and four latent construct effects are specified here (see Figure 2b). All of the exogenous effects specified have been reasonably well established in earlier studies. Immigrants tend to cluster in urban neighborhoods (see Bhat *et al.*, 2013), while those who own households are less likely to reside in urban neighborhoods. There is also evidence that individuals with children tend to favor suburban neighborhoods due to the open spaces and good quality schools (Aditjandra *et al.*, 2012), as do households with many young active adults (Brownstone and Golob, 2009). Further, as has been found in many earlier studies, immigrants, more so than US-born individuals, tend to use public transportation for their commute. In addition to the variable effects above, we also allow constants in two of the utilities for residential location and two of the utilities for commute mode. Specifically, we use a constant effect of 0.2 in the urban location utility and 0.3 in the suburban location utility (with the rural constant specified to be zero for

identification). Also, we use a constant effect of -0.5 for the PT mode, and -0.2 for the NM mode (with the MT mode constant specified to be zero).

The latent construct effects specified are rather intuitive. These are specified to shift the utility of specific alternatives of the nominal variables. Essentially, then, in the notation of Section 2.2, $\boldsymbol{\omega}_g = \boldsymbol{\vartheta}_g$, because $\boldsymbol{\beta}_g$ is an identity matrix. Thus, for convenience, we will refer to the parameters to be estimated as being elements of $\boldsymbol{\omega}_g$, which are the same as the elements of $\boldsymbol{\vartheta}_g$. For the residential location nominal outcome, individuals with a green lifestyle propensity tend to reside in urban neighborhoods, so that they can pursue their desired lifestyles due to greater opportunities to pursue city life while adopting green modes of transportation (Schwanen and Mokhtarian, 2007). For the commute mode nominal outcome, green lifestyle propensity is specified to positively affect the use of PT and NM modes, while travel freedom/privacy affinity increases the propensity to use the MT mode.

The covariance matrix of $\boldsymbol{\zeta}$ is specified as follows.

$$\begin{aligned} \text{Var}(\boldsymbol{\zeta}) = \boldsymbol{\Lambda} &= \begin{bmatrix} 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 \\ 0.00 & 1.00 & 0.70 & 0.00 & 0.00 & 0.00 \\ 0.00 & 0.70 & 1.49 & 0.00 & 0.00 & 0.00 \\ 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 \\ 0.00 & 0.00 & 0.00 & 0.00 & 1.00 & 0.60 \\ 0.00 & 0.00 & 0.00 & 0.00 & 0.60 & 1.36 \end{bmatrix} \\ &= \mathbf{L}_\Lambda \mathbf{L}'_\Lambda = \begin{bmatrix} 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 \\ 0.00 & 1.00 & 0.00 & 0.00 & 0.00 & 0.00 \\ 0.00 & 0.70 & 1.00 & 0.00 & 0.00 & 0.00 \\ 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 \\ 0.00 & 0.00 & 0.00 & 0.00 & 1.00 & 0.00 \\ 0.00 & 0.00 & 0.00 & 0.00 & 0.60 & 1.00 \end{bmatrix} \begin{bmatrix} 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 \\ 0.00 & 1.00 & 0.70 & 0.00 & 0.00 & 0.00 \\ 0.00 & 0.00 & 1.00 & 0.00 & 0.00 & 0.00 \\ 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 \\ 0.00 & 0.00 & 0.00 & 0.00 & 1.00 & 0.60 \\ 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 1.00 \end{bmatrix} \end{aligned} \quad (30)$$

In the matrix $\boldsymbol{\Lambda}$, four elements are to be estimated ($l_{\Lambda32} = 0.70, l_{\Lambda33} = 1.49, l_{\Lambda65} = 0.60, l_{\Lambda66} = 1.36$).

4.3.3 Endogenous Outcome Effects

These effects correspond to recursive effects among the endogenous outcomes, as discussed just before Section 4.1. These are parts of the $\tilde{\boldsymbol{\gamma}}$ matrix (for the continuous/ordinal outcomes), the $\tilde{\boldsymbol{\gamma}}$ matrix (for the count outcomes), and the \mathbf{b} matrix (for the nominal outcomes). The important point is that these are “cleansed” effects after accommodating unobserved covariance effects among the endogenous outcomes engendered by the presence of latent constructs, as discussed in the previous two sections. Figure 2c provides a pictorial representation for these endogenous effects. For the continuous/ordinal outcomes, we specify that urban dwelling leads to a shorter

commute distance, and more non-commute travel occurrences by the NM and PT modes (see Paleti *et al.*, 2013). For the auto count variable, several earlier studies have established that urban dwellers tend to own fewer vehicles even after accounting for any residential self-selection effects (see, for example, Bhat and Guo, 2007). This effect is specified through the threshold in the count model; that is, in the \mathbf{x} vector with a corresponding coefficient vector $\tilde{\gamma}$ (the $\tilde{\gamma}$ matrix becomes a vector in our simulation design because there is only one count variable). In particular, in our formulation of the count model, a positive coefficient element in $\tilde{\gamma}$ implies that an increase in the corresponding element of \mathbf{x} shifts all the thresholds toward the left of the auto ownership propensity scale (see Castro *et al.*, 2011), which has the effect of reducing the probability of zero cars, while a negative coefficient in $\tilde{\gamma}$ implies that an increase in the corresponding element of \mathbf{x} shifts all the thresholds toward the right of the auto ownership propensity scale, which has the effect of increasing the probability of zero cars. In our simulation design, we impose a negative coefficient of -0.5.

For the nominal variables, our design specifies a positive effect of urban dwelling on the propensity to use PT as the commute mode, and a negative effect of car ownership and commute distance on the use of the NM mode for the commute.

4.3.4 Overall Measurement Equation System

The overall measurement equation for the vector $\mathbf{y}U = [\tilde{\mathbf{y}}', U']$ takes the following mathematical form:

$$\begin{array}{l}
\left[\begin{array}{l}
y_1 = \log(\text{commute dist}) \\
\tilde{y}_1^* = \text{NC propensity by NM} \\
\tilde{y}_2^* = \text{NC propensity by PT} \\
\tilde{y}_3^* = \text{NC propensity by MT} \\
\tilde{y}_1^* = \text{autoown. propensity} \\
U_{1,\text{urban}} \\
U_{1,\text{suburban}} \\
U_{1,\text{rural}} \\
U_{2,\text{MT}} \\
U_{2,\text{PT}} \\
U_{2,\text{NMT}}
\end{array} \right] = \begin{bmatrix}
1.0 & 0.5 & 0.0 & 0.0 & 0.0 & 0.0 & 0.0 & -0.3 \\
1.0 & 0.0 & 0.0 & -0.2 & 0.0 & 0.0 & 0.0 & 0.6 \\
1.0 & 0.0 & 0.0 & 0.0 & 0.0 & 0.0 & 0.0 & 0.2 \\
1.0 & 0.0 & 0.0 & 0.4 & -0.3 & 0.0 & 0.0 & 0.0 \\
0.0 & 0.0 & 0.0 & 0.0 & 0.0 & 0.0 & 0.0 & -0.5 \\
0.2 & 0.4 & -0.5 & 0.0 & 0.0 & 0.0 & 0.0 & 0.0 \\
0.3 & 0.0 & 0.0 & 0.2 & 0.3 & 0.0 & 0.0 & 0.0 \\
0.0 & 0.0 & 0.0 & 0.0 & 0.0 & 0.0 & 0.0 & 0.0 \\
0.0 & 0.0 & 0.0 & 0.0 & 0.0 & 0.0 & 0.0 & 0.0 \\
-0.5 & 0.3 & 0.0 & 0.0 & 0.0 & 0.0 & 0.0 & 0.2 \\
-0.2 & 0.0 & 0.0 & 0.0 & 0.0 & -0.6 & -0.4 & 0.0
\end{bmatrix} \left[\begin{array}{l}
\text{Constant} \\
\text{Immigrant household} \\
\text{Own house} \\
\# \text{ of Children} < 11 \text{ yrs} \\
\# \text{ of young active adults} \\
\text{Commute distance} \\
\text{auto ownership} \\
\text{urban dwelling}
\end{array} \right] \\
+ \begin{bmatrix}
0.0 & 0.2 \\
0.6 & 0.0 \\
0.2 & 0.0 \\
0.0 & 0.3 \\
-0.5 & 0.5 \\
0.4 & 0.0 \\
0.0 & 0.0 \\
0.0 & 0.0 \\
0.0 & 0.2 \\
0.4 & 0.0 \\
0.6 & 0.0
\end{bmatrix} \begin{bmatrix}
z_1^* = \text{GLP} \\
z_2^* = \text{TFA}
\end{bmatrix} + \begin{bmatrix}
\varepsilon_1 \\
\tilde{\varepsilon}_1 \\
\tilde{\varepsilon}_2 \\
\tilde{\varepsilon}_3 \\
\tilde{\varepsilon}_1 \\
\varsigma_{11} \\
\varsigma_{12} \\
\varsigma_{13} \\
\varsigma_{21} \\
\varsigma_{22} \\
\varsigma_{23}
\end{bmatrix}
\end{array}$$

