**Online supplement for the paper**

**"Sharing the Road with Autonomous Vehicles: Perceived Safety and Regulatory Preferences"**

Gopindra S. Nair, Chandra R. Bhat (corresponding author)

**Mathematical formulation of GHDM model with only ordinal outcomes**

Consider the case of an individual $q \in \{1, 2, ..., Q\}$. Let $l \in \{1, 2, ..., L\}$ be the index of the latent constructs and let $z_{ql}^{*}$ be the value of the latent variable $l$ for the individual $q$. $z_{ql}^{*}$ is expressed as a function of its explanatory variables as,

$$z_{ql}^{*} = \boldsymbol{w}_{ql}^{\mathrm{T}}\boldsymbol{\alpha} + \eta_{ql},\tag{1}$$

where $\boldsymbol{w}_{ql}$ $(D \times 1)$ is a column vector of the explanatory variables of latent variable $l$ and $\boldsymbol{\alpha}$ $(D \times 1)$ is a vector of its coefficients. $\eta_{ql}$ is the unexplained error term and is assumed to follow a standard normal distribution. Equation (1) can be expressed in the matrix form as,

$$\boldsymbol{z}_{q}^{*} = \boldsymbol{w}_{q}\boldsymbol{\alpha} + \boldsymbol{\eta}_{q},\tag{2}$$

where $\boldsymbol{z}_{q}^{*}$ $(L \times 1)$ is a column vector of all the latent variables, $\boldsymbol{w}_{q}$ $(L \times D)$ is a matrix formed by vertically stacking the vectors $(\boldsymbol{w}_{q1}^{\mathrm{T}}, \boldsymbol{w}_{q2}^{\mathrm{T}}, ..., \boldsymbol{w}_{qL}^{\mathrm{T}})$ and $\boldsymbol{\eta}_{q}$ $(D \times 1)$ is formed by vertically stacking $(\eta_{q1}, \eta_{q2}, ..., \eta_{qL})$. $\boldsymbol{\eta}_{q}$ follows a multivariate normal distribution centered at the origin and having a correlation matrix of $\boldsymbol{\Gamma}$ $(L \times L)$, i.e., $\boldsymbol{\eta}_{q} \sim MVN_{L}(\boldsymbol{0}_{L}, \boldsymbol{\Gamma})$, where $\boldsymbol{0}_{L}$ is a vector of zeros. The variance of all the elements in $\boldsymbol{\eta}_{q}$ is fixed as unity because it is not possible to uniquely identify a scale for the latent variables. Equation (2) constitutes the SEM component of the framework.

Let $j \in \{1, 2, ..., J\}$ denote the index of the outcome variables (including the indicator variables). Let $y_{qj}^{*}$ be the underlying continuous measure associated with the outcome variable $y_{qj}$. Then,

$$y_{qj} = k \text{ if } t_{jk} < y_{qj}^{*} \leq t_{j(k+1)},\tag{3}$$

where $k \in \{1, 2, ..., K_{j}\}$ denotes the ordinal category assumed by $y_{qj}$ and $t_{jk}$ denotes the lower boundary of the $k^{\mathrm{th}}$ discrete interval of the continous measure associated with the $j^{\mathrm{th}}$ outcome.

$t_{jk} < t_{j(k+1)}$ for all $j$ and all $k$. Since $y_j^*$ may take any value in $(-\infty, \infty)$, we fix the value of $t_{j1} = -\infty$ and $t_{j(K_j+1)} = \infty$ for all $j$. Since the location of the thresholds on the real-line is not uniquely identifiable, we also set $t_{j2} = 0$. $y_j^*$ is expressed as a function of its explanatory variables as,

$$y_{qj}^* = x_{qj}^{\mathrm{T}}\boldsymbol{\beta} + z_q^{*\mathrm{T}}\boldsymbol{d}_j + \xi_{qj}, \tag{4}$$

where $x_{qj}(E\times1)$ is a vector of size of explanatory variables for the continuous measure $y_{qj}^*$. $\boldsymbol{\beta}$ $(E\times1)$ is a column vector of the coefficients associated with $x_{qj}$ and $\boldsymbol{d}_j$ $(L\times1)$ is the vector of coefficeints of the latent variables for outcome $j$. $\xi_{qj}$ is a stochastic error term that captures the effect of unobserved variables on the value of $y_{qj}^*$. $\xi_{qj}$ is assumed to follow a standard normal distribution. Jointly, the continuous measures of the $J$ outcome variables may be expressed as,