Based on the above, and using the notations employed in Section 2.2., the parameters to be estimated in the measurement equation above include the following:

$$\text{Vech}(\tilde{\gamma}) = [\gamma_{11} = 1, \gamma_{12} = 0.5, \gamma_{18} = -0.3, \tilde{\gamma}_{11} = 1, \tilde{\gamma}_{14} = -0.2, \tilde{\gamma}_{18} = 0.6, \tilde{\gamma}_{21} = 1, \tilde{\gamma}_{28} = 0.2, \tilde{\gamma}_{31} = 1, \tilde{\gamma}_{34} = 0.4, \tilde{\gamma}_{35} = -0.3],$$

$\text{Vech}(\tilde{\gamma}_I) = [\tilde{\gamma}_{11} = 1, \tilde{\gamma}_{18} = -0.5]$ (this is the vector corresponding to the coefficients on the constant and the urban dwelling variable embedded in the threshold in the auto ownership count model),

$$\text{Vech}(\mathbf{b}) = [b_{111} = 0.2, b_{112} = 0.4, b_{113} = -0.5, b_{121} = 0.5, b_{124} = 0.2, b_{125} = 0.3, b_{221} = -0.5, b_{222} = 0.3, b_{228} = 0.2, b_{231} = -0.2, b_{236} = -0.6, b_{237} = -0.4],$$

$$\text{Vech}(\tilde{\mathbf{d}}) = [d_{12} = 0.2, \tilde{d}_{11} = 0.6, \tilde{d}_{21} = 0.2, \tilde{d}_{32} = 0.3, \tilde{d}_{11} = -0.5, \tilde{d}_{12} = 0.5], \text{ and}$$

$$\text{Vech}(\boldsymbol{\omega}) = [\omega_{111} = 0.4, \omega_{212} = 0.2, \omega_{221} = 0.4, \omega_{231} = 0.6].$$

In addition, we have the variance component for the continuous outcome $l_{\Sigma} = 1.25$. the flexibility parameter $\varphi_1 = 0.75$ and the dispersion parameter vector $\theta = 2.0$ for the auto ownership count outcome, the single element ($l_{\Gamma} = -0.6$) in the covariance matrix of the

error terms in the structural equation system, and the parameters for the covariance matrix of the nominal outcomes: $I_{\Lambda32} = 0.70, I_{\Lambda33} = 1.49, I_{\Lambda65} = 0.60, I_{\Lambda66} = 1.36$.

4.4 Data Generation Process

To generate the simulated dataset, the first step is to develop values for the exogenous variables in the vectors w and x . There are six dummy variables in these two vectors, corresponding to bachelor's degree or higher (w_1), person lives alone (w_2), male (w_3), high income (w_4), immigrant (x_1), and own household (x_2). To construct these dummy variables, independent values were drawn from the standard uniform distribution. If the value drawn was less than 0.5, the value of '0' was assigned for the dummy variable. Otherwise, the value of '1' was assigned. For the two count exogenous variables corresponding to the number of children less than 11 years of age and the number of young active adults, a maximum value for each variable was first assigned (three for the first, and five for the second). Then, the range of the uniform distribution (0 to 1) was divided into as many equal ranges as the maximum value for the count plus one. Independent draws for the two count variables were made from the uniform distribution, and the value assigned of the count was based on the range in which a draw fell. For example, for the "number of children less than 11 years" variable, four equal intervals were created: $[0.00, 0.25)$, $[0.25, 0.50)$, $[0.50, 0.75)$, or $[0.75, 1.00]$. If a draw was between 0.00 and 0.25 (but not including 0.25 exactly), a value of 0 was assigned for the variable; if a draw was between 0.25 and 0.5 (but not including 0.50 exactly), a value of 1 was assigned and so on.

The procedure above is used to construct a synthetic sample of $Q=1000, 2000$, and 3000 realizations of the exogenous variables. We consider different samples sizes to assess the accuracy and appropriateness of the asymptotic properties of the MACML estimator for finite sample sizes. Once drawn, the exogenous variables are held fixed for the rest of the simulation exercise. In the rest of this section, we will discuss the procedure to generate the data set assuming $Q=1000$ observations (the same procedure may be applied for $Q=2000$ and $Q=3000$ observations). For each of the 1000 observations, a specific realization of the vector $(\bar{\epsilon}', \zeta')' [(E + \bar{G}) \times 1]$ is drawn from the multivariate distribution with mean $\mathbf{0}_{11}$ (a column vector of zero values of dimension 11) and covariance structure given by $\mathbf{\Omega}$ in Equation (16). The sub-vector of the mean vector \mathbf{B}_2 that corresponds to the utilities of the three residential choice alternatives is also computed using the expression in Equation (16). Then, the realization corresponding to $\zeta_1 = (\zeta_{11}, \zeta_{12}, \zeta_{13})'$ (the error terms drawn for the three residential choice alternatives) is added to the mean vector for the three residential choice alternatives to obtain the realization of $U_1 = (U_{1,urban}, U_{1,suburban}, U_{1,rural})'$ for each observation. The alternative with the highest utility value is then picked, and identified as the chosen residential choice alternative for each observation. Next, the continuous outcome y_1 is generated based on the exogenous variables, the design parameters, and the realization of the value of ϵ_1 from earlier. Similarly, the latent continuous values for the ordinal indicators are also generated, and then translated into ordinal outcomes based on comparison with the corresponding design

thresholds. For the auto ownership count outcome, the latent continuous value is generated exactly as for the ordinal indicators. However, the thresholds also need to be computed based on the design parameters as well as the realized actual value of the urban residential choice outcome. Then, the latent continuous value for the count outcome is translated into an actual count outcomes based on a comparison with the computed thresholds. Finally, the utilities for the commute mode choice alternatives are computed based on exogenous variables, all realized values of the other endogenous outcomes, as well as the realization corresponding to $\zeta_2 = (\zeta_{21}, \zeta_{22}, \zeta_{23})'$ from earlier (the error terms drawn for the three commute mode choice alternatives).

The above data generation process is undertaken 200 times with different realizations of the random error components to generate 200 datasets for each sample size. The MACML estimator is applied to each dataset to estimate the 57 underlying parameters. A single random permutation is generated for each individual (the random permutation varies across individuals, but is the same across iterations for a given individual) to decompose the MVNCD function into a product sequence of marginal and conditional probabilities (see Section 2.1 of Bhat, 2011)⁵. In order to obtain a sense of the approximation error (explained in the following subsection), 10 datasets are randomly selected from the 200 datasets for each sample size (*i.e.*, $N=1000, 2000, \text{ and } 3000$). Then the estimator is applied to each dataset 10 times with different permutations. Based on the 100 estimations (10 datasets \times 10 runs with different permutations per dataset) for each sample size, the estimates of approximation error are derived.