$$\boldsymbol{y}_q^* = x_q\boldsymbol{\beta} + \boldsymbol{d}z_q^* + \boldsymbol{\xi}_q, \tag{5}$$

where $\boldsymbol{y}_q^*$ $(J\times1)$ and $\boldsymbol{\xi}_q$ $(J\times1)$ are the vectors formed by vertically stacking $y_{qj}^*$ and $\xi_{qj}$ respectively of the $J$ dependent variables. $x_q$ $(J\times E)$ is a matrix formed by vertically stacking the vectors $\left(x_{q1}^{\mathrm{T}}, x_{q2}^{\mathrm{T}}, ..., x_{qJ}^{\mathrm{T}}\right)$ and $\boldsymbol{d}$ $(J\times L)$ is a matrix formed by vertically stacking $\left(\boldsymbol{d}_1^{\mathrm{T}}, \boldsymbol{d}_2^{\mathrm{T}}, ..., \boldsymbol{d}_J^{\mathrm{T}}\right)$. $\boldsymbol{\xi}_q$ follows a multivariate normal distribution centered at the origin with an identity matrix as the covariance matrix (independent error terms). $\boldsymbol{\xi}_q \sim MVN_J(\boldsymbol{0}_J, \mathbf{I}_J)$. We assume the terms in $\boldsymbol{\xi}_q$ to be independent because it is not possible to uniquely identify all the correlations between the elements in $\boldsymbol{\eta}_q$ and all the correlations between the elements in $\boldsymbol{\xi}_q$. Further, because of the ordinal nature of the outcome variables, the scale of $\boldsymbol{y}_q^*$ cannot be uniquely identified. Therefore, the variances of all elements in $\boldsymbol{\xi}_q$ is fixed to one. The reader is referred to Bhat (2015) for further nuances regarding the identification of coefficients in the GHDM framework.

Substituting Equation (2) in Equation (5), $\boldsymbol{y}_q^*$ can be expressed in the reduced form as,

$$\boldsymbol{y}_q^* = x_q\boldsymbol{\beta} + \boldsymbol{d}\left(w_q\boldsymbol{\alpha} + \boldsymbol{\eta}_q\right) + \boldsymbol{\xi}_q, \tag{6}$$

$$\boldsymbol{y}_q^* = x_q\boldsymbol{\beta} + \boldsymbol{d}w_q\boldsymbol{\alpha} + \boldsymbol{d}\boldsymbol{\eta}_q + \boldsymbol{\xi}_q. \tag{7}$$

In the R.H.S. of Equation (7), $\eta_q$ and $\xi_q$ are random vectors that follow the multivariate normal distribution and the other variables are constants. Therefore, $y_q^*$ also follows the multivariate normal distribution with a mean of $b = x_q\beta + dw_q\alpha$ (all the elements of $\eta_q$ and $\xi_q$ have a mean of zero) and a covariance matrix of $\Sigma = d\Gamma d^{\mathrm{T}} + I_J$.

$$y_q^* \sim MVN_J(b, \Sigma).\tag{8}$$

The parameters that are to be estimated are the elements of $\alpha$, strictly upper triangular elements of $\Gamma$, elements of $\beta$, elements of $d$ and $t_{jk}$ for all $j$ and $k \in \{3,4,...,K_j\}$. Let $\theta$ be a vector of all the parameters that need to be estimated. The maximum likelihood approach can be used for estimating these parameters. The likelihood of the $q^{\text{th}}$ observation will be,

$$L_q(\theta) = \int_{v_1 = t_{1y_{q1}} - b_1}^{v_1 = t_{1(y_{q1}+1)} - b_1} \int_{v_2 = t_{2y_{q2}} - b_2}^{v_2 = t_{2(y_{q2}+1)} - b_2} \cdots \int_{v_J = t_{Jy_{qJ}} - b_J}^{v_J = t_{J(y_{qJ}+1)} - b_J} \phi_J(v_1, v_2, \ldots, v_J \mid \Sigma) dv_1 dv_2 \ldots dv_J,\tag{9}$$

where, $\phi_J(v_1, v_2, \ldots, v_J \mid \Sigma)$ denotes the probability density of a $J$ dimensional multivariate normal distribution centered at the origin with a covariance matrix $\Sigma$ at the point $(v_1, v_2, \ldots, v_J)$. Since a closed form expression does not exist for this integral and evaluation using simulation techniques can be time consuming, we used the One-variate Univariate Screening technique proposed by Bhat (2018) for approximating this integral. The estimation of parameters was carried out using the *maxlik* library in the GAUSS matrix programming language.