4.5 Performance Evaluation

The performance of the MACML inference approach in estimating the parameters of the GHDM and the corresponding standard errors is evaluated as follows (the discussion below is for a specific sample size; the same procedure is applied for evaluating performance with the different sample sizes of 1000, 2000, and 3000).

- (1) Estimate the MACML parameters for the 200 datasets. Estimate the standard errors using the Godambe (sandwich) estimator.
- (2) Compute the mean for each model parameter across the 200 datasets to obtain a mean estimate. Compute the **absolute percentage (finite sample) bias (APB)** of the estimator as:

$$APB = \left| \frac{\text{mean estimate} - \text{true value}}{\text{true value}} \right| \times 100 \quad (31)$$

- (3) Compute the standard deviation of the mean estimate across the 200 datasets, and label this as the **finite sample standard deviation or FSSD** (essentially, this is the empirical standard error).

⁵ Technically, the MVNCD approximation should improve with a higher number of permutations in the MACML approach. However, when we investigated the effect of different numbers of random permutations per individual, we noticed little difference in the estimation results between using a single permutation and higher numbers of permutations, and hence we settled with a single permutation per individual.

- (4) Compute the mean standard error for each model parameter across the 200 datasets, and label this as **the asymptotic standard error or ASE** (essentially this is the standard error of the distribution of the estimator as the sample size gets large). Compute the ASE as a percentage of the mean estimate.
- (5) Next, to evaluate the accuracy of the asymptotic standard error formula as computed using the MACML inference approach for the finite sample size used, compute the **absolute percentage bias of the asymptotic standard error (APBASE)** for each parameter relative to the corresponding finite sample standard deviation.

$$APBASE = \left| \frac{ASE - FSSD}{FSSD} \right| \times 100$$

- (6) For each of the randomly selected 10 datasets (out of the 200 datasets), compute the mean estimate (10ME) for each model parameter across the 10 random permutations used for that dataset (to evaluate the MVNCD function). Then, for each of the 10 datasets, compute the standard deviation of the parameter values (across permutations) around the 10ME value. Take the mean of the standard deviation value across all the 10 datasets, and label this as the **approximation error (APERR)**.

4.6 Simulation Results

The simulation results for $Q=1000$, 2000, and 3000 are presented in Tables 2, 3, and 4, respectively. The tables provide the true value of the parameters (second column), followed by the parameter estimate results and the standard error estimate results.

A number of observations may be made from the tables. First, the ability of the MACML approach to recover the parameters underlying the GHDM model is pretty good, as may be observed from the magnitude of the absolute percentage bias (APB) values. In particular, the mean APB value (see the bottom row of the third column under “Parameter Estimates”) is 9.28% with 1000 observations, reducing to 8.39% with 2000 observations and further to 6.29% with 3000 observations. Overall, the difference between 1000 and 2000 observations in more accurately recovering parameters is moderate. But there is a larger difference in the APB values appears when moving from 2000 observations to 3000 observations, suggesting that there are critical thresholds in the number of observations in terms of recovering parameters well. Second, the parameters corresponding to the effects of exogenous variables on the latent variables (that is, the elements of $\text{Vech}(\alpha)$), the effects of the latent variables on the non-nominal outcomes (that is, the elements of $\text{Vech}(\vec{d})$), and the effects of the latent variables on the nominal outcomes (that is, the elements of $\text{Vech}(\vartheta)$) are generally relatively more difficult to accurately estimate compared to other parameters. Thus, for the case of $Q=1000$ observations, the APB value for the $\text{Vech}(\alpha)$ elements range from 1.006% to 28.663% with a mean APB of 14.34), the APB value for the $\text{Vech}(\vec{d})$ elements range between 6.261% and 47.373% (with a mean APB of 21.16%), and the APB values for the $\text{Vech}(\vartheta)$ elements range from 1.429% to 33.50% (with a mean of 12.43%). For datasets with 1000, 2000, and 3000 observations, the mean APB values for (a) the $\text{Vech}(\alpha)$

elements are 14.34%, 14.79%, and 7.42%, respectively, (b) the $\text{Vech}(\vec{d})$ elements are 21.16%, 20.27%, and 15.34%, respectively, and (c) for the $\text{Vech}(\varpi)$ elements are 12.43%, 7.03%, and 10.87%, respectively. The relatively less accurate recovery of these sets of parameters is intuitive. As one can notice from Equations (15) and (16), the only way to disentangle the effects of the \vec{d} matrix and the α matrix in the first (non-nominal) part of Equation (15) is through the identification of the \vec{d} matrix elements from the covariance matrix Ω . Similarly, the only way to disentangle the effects of the ϖ matrix and the α matrix in the second (nominal) part of Equation (15) is through the identification of the ϖ matrix elements from the covariance matrix Ω . As such, the \vec{d} matrix elements and the ϖ matrix elements enter into the covariance matrix Ω in a non-linear fashion (see Equation 16), and Ω itself enters into the composite likelihood function (Equation 21) in a complex manner. It is also interesting to note that the improvement in the accuracy of recovery is substantial for the $\text{Vech}(\alpha)$ and $\text{Vech}(\vec{d})$ parameters as one goes from 2000 to 3000 observations, which is essentially driving the substantially overall improved performance with 3000 observations relative to 2000 observations as pointed out earlier. An additional point to note here is that, while there are some variations in the ability to recover the latent variable loadings on different kinds of variables (continuous, ordinal, count, and nominal variables), there were no clear systematic patterns in the level of accuracy in estimating the latent factor loadings for different types of dependent variables. Third, the effects of exogenous and endogenous variables on the different kinds of variables (corresponding to $\text{Vech}(\vec{\gamma})$, $\text{Vech}(\tilde{\gamma})$, and $\text{Vech}(b)$) are accurately recovered. In general, it appears that these effects are less accurately recovered for the continuous dependent variable, relative to other types of variables (see the higher APB value for the γ_{11} , γ_{12} , and γ_{18} elements relative to other γ and b parameters in the tables). Fourth, and moving on to the standard error estimates, the entries in the “finite sample standard error (FSSE)” column indicate that the empirical ability of the MACML estimator to pin down parameters (that is, the precision of parameter recovery) is quite good. In particular, as a percentage of the true values, the mean FSSE values across all parameters are 34.09, 22.54, and 18.97 for 1000, 2000, and 3000 observations, respectively (see the last row of the sub-column entitled “% of true value” under the FSSE column). However, once again, and for the same reason that it is difficult to accurately recover the parameters of $\text{Vech}(\alpha)$, $\text{Vech}(\vec{d})$, and $\text{Vech}(\varpi)$, the FSSE values are relatively higher for these sets of parameters than for all parameters as a whole. For datasets with 1000, 2000, and 3000 observations, the FSSE values as a percentage of the true values for (a) the $\text{Vech}(\alpha)$ elements are 40%, 29%, and 20.6%, respectively, (b) the $\text{Vech}(\vec{d})$ elements are 40.9%, 25.1%, and 23.4%, respectively, and (c) for the $\text{Vech}(\varpi)$ elements are 41.8%, 33.6%, and 29.2%, respectively. Overall, it is difficult to both accurately and precisely recover the effects of exogenous variables on the latent variables (in the structural equation system) as well as the effects of the latent variables on the outcomes (in the measurement equation system). The suggestion is the exercise of caution when GHDM models are being estimated with few observations. Our results suggest that there may be a need for 3000 observations or so for good accuracy and precision in the estimated coefficients. Of course, the situation is likely to be context-

specific, but our simulation analysis does provide some guidance. Interestingly, the FSSE values as a percentage of true values are also rather high for the effects of the exogenous and non-nominal endogenous variables on the utility functions of the nominal variables (that is, the elements of the \mathbf{b} matrix). The FSSE values are 45.4%, 30.5%, and 30.7% for the 1000, 2000, and 3000 observation cases, respectively. This is a case where the APB is very low (accuracy is high) for the elements of the \mathbf{b} matrix, but the precision of estimates is not very good. The relatively poor precision of estimates in the nominal variable equation is not all that surprising, given that multiple latent variables (corresponding to the utilities of alternatives) are used to characterize a nominal outcome, unlike the case of the non-nominal outcomes where a single underlying (observed or latent) variable is used to characterize the observed outcomes. Fifth, the asymptotic formula of the CML approach performs reasonably well in estimating the FSSEs, based on the APBASE values. The mean APBASE values are 25.02%, 16.20%, and 22.69%. While these may not seem small, one should keep in mind that the FSSE values themselves are quite small, leading to rather high APBASE values even if the ASE value is close to the FSSE value in actual magnitude. Further, the APBASE value does not show a decrease as the number of observations increases because the FSSE value itself keeps decreasing as the number of observations increase. In general, the FSSE and the ASE values are not too different from one another regardless of sample size, indicating that the asymptotic formula is performing quite well in estimating the finite sample standard error even for a sample size of the order of 1000. Finally, the APERR in the last column of all three tables indicates that even a single permutation (for each observation) of the approximation approach used to evaluate the MVNCD function provides adequate precision. For the case with 1000 observations, the values of the APERR range between 0.00007 and 0.00721, and the mean APERR is 0.00124. At $Q=2000$, the minimum and maximum APERR values are 0.00010 and 0.00604, respectively, with the mean APERR decreasing to 0.00083. When $Q=3000$, the minimum and maximum APERR values are 0.00004 and 0.00150, respectively, with the mean APERR decreasing further to 0.00032. More importantly, the approximation error (as a percentage of the FSSE), averaged across all the parameters, is of the order of 0.73%, 0.75%, and 0.37% for 1000, 2000, and 3000 observations, respectively. This is clear evidence that the convergent values are about the same for a given data set regardless of the permutation used for the decomposition of the multivariate probability expression.

Table 2: Simulation Results for the 1000-Observations Case with 200 Datasets

Parameters	True Value	Parameters Estimates			Standard Error Estimates					
		Mean Est.	Abs. Bias	APB	FSSE		ASE		APBASE (%)	APERR
					Value	% of true value	Value	% of true value		
α_{11}	0.80	0.790	0.010	1.263	0.160	20.000	0.142	17.750	11.250	0.00082
α_{12}	-0.30	-0.297	0.003	1.006	0.135	45.000	0.094	31.333	30.370	0.00091
α_{22}	0.20	0.147	0.053	26.418	0.126	63.000	0.094	47.000	25.397	0.00074
α_{23}	0.50	0.357	0.143	28.663	0.158	31.600	0.104	20.800	34.177	0.00088
l_{Γ}	-0.60	-0.517	0.083	13.833	0.322	53.667	0.218	36.333	32.298	0.00150
γ_{11}	1.00	1.059	0.059	5.900	0.063	6.300	0.116	11.600	84.127	0.00014
γ_{12}	0.50	0.411	0.089	17.742	0.067	13.400	0.118	23.600	76.119	0.00022
γ_{18}	-0.30	-0.244	0.056	18.505	0.061	20.333	0.052	17.333	14.754	0.00019
$\tilde{\gamma}_{11}$	1.00	0.865	0.135	13.500	0.121	12.100	0.101	10.100	16.529	0.00035
$\tilde{\gamma}_{14}$	-0.20	-0.201	0.001	0.587	0.035	17.500	0.040	20.000	14.286	0.00016
$\tilde{\gamma}_{18}$	0.60	0.606	0.006	1.069	0.102	17.000	0.095	15.833	6.863	0.00041
$\tilde{\gamma}_{21}$	1.00	0.836	0.164	16.400	0.064	6.400	0.039	3.900	39.063	0.00014
$\tilde{\gamma}_{28}$	0.20	0.197	0.003	1.721	0.069	34.500	0.072	36.000	4.348	0.00017
$\tilde{\gamma}_{31}$	1.00	0.847	0.153	15.300	0.112	11.200	0.100	10.000	10.714	0.00010
$\tilde{\gamma}_{34}$	0.40	0.423	0.023	5.650	0.050	12.500	0.043	10.750	14.000	0.00015
$\tilde{\gamma}_{35}$	-0.30	-0.315	0.015	4.868	0.043	14.333	0.039	13.000	9.302	0.00007
$\tilde{\gamma}_{11}$	1.00	0.875	0.125	12.500	0.136	13.600	0.093	9.300	31.618	0.00043
$\tilde{\gamma}_{18}$	-0.50	-0.535	0.035	7.099	0.090	18.000	0.067	13.400	25.556	0.00033
b_{111}	0.20	0.197	0.003	1.438	0.153	76.500	0.115	57.500	24.837	0.00160
b_{112}	0.40	0.398	0.002	0.395	0.125	31.250	0.101	25.250	19.200	0.00098
b_{113}	-0.50	-0.491	0.009	1.700	0.134	26.800	0.112	22.400	16.418	0.00127
b_{121}	0.30	0.320	0.020	6.664	0.172	57.333	0.134	44.667	22.093	0.00074
b_{124}	0.20	0.190	0.010	5.242	0.069	34.500	0.063	31.500	8.696	0.00036
b_{125}	0.30	0.291	0.009	3.034	0.107	35.667	0.090	30.000	15.888	0.00044
b_{221}	-0.50	-0.513	0.013	2.575	0.123	24.600	0.090	18.000	26.829	0.00057
b_{222}	0.30	0.300	0.000	0.105	0.097	32.333	0.075	25.000	22.680	0.00086
b_{228}	0.20	0.215	0.015	7.303	0.100	50.000	0.071	35.500	29.000	0.00071
b_{231}	-0.20	-0.197	0.003	1.595	0.160	80.000	0.134	67.000	16.250	0.00414

Table 2 (Cont.): Simulation Results for the 1000-Observations Case with 200 Datasets

Parameters	True Value	Parameters Estimates			Standard Error Estimates					
		Mean Est.	Abs. Bias	APB	FSSE		ASE		APBASE (%)	APERR
					Value	% of true value	Value	% of true value		
b_{236}	-0.60	-0.591	0.009	1.481	0.287	47.833	0.345	57.500	20.209	0.00201
b_{237}	-0.40	-0.405	0.005	1.329	0.193	48.250	0.240	60.000	24.352	0.00157
d_{12}	0.20	0.173	0.027	13.280	0.061	30.500	0.043	21.500	29.508	0.00058
\tilde{d}_{11}	0.60	0.639	0.039	6.544	0.187	31.167	0.147	24.500	21.390	0.00093
\tilde{d}_{21}	0.20	0.213	0.013	6.261	0.070	35.000	0.078	39.000	11.429	0.00226
\tilde{d}_{32}	0.30	0.442	0.142	47.373	0.173	57.667	0.127	42.333	26.590	0.00083
\tilde{d}_{11}	-0.50	-0.435	0.065	12.970	0.133	26.600	0.096	19.200	27.820	0.00078
\tilde{d}_{12}	0.50	0.703	0.203	40.513	0.322	64.400	0.191	38.200	40.683	0.00059
ϖ_{111}	0.40	0.406	0.006	1.429	0.120	30.000	0.134	33.500	11.667	0.00428
ϖ_{212}	0.20	0.267	0.067	33.500	0.169	84.500	0.096	48.000	43.195	0.00115
ϖ_{221}	0.40	0.424	0.024	5.899	0.129	32.250	0.120	30.000	6.977	0.00262
ϖ_{231}	0.60	0.653	0.053	8.900	0.301	50.167	0.325	54.167	7.973	0.00263
l_{Σ}	1.25	1.049	0.201	16.080	0.042	3.360	0.047	3.760	11.905	0.00040
ψ_{12}	1.50	1.472	0.028	1.894	0.158	10.533	0.119	7.933	24.684	0.00075
ψ_{22}	1.50	1.453	0.047	3.119	0.064	4.267	0.038	2.533	40.625	0.00089
ψ_{32}	1.50	1.524	0.024	1.631	0.152	10.133	0.102	6.800	32.895	0.00035
φ_1	0.75	0.703	0.047	6.202	0.161	21.467	0.087	11.600	45.963	0.00026
θ	2.00	1.680	0.320	16.000	0.719	35.950	0.347	17.350	51.739	0.00062
$l_{\Lambda 32}$	0.70	0.715	0.015	2.213	0.302	43.143	0.231	33.000	23.510	0.00235
$l_{\Lambda 33}$	1.49	1.577	0.087	5.871	1.003	67.315	0.505	33.893	49.651	0.00549
$l_{\Lambda 65}$	0.60	0.604	0.004	0.632	0.417	69.500	0.380	63.333	8.873	0.00721
$l_{\Lambda 66}$	1.36	1.481	0.121	8.894	1.046	76.912	0.976	71.765	6.692	0.00392
Overall mean value across parameters			0.056	9.28	0.187	34.09	0.148	28.49	25.02	0.00124