**Output predictions**

To predict the outcome for an individual, first the random error terms for the latent variables are drawn from a multivariate normal distribution that has a mean of $0_L$ and a covariance matrix of $\Gamma$. Following this the latent variables are predicted using Equation (2). Then the random error terms of the outcome variables are drawn from independent standard normal distributions. The outcomes are then predicted in a sequential manner using Equation (3) and Equation (4). Since some of the outcome variables are used as explanatory variables for other outcome variables, the sequence used for estimating the outcomes is such that the outcome that is predicted $i^{\text{th}}$ will only have explanatory variables that are exogenous or that have been predicted in one of the previous $(i-1)$ steps.

## Average treatment effects

The expression for computing the ATE is as follows,

$$A\hat{T}E_{jiAB} = \frac{1}{Q}\sum_{q=1}^{Q}\left(P_q\left(y_j = i \mid B\right) - P_q\left(y_j = i \mid A\right)\right),$$ (10)

where $A\hat{T}E_{jiAB}$ denotes the ATE on the $i^{\text{th}}$ level of the $j^{\text{th}}$ outcome variable by applying a treatment that changes condition of individuals from $A$ to $B$, $P_q(.)$ is a function that computes the probability for the individual $q$. Since the computation of this probability is cumbersome when the outcome variable under consideration is explained by other outcome variables that are also affected by the treatment, we use Monte Carlo methods to simulate the required probabilities. For approximating the probability $P_q\left(y_j = i \mid B\right)$, we make 1000 predictions for the individual $q$ where each prediction is made using a different random draw. The process for making the prediction is the same as that described in the previous section except that the variables that are affected by the condition $B$ is set as per the condition irrespective of the value generated by the prediction process. The estimate for $P_q\left(y_j = i \mid B\right)$ will then be the proportion of the number of cases where $y_j$ assumes the value $i$ to the total number of predictions.

## References

Bhat, C.R. (2015). A new generalized heterogeneous data model (GHDM) to jointly model mixed types of dependent variables. *Transportation Research Part B*, 79, 50–77. https://doi.org/10.1016/j.trb.2015.05.017

Bhat, C.R. (2018). New matrix-based methods for the analytic evaluation of the multivariate cumulative normal distribution function. *Transportation Research Part B*, 109, 238–256. https://doi.org/10.1016/j.trb.2018.01.011

**Latent variable loadings and thresholds of indicators**

| Variable (Minimum level – Maximum level) | Loading | | Constant | | Threshold 2\|3 | | Threshold 3\|4 | |
|---|---|---|---|---|---|---|---|---|
| | Estimate | t-stat | Estimate | t-stat | Estimate | t-stat | Estimate | t-stat |
| Tech-savviness indicators | | | | | | | | |
| Computer use frequency (Never/Almost never – Everyday) | 0.803 | 17.59 | 2.620 | 27.86 | 0.952 | 17.43 | 1.727 | 25.49 |
| Frequency of internet use (Less than once a day – Many times a day) | 0.853 | 18.83 | 1.794 | 24.99 | 1.272 | 27.61 | | |
| Uses laptop (No – Yes) | 0.494 | 15.99 | 0.769 | 20.29 | | | | |
| Uses tablet or e-book reader (No – Yes) | 0.287 | 12.23 | 0.049 | 1.85 | | | | |
| Used voice activated digital assistant in smart devices (No – Yes) | 0.176 | 7.26 | -0.669 | -26.41 | | | | |
| Enthusiasm about riding in / sharing the road with AVs | | | | | | | | |
| Enthusiasm about the development of AVs (Not at all enthusiastic – Very enthusiastic) | 1.850 | 18.30 | 2.574 | 19.82 | 2.420 | 21.78 | 4.599 | 23.23 |
| Anxiety about riding in / sharing road with AVs | | | | | | | | |
| Worried about the development of AVs (Not at all worried – Very worried) | 1.229 | 7.73 | 1.487 | 12.18 | 1.848 | 12.49 | 3.886 | 12.72 |