Table 3: Simulation Results for the 2000-Observations Case with 200 Datasets

Parameters	True Value	Parameters Estimates			Standard Error Estimates					
		Mean Est.	Abs. Bias	APB	FSSE		ASE		APBASE (%)	APERR
					Value	% of true value	Value	% of true value		
α_{11}	0.80	0.859	0.059	7.417	0.132	16.500	0.121	15.125	8.333	0.00060
α_{12}	-0.30	-0.303	0.003	1.135	0.105	35.000	0.072	24.000	31.429	0.00060
α_{22}	0.20	0.160	0.040	19.908	0.083	41.500	0.064	32.000	22.892	0.00030
α_{23}	0.50	0.347	0.153	30.681	0.116	23.200	0.070	14.000	39.655	0.00037
l_{Γ}	-0.60	-0.552	0.048	8.000	0.234	39.000	0.182	30.333	22.222	0.00060
γ_{11}	1.00	1.066	0.066	6.600	0.048	4.800	0.042	4.200	12.500	0.00039
γ_{12}	0.50	0.407	0.093	18.571	0.047	9.400	0.043	8.600	8.511	0.00043
γ_{18}	-0.30	-0.255	0.045	14.998	0.048	16.000	0.044	14.667	8.333	0.00021
$\tilde{\gamma}_{11}$	1.00	0.851	0.149	14.900	0.083	8.300	0.069	6.900	16.867	0.00028
$\tilde{\gamma}_{14}$	-0.20	-0.192	0.008	4.002	0.027	13.500	0.016	8.000	40.741	0.00016
$\tilde{\gamma}_{18}$	0.60	0.572	0.028	4.608	0.070	11.667	0.063	10.500	10.000	0.00032
$\tilde{\gamma}_{21}$	1.00	0.876	0.124	12.400	0.045	4.500	0.028	2.800	37.778	0.00012
$\tilde{\gamma}_{28}$	0.20	0.191	0.009	4.429	0.049	24.500	0.051	25.500	4.082	0.00011
$\tilde{\gamma}_{31}$	1.00	0.856	0.144	14.400	0.073	7.300	0.068	6.800	6.849	0.00011
$\tilde{\gamma}_{34}$	0.40	0.407	0.007	1.713	0.029	7.250	0.028	7.000	3.448	0.00011
$\tilde{\gamma}_{35}$	-0.30	-0.306	0.006	1.944	0.027	9.000	0.026	8.667	3.704	0.00010
$\tilde{\gamma}_{11}$	1.00	0.852	0.148	14.800	0.093	9.300	0.065	6.500	30.108	0.00026
$\tilde{\gamma}_{18}$	-0.50	-0.528	0.028	5.560	0.069	13.800	0.046	9.200	33.333	0.00016
b_{111}	0.20	0.193	0.007	3.699	0.142	71.000	0.135	67.500	4.930	0.00099
b_{112}	0.40	0.394	0.006	1.518	0.083	20.750	0.070	17.500	15.663	0.00081
b_{113}	-0.50	-0.497	0.003	0.546	0.090	18.000	0.078	15.600	13.333	0.00079
b_{121}	0.30	0.305	0.005	1.548	0.105	35.000	0.093	31.000	11.429	0.00074
b_{124}	0.20	0.195	0.005	2.424	0.043	21.500	0.044	22.000	2.326	0.00036
b_{125}	0.30	0.301	0.001	0.406	0.059	19.667	0.064	21.333	8.475	0.00035
b_{221}	-0.50	-0.517	0.017	3.445	0.081	16.200	0.061	12.200	24.691	0.00091
b_{222}	0.30	0.297	0.003	0.849	0.059	19.667	0.052	17.333	11.864	0.00042
b_{228}	0.20	0.201	0.001	0.524	0.059	29.500	0.049	24.500	16.949	0.00043
b_{231}	-0.20	-0.223	0.023	11.265	0.139	69.500	0.142	71.000	2.158	0.00145

Table 3 (Cont.): Simulation Results for the 2000-Observations Case with 200 Datasets

Parameters	True Value	Parameters Estimates			Standard Error Estimates					
		Mean Est.	Abs. Bias	APB	FSSE		ASE		APBASE (%)	APERR
					Value	% of true value	Value	% of true value		
b_{236}	-0.60	-0.612	0.012	2.011	0.141	23.500	0.148	24.667	4.965	0.00158
b_{237}	-0.40	-0.408	0.008	1.983	0.085	21.250	0.099	24.750	16.471	0.00132
d_{12}	0.20	0.164	0.036	18.018	0.036	18.000	0.029	14.500	19.444	0.00027
\tilde{d}_{11}	0.60	0.565	0.035	5.802	0.130	21.667	0.100	16.667	23.077	0.00089
\tilde{d}_{21}	0.20	0.192	0.008	4.200	0.057	28.500	0.053	26.500	7.018	0.00123
\tilde{d}_{32}	0.30	0.419	0.119	39.685	0.088	29.333	0.084	28.000	4.545	0.00068
\tilde{d}_{11}	-0.50	-0.394	0.106	21.157	0.079	15.800	0.063	12.600	20.253	0.00046
\tilde{d}_{12}	0.50	0.664	0.164	32.737	0.187	37.400	0.129	25.800	31.016	0.00062
ϖ_{111}	0.40	0.396	0.004	0.947	0.084	21.000	0.092	23.000	9.524	0.00213
ϖ_{212}	0.20	0.245	0.045	22.500	0.134	67.000	0.118	59.000	11.940	0.00153
ϖ_{221}	0.40	0.386	0.014	3.470	0.083	20.750	0.081	20.250	2.410	0.00218
ϖ_{231}	0.60	0.593	0.007	1.222	0.152	25.333	0.123	20.500	19.079	0.00168
l_{Σ}	1.25	1.099	0.151	12.080	0.028	2.240	0.033	2.640	17.857	0.00012
ψ_{12}	1.50	1.415	0.085	5.680	0.098	6.533	0.076	5.067	22.449	0.00064
ψ_{22}	1.50	1.447	0.053	3.535	0.042	2.800	0.043	2.867	2.381	0.00040
ψ_{32}	1.50	1.501	0.001	0.055	0.067	4.467	0.065	4.333	2.985	0.00038
φ_1	0.75	0.697	0.053	7.117	0.092	12.267	0.057	7.600	38.043	0.00017
θ	2.00	1.764	0.236	11.800	0.398	19.900	0.167	8.350	58.040	0.00055
$l_{\Lambda 32}$	0.70	0.704	0.004	0.546	0.159	22.714	0.158	22.571	0.629	0.00120
$l_{\Lambda 33}$	1.49	1.512	0.022	1.458	0.486	32.617	0.564	37.852	16.049	0.00313
$l_{\Lambda 65}$	0.60	0.621	0.021	3.429	0.227	37.833	0.248	41.333	9.251	0.00186
$l_{\Lambda 66}$	1.36	1.468	0.108	7.945	0.555	40.809	0.665	48.897	19.820	0.00604
Overall mean value across parameters			0.050	8.39	0.113	22.54	0.102	20.25	16.20	0.00083

Table 4: Simulation Results for the 3000-Observations Case with 200 Datasets

Parameters	True Value	Parameters Estimates			Standard Error Estimates					
		Mean Est.	Abs. Bias	APB	FSSE		ASE		APBASE (%)	APERR
					Value	% of true value	Value	% of true value		
α_{11}	0.80	0.763	0.037	4.641	0.081	10.125	0.071	8.875	12.346	0.00027
α_{12}	-0.30	-0.275	0.025	8.452	0.071	23.667	0.050	16.667	29.577	0.00027
α_{22}	0.20	0.194	0.006	3.000	0.061	30.500	0.048	24.000	21.311	0.00017
α_{23}	0.50	0.432	0.068	13.600	0.090	18.000	0.052	10.400	42.222	0.00013
l_{Γ}	-0.60	-0.583	0.017	2.833	0.087	14.500	0.115	19.167	32.184	0.00037
γ_{11}	1.00	1.068	0.068	6.800	0.036	3.600	0.033	3.300	8.333	0.00004
γ_{12}	0.50	0.406	0.094	18.792	0.036	7.200	0.035	7.000	2.778	0.00005
γ_{18}	-0.30	-0.249	0.051	16.857	0.039	13.000	0.036	12.000	7.692	0.00006
$\tilde{\gamma}_{11}$	1.00	0.957	0.043	4.300	0.067	6.700	0.059	5.900	11.940	0.00015
$\tilde{\gamma}_{14}$	-0.20	-0.200	0.000	0.057	0.023	11.500	0.014	7.000	39.130	0.00010
$\tilde{\gamma}_{18}$	0.60	0.606	0.006	1.008	0.068	11.333	0.055	9.167	19.118	0.00017
$\tilde{\gamma}_{21}$	1.00	0.963	0.037	3.700	0.040	4.000	0.037	3.700	7.500	0.00005
$\tilde{\gamma}_{28}$	0.20	0.201	0.001	0.332	0.040	20.000	0.042	21.000	5.000	0.00004
$\tilde{\gamma}_{31}$	1.00	0.965	0.035	3.500	0.059	5.900	0.055	5.500	6.780	0.00005
$\tilde{\gamma}_{34}$	0.40	0.412	0.012	3.114	0.024	6.000	0.014	3.500	41.667	0.00004
$\tilde{\gamma}_{35}$	-0.30	-0.310	0.010	3.392	0.022	7.333	0.013	4.333	40.909	0.00004
$\tilde{\gamma}_{11}$	1.00	0.956	0.044	4.400	0.090	9.000	0.055	5.500	38.889	0.00013
$\tilde{\gamma}_{18}$	-0.50	-0.527	0.027	5.323	0.058	11.600	0.038	7.600	34.483	0.00010
b_{111}	0.20	0.211	0.011	5.429	0.145	72.500	0.110	55.000	24.138	0.00030
b_{112}	0.40	0.394	0.006	1.576	0.075	18.750	0.058	14.500	22.667	0.00034
b_{113}	-0.50	-0.495	0.005	0.914	0.084	16.800	0.065	13.000	22.619	0.00044
b_{121}	0.30	0.317	0.017	5.553	0.101	33.667	0.075	25.000	25.743	0.00020
b_{124}	0.20	0.194	0.006	3.239	0.047	23.500	0.037	18.500	21.277	0.00014
b_{125}	0.30	0.294	0.006	2.053	0.065	21.667	0.052	17.333	20.000	0.00016
b_{221}	-0.50	-0.512	0.012	2.379	0.066	13.200	0.049	9.800	25.758	0.00035
b_{222}	0.30	0.297	0.003	1.094	0.051	17.000	0.043	14.333	15.686	0.00015
b_{228}	0.20	0.205	0.005	2.471	0.052	26.000	0.040	20.000	23.077	0.00022
b_{231}	-0.20	-0.206	0.006	3.011	0.146	73.000	0.188	94.000	28.767	0.00038

Table 4 (Cont.): Simulation Results for the 3000-Observations Case with 200 Datasets

Parameters	True Value	Parameters Estimates			Standard Error Estimates					
		Mean Est.	Abs. Bias	APB	FSSE		ASE		APBASE (%)	APERR
					Value	% of true value	Value	% of true value		
b_{236}	-0.60	-0.609	0.009	1.537	0.159	26.500	0.201	33.500	26.415	0.00037
b_{237}	-0.40	-0.413	0.013	3.272	0.105	26.250	0.137	34.250	30.476	0.00041
d_{12}	0.20	0.155	0.045	22.334	0.036	18.000	0.023	11.500	36.111	0.00018
\tilde{d}_{11}	0.60	0.668	0.068	11.286	0.119	19.833	0.087	14.500	26.891	0.00060
\tilde{d}_{21}	0.20	0.221	0.021	10.576	0.047	23.500	0.027	13.500	42.553	0.00052
\tilde{d}_{32}	0.30	0.352	0.052	17.333	0.078	26.000	0.064	21.333	17.949	0.00035
\tilde{d}_{11}	-0.50	-0.426	0.074	14.881	0.071	14.200	0.052	10.400	26.761	0.00022
\tilde{d}_{12}	0.50	0.578	0.078	15.600	0.196	39.200	0.117	23.400	40.306	0.00018
ϖ_{111}	0.40	0.423	0.023	5.775	0.083	20.750	0.077	19.250	7.229	0.00090
ϖ_{212}	0.20	0.164	0.036	18.000	0.110	55.000	0.100	50.000	9.091	0.00051
ϖ_{221}	0.40	0.436	0.036	9.100	0.067	16.750	0.068	17.000	1.493	0.00086
ϖ_{231}	0.60	0.664	0.064	10.593	0.145	24.167	0.107	17.833	26.207	0.00100
l_{Σ}	1.25	1.119	0.131	10.480	0.025	2.000	0.027	2.160	8.000	0.00007
ψ_{12}	1.50	1.481	0.019	1.253	0.090	6.000	0.071	4.733	21.111	0.00038
ψ_{22}	1.50	1.450	0.050	3.307	0.032	2.133	0.036	2.400	12.500	0.00018
ψ_{32}	1.50	1.511	0.011	0.748	0.063	4.200	0.057	3.800	9.524	0.00015
φ_1	0.75	0.703	0.047	6.275	0.088	11.733	0.050	6.667	43.182	0.00008
θ	2.00	1.855	0.145	7.250	0.173	8.650	0.142	7.100	17.919	0.00017
$l_{\Lambda 32}$	0.70	0.718	0.018	2.528	0.165	23.571	0.128	18.286	22.424	0.00064
$l_{\Lambda 33}$	1.49	1.503	0.013	0.894	0.191	12.819	0.125	8.389	34.555	0.00143
$l_{\Lambda 65}$	0.60	0.612	0.012	2.038	0.116	19.333	0.081	13.500	30.172	0.00055
$l_{\Lambda 66}$	1.36	1.465	0.105	7.711	0.242	17.794	0.271	19.926	11.983	0.00150
Overall mean value across parameters			0.035	6.29	0.085	18.97	0.072	16.19	22.69	0.00032

4.6.1 Effects of Ignoring Latent Construct Effects

This section presents the results of the estimation when the latent variables are ignored, and the resulting dependencies among the multidimensional outcomes are not considered. As discussed earlier in the first part of Section 4, this is equivalent to ignoring all potential self-selection effects, which then should corrupt all endogenous variable effects discussed in Section 4.3.3, and lead to inaccurate and inefficient estimation of other parameters as well. Ignoring the presence of latent variables is tantamount to the restriction in the GHDM model that all elements of the \vec{d} matrix and the $\vec{\omega}$ matrix in Equation (15) are zero (no effects of latent variables on any (and all) outcome(s)). But doing so immediately renders all elements of α and Γ unidentifiable, because the only way these elements are identified is by the relationship between the latent variable vector z^* and the observed outcomes. Thus, we also essentially are setting all elements of α and Γ to zero in the restricted model. The resulting equivalent of Equation (15), which we will refer to as the independent model for ease, can be compared with the GHDM model using the adjusted composite log-likelihood ratio test (ADCLRT) value (see Pace *et al.*, 2011 and Bhat, 2011 for more details on the ADCLRT statistic, which is the equivalent of the log-likelihood ratio test statistic when a composite marginal likelihood inference approach is used; this statistic has an approximate chi-squared asymptotic distribution).

For the comparison of the GHDM and independent model coefficient estimates (vis-à-vis the true values of the experimental design), we estimate the independent model on the same 200 datasets as we estimated the GHDM model on earlier. Based on the results for the GHDM model, we decided to undertake this comparison only for the case of $Q=3000$ observations. For each of the 200 data sets, we use the same set of permutations for the joint model and the independent model, so that we are able to appropriately compare the ability to recover parameters from the two models. We made this comparison between the two models only for those coefficients estimated in the independent model. The GHDM model mean APB is 4.19 relative to the independent model mean APB of 16.03 (the complete table results are available from the author). In addition to an APB comparison between the joint model and the independent model, we also compare the performance of the two models using the ADCLRT test. The ADCLRT statistic for the test between the two models has an approximate chi-squared distribution with 15 degrees of freedom. The corresponding table value for the chi-squared distribution is 32.8 at the 0.5% level of significance. In this paper, we identify the number of times (corresponding to the 200 data sets) that the ADCLRT value rejects the independent model in favor of the joint model. The result indicates that the joint model rejects the independent model in all the 200 data sets, further reinforcing the need to consider the GHDM model.

4.7 Procedure for Treatment Effects Based on Residential Choice

The estimation results from the simulation experiment may be used to examine the differences between the GHDM and independent models as they relate to the implied effects of one outcome variable on another. To demonstrate the potential problems of ignoring latent variables, we examine the impact of residential location choice on auto ownership (other outcome effects may also be computed, but, because this is only a

simulation effort, we focus on one effect to demonstrate the potential biases accruing from ignoring jointness). This is helpful to obtain insights regarding whether, and how much, an independent model can bias the influence of an urban-like high density design on travel-related behaviors. An important approach to do so is the Average Treatment Effect (ATE) (see Heckman and Vytlačil, 2000 and Heckman *et al.*, 2001).

In the context of motorized vehicle ownership, the ATE measure provides the expected difference in motorized vehicle ownership for a random individual if s/he were located in a specific density configuration i as opposed to another density configuration $i' \neq i$. The measure is estimated as follows:

$$\hat{ATE}_{ii'} = \frac{1}{Q} \sum_{q=1}^Q \left(\sum_{j_1=0}^{\infty} k_{q1} \cdot [P(k_{q1} | a_{qi} = 1) - P(k_{q1} | a_{qi'} = 1)] \right),$$

where a_{qi} is the dummy variable for the density category i for the individual q , and k_{q1} is an index for auto ownership k_{q1} ($k_{q1} = 0, 1, 2, \dots, \infty$) (the subscript '1', consistent with the notation used earlier, indicates that auto ownership is the first count variable in the model system). Although the summation in the equation above extends until infinity, we consider counts only up to $k_{q1} = 10$. This should not affect the computations because the probabilities associated with higher motorized vehicle ownership levels are very close to zero.

The analyst can compute the ATE measures for all the pairwise combinations of residential density category relocations. Here, we focus on the case when an individual in a rural location is transplanted to an urban location. The standard error of the ATE measure is obtained using bootstraps from the sampling distributions of the estimated parameters. The GHDM model estimates an ATE of -0.178 (standard error of 0.013), which implies that a random household that is shifted from a rural location to an urban location will, on average, reduce its motorized vehicle ownership level by 0.178 vehicles. The corresponding independent model estimate is much higher with an ATE of -0.338 (standard error of 0.011), which indicates a much higher reduction in auto ownership because of a household move from a rural area to an urban area. This overestimation in the independent model is because the probability of residing in an urban area and the propensity to own autos are negatively correlated because of the latent green lifestyle propensity (GLP) latent construct (note that, in Figure 2b, GLP has a positive effect on the utility of residing in an urban area, and, in Figure 2a, GLP has a negative effect on auto ownership propensity). If this GLP construct is ignored (as in the independent model), the result is a transfer of the negative covariance due to the GLP construct to a much higher negative (and biased) ATE of urban dwelling on auto ownership count. Thus, accounting for endogeneity effects is not simply of academic interest, but can have substantial real implications for variable effects and subsequent policy analysis.

References

- Aditjandra, P. T., Cao, X. J., and Mulley, C. (2012). Understanding neighbourhood design impact on travel behaviour: An application of structural equations model to a British metropolitan data. *Transportation Research Part A*, 46(1), 22-32.
- Bartholomew, K.J., Ntoumanis, N., Ryan, R.M., Bosch, J.A., and Thøgersen-Ntoumani, C. (2011). Self-determination theory and diminished functioning the role of interpersonal control and psychological need thwarting. *Personality and Social Psychology Bulletin*, 37(11), 1459-1473.
- Ben-Akiva, M., McFadden, D., Train, K., Walker, J., Bhat, C., Bierlaire, M., . . . Bunch, D.S. (2002). Hybrid choice models: Progress and challenges. *Marketing Letters*, 13(3), 163-175.
- Bhat, C.R. (2011). The maximum approximate composite marginal likelihood (MACML) estimation of multinomial probit-based unordered response choice models. *Transportation Research Part B*, 45(7), 923-939.
- Bhat, C.R. (2014). The composite marginal likelihood (CML) inference approach with applications to discrete and mixed dependent variable models. *Foundations and Trends in Econometrics*, 7(1), 1-117.
- Bhat, C.R., and Dubey, S.K. (2014). A new estimation approach to integrate latent psychological constructs in choice modeling. *Transportation Research Part B*, 67, 68-85.
- Bhat, C.R., and Guo, J.Y. (2007). A comprehensive analysis of built environment characteristics on household residential choice and auto ownership levels. *Transportation Research Part B*, 41(5), 506-526.
- Bhat, C.R., and Sidharthan, R. (2011). A simulation evaluation of the maximum approximate composite marginal likelihood (MACML) estimator for mixed multinomial probit models. *Transportation Research Part B*, 45(7), 940-953.
- Bhat, C.R., Paleti, R., Pendyala, R.M., Lorenzini, K., and Konduri, K.C. (2013). Accommodating immigration status and self-selection effects in a joint model of household auto ownership and residential location choice. *Transportation Research Record: Journal of the Transportation Research Board*, 2382(1), 142-150.
- Bhat, C.R., Astroza, S., Sidharthan, R., Jobair Bin Alam, M., and Khushefati, W.H. (2014a). A joint count-continuous model of travel behavior with selection based on a multinomial probit residential density choice model. *Transportation Research Part B*, 68, 31-51.
- Bhat, C.R., Paleti, R., and Singh, P. (2014b). A spatial multivariate count model for firm location decisions. *Journal of Regional Science*, 54(3), 462-502.
- Bolduc, D., Ben-Akiva, M., Walker, J., Michaud, A., 2005. Hybrid choice models with logit kernel: applicability to large scale models. In: Lee-Gosselin, M., Doherty, S. (eds.) *Integrated Land-Use and Transportation Models: Behavioral Foundations*, Elsevier, Oxford, 275-302.
- Bollen, K. A. (1989). *Structural Equations with Latent Variables*, Wiley Series in Probability and Mathematical Statistics, John Wiley & Sons.

- Brownstone, D., and Golob, T.F. (2009). The impact of residential density on vehicle usage and energy consumption. *Journal of Urban Economics*, 65(1), 91-98.
- Cao, X., and Fan, Y. (2012). Exploring the influences of density on travel behavior using propensity score matching. *Environment and Planning-Part B*, 39(3), 459.
- Castro, M., Paleti, R., and Bhat, C.R. (2012). A latent variable representation of count data models to accommodate spatial and temporal dependence: Application to predicting crash frequency at intersections. *Transportation Research Part B*, 46(1), 253-272.
- Castro, M., Eluru, N., Bhat, C.R., and Pendyala, R.M. (2011). Joint model of participation in nonwork activities and time-of-day choice set formation for workers. *Transportation Research Record: Journal of the Transportation Research Board*, 2254, 140-150.
- Clark, W. A., Huang, Y., and Withers, S. (2003). Does commuting distance matter?: Commuting tolerance and residential change. *Regional Science and Urban Economics*, 33(2), 199-221.
- Day, L.L. (2000). Choosing a house: the relationship between dwelling type, perception of privacy and residential satisfaction. *Journal of Planning Education and Research*, 19(3), 265-275.
- Daziano, R.A., and Bolduc, D. (2013). Incorporating pro-environmental preferences towards green automobile technologies through a Bayesian hybrid choice model. *Transportmetrica A: Transport Science*, 9(1), 74-106.
- De Leon, A.R., and Carrière, K. (2007). General mixed-data model: Extension of general location and grouped continuous models. *Canadian Journal of Statistics*, 35(4), 533-548.
- De Leon, A.R., and Chough, K.C. (2013). *Analysis of Mixed Data: Methods & Applications*, CRC Press, Taylor & Francis Group, Boca Raton, FL.
- De Leon, A.R., and Zhu, Y. (2008). ANOVA extensions for mixed discrete and continuous data. *Computational Statistics & Data Analysis*, 52(4), 2218-2227.
- De Leon, A., Soo, A., and Williamson, T. (2011). Classification with discrete and continuous variables via general mixed-data models. *Journal of Applied Statistics*, 38(5), 1021-1032.
- Faes, C., Geys, H., and Catalano, P. (2009). Joint models for continuous and discrete longitudinal data. *Longitudinal Data Analysis*, 327-348.
- Feddag, M.-L. (2013). Composite likelihood estimation for multivariate probit latent traits models. *Communications in Statistics-Theory and Methods*, 42(14), 2551-2566.
- Gates, K.M., Molenaar, P., Hillary, F.G., and Slobounov, S. (2011). Extended unified SEM approach for modeling event-related fMRI data. *NeuroImage*, 54(2), 1151-1158.
- Gifford, R., Nilsson, A. (2014). Personal and social factors that influence pro environmental concern and behaviour: A review. *International Journal of Psychology*, 49(3), 141-157.
- Godambe, V.P., 1960. An optimum property of regular maximum likelihood estimation. *The Annals of Mathematical Statistics*, 31(4), 1208-1211.
- Gueorguieva, R., and Sanacora, G. (2006). Joint analysis of repeatedly observed continuous and ordinal measures of disease severity. *Statistics in Medicine*, 25(8), 1307-1322.

- Heckman, J.J., and Vytlacil, E.J. (2000). The relationship between treatment parameters within a latent variable framework. *Economics Letters*, 66(1), 33-39.
- Heckman, J., Tobias, J.L., and Vytlacil, E. (2001). Four parameters of interest in the evaluation of social programs. *Southern Economic Journal*, 211-223.
- Hoshino, T., and Bentler, P.M. (2013). Bias in factor score regression and a simple solution. In: De Leon, A.R., and Chough, K.C. (eds.), *Analysis of Mixed Data: Methods & Applications*, CRC Press, Taylor & Francis Group, Boca Raton, FL, 43-61.
- Jansen, S.J. (2012). What is the worth of values in guiding residential preferences and choices? *Journal of Housing and the built Environment*, 27(3), 273-300.
- Jöreskog, K.G. (1977). Factor analysis by least squares and maximum likelihood methods. In: Enslein, K., Ralston, A., and Wilf, H.S. (eds), *Statistical Methods for Digital Computers*, John Wiley & Sons, New York.
- Keane, M.P. (1992). A note on identification in the multinomial probit model. *Journal of Business & Economic Statistics*, 10(2), 193-200.
- Kim, J., and Brownstone, D. (2013). The impact of residential density on vehicle usage and fuel consumption: Evidence from national samples. *Energy Economics*, 40, 196-206.
- Liu, X., Vedlitz, A., Shi, L. (2014). Examining the determinants of public environmental concern: Evidence from national public surveys. *Environmental Science & Policy*, 39, 77-94.
- Maddala, G. (1983). *Limited-Dependent and Qualitative Variables in Econometrics*, Cambridge University Press, Cambridge, UK.
- Mokhtarian, P.L., and Cao, X. (2008). Examining the impacts of residential self-selection on travel behavior: A focus on methodologies. *Transportation Research Part B*, 42(3), 204-228.
- Munkin, M.K., and Trivedi, P.K. (2008). Bayesian analysis of the ordered probit model with endogenous selection. *Journal of Econometrics*, 143(2), 334-348.
- O'Brien, R.M. (1994). Identification of simple measurement models with multiple latent variables and correlated errors. *Sociological Methodology*, 24, 137-170.
- Pace, L., Salvani, A., and Sartori, N. (2011). Adjusting composite likelihood ratio statistics. *Statistica Sinica*, 21(1), 129.
- Paleti, R., Bhat, C.R., and Pendyala, R.M. (2013). Integrated Model of Residential Location, Work Location, Vehicle Ownership, and Commute Tour Characteristics. *Transportation Research Record: Journal of the Transportation Research Board*, 2382, 162-172.
- Pinjari, A. R., Eluru, N., Bhat, C.R., Pendyala, R.M., and Spissu, E. (2008). Joint model of choice of residential neighborhood and bicycle ownership: accounting for self-selection and unobserved heterogeneity. *Transportation Research Record: Journal of the Transportation Research Board*, 2082, 17-26.
- Pinjari, A.R., Pendyala, R.M., Bhat, C.R., and Waddell, P.A. (2011). Modeling the choice continuum: an integrated model of residential location, auto ownership, bicycle ownership, and commute tour mode choice decisions. *Transportation*, 38(6), 933-958.

- Rashidi, T.H., Auld, J., and Mohammadian, A.K. (2012). A behavioral housing search model: Two-stage hazard-based and multinomial logit approach to choice-set formation and location selection. *Transportation Research Part A*, 46(7), 1097-1107.
- Reilly, T., and O'Brien, R.M. (1996). Identification of confirmatory factor analysis models of arbitrary complexity the side-by-side rule. *Sociological Methods & Research*, 24(4), 473-491.
- Schwanen, T., and Mokhtarian, P.L. (2007). Attitudes toward travel and land use and choice of residential neighborhood type: Evidence from the San Francisco bay area. *Housing Policy Debate*, 18(1), 171-207.
- Sener, I.N., Eluru, N., and Bhat, C.R. (2009). An analysis of bicycle route choice preferences in Texas, US. *Transportation*, 36(5), 511-539.
- Shiftan, Y., Outwater, M. L., and Zhou, Y. (2008). Transit market research using structural equation modeling and attitudinal market segmentation. *Transport Policy*, 15(3), 186-195.
- Stapleton, D.C. (1978). Analyzing political participation data with a MIMIC Model. *Sociological Methodology*, 52-74.
- Sutton, A.J., Abrams, K.R., Jones, D.R., Sheldon, T.A., and Song, F. (2000). *Methods for Meta-Analysis in Medical Research*. Wiley Series in Probability and Statistics, John Wiley & Sons, Chichester, England.
- Teixeira-Pinto, A., and Harezlak, J. (2013). Factorization and latent variable models for joint analysis of binary and continuous outcomes. In: De Leon, A.R., and Chough, K.C. (eds.), *Analysis of Mixed Data: Methods & Applications*, CRC Press, Taylor & Francis Group, Boca Raton, FL, 81-91.
- Temme, D., Paulssen, M., and Dannewald, T. (2008). Incorporating latent variables into discrete choice models-A simultaneous estimation approach using SEM software. *Business Research*, 1(2).
- Wu, B., de Leon, A., and Withanage, N. (2013). Joint analysis of mixed discrete and continuous outcomes via copulas. In: De Leon, A.R., and Chough, K.C. (eds.), *Analysis of Mixed Data: Methods & Applications*, CRC Press, Taylor & Francis Group, Boca Raton, FL, 139-156.
- Zhao, Y., and Joe, H. (2005). Composite likelihood estimation in multivariate data analysis. *Canadian Journal of Statistics*, 33(3), 335-356.