

1 **A PRACTICAL METHOD TO TEST THE VALIDITY OF THE STANDARD GUMBEL DISTRIBUTION IN LOGIT-**
2 **BASED MULTINOMIAL CHOICE MODELS OF HUMAN TRAVEL BEHAVIOR**

3
4
5 **Xin Ye** (*corresponding author*)

6 Tongji University, College of Transportation Engineering,
7 Key Laboratory of Road and Traffic Engineering of Ministry of Education
8 4800 Cao'an Road, Shanghai, China, 201804.
9 Tel: +86-21-59946091
10 Email: xye@tongji.edu.cn

11
12 **Venu M. Garikapati**

13 Arizona State University, School of Sustainable Engineering and the Built Environment
14 660 S. College Avenue, Tempe, AZ 85281
15 Tel: 480-522-8067; Email: venu.garikapati@asu.edu

16
17 **Daehyun You**

18 Maricopa Association of Governments
19 302 N. First Avenue, Suite 300, Phoenix, AZ 85003
20 Tel: (602) 254-6300; Email: dyou@azmag.gov

21
22 **Ram M. Pendyala**

23 Arizona State University, School of Sustainable Engineering and the Built Environment
24 660 S. College Avenue, Tempe, AZ 85281
25 Tel: 480-965-3589; Email: ram.pendyala@asu.edu

26
27
28
29
30 *Submitted for Presentation only*

31
32 Word count: 6450 text + 6 tables/figures x 250 = 7,950 words

33 96th Annual Meeting of the Transportation Research Board

34
35 Committee on Transportation Demand Forecasting (ADB40)

36
37 August 2016
38

1 **ABSTRACT**

2 Most multinomial choice models, particularly in practice (e.g., multinomial logit model), assume an
3 extreme-value Gumbel distribution for the random components of utility functions. The use of this
4 distribution offers a closed-form likelihood expression when the utility maximization principle is
5 applied to model choice behaviors. The maximum likelihood estimation method can be easily
6 applied to estimate model coefficients. However, maximum likelihood estimators are consistent and
7 efficient only if distributional assumptions on the random error terms are valid. It is therefore
8 important to test the validity of underlying distributional assumptions that form the basis of parameter
9 estimation and policy evaluation. In this paper, a practical but strict method is proposed to test the
10 distributional assumption of the random component of utility functions in both the multinomial logit
11 (MNL) model and multiple discrete-continuous extreme value (MDCEV) model. Based on a semi-
12 nonparametric approach, a closed-form likelihood function that nests the MNL or MDCEV model being
13 tested is derived. Then, the traditional likelihood ratio test can be applied to test violations of the
14 standard Gumbel distribution assumption. Simulation experiments are conducted to show that the
15 test yields acceptable Type-I and Type-II error probabilities at commonly available sample sizes. The
16 test is then applied to three real-world discrete and discrete-continuous choice models. For all three
17 models, the proposed test rejects the validity of the standard Gumbel distribution in most utility
18 functions, calling for the development of approaches that overcome adverse effects of violations of
19 distributional assumptions.

20

21 **Keywords:** travel behavior models, discrete choice models, violations of distributional assumptions,
22 test of validity of distributional assumption, multinomial logit model, multiple discrete-continuous
23 extreme value model

1. INTRODUCTION

The Gumbel distribution (also called Type-I extreme value distribution) plays a central role in travel choice models, including both discrete choice models (McFadden, 1974) and multiple discrete-continuous choice models (e.g., Bhat, 2005 and Bhat, 2008). This can be attributed to two main reasons. First, the Gumbel distribution is close to the normal distribution; in the absence of any specific information about the behavioral phenomenon under investigation, it is often assumed in econometric choice models that the random disturbance term which captures the overall impact of unobserved factors is normally distributed. Second, when the Gumbel distribution is assumed for random components in utility functions, a closed-form expression for the likelihood function is obtained when the utility maximization principle is applied. With a neat closed-form expression for the likelihood function, maximum likelihood estimation (MLE) methods can be easily applied to estimate model coefficients consistently and efficiently.

Given these advantages, the Multinomial Logit (MNL) model and Multiple Discrete-Continuous Extreme Value (MDCEV) model, both of which are based on the Gumbel distribution assumption for the random error components, are widely used in practice. Although strides have been made in estimating model formulations that assume a normal distribution for the random error components, namely, the Multinomial Probit Model (Train, 2009) and Multiple Discrete-Continuous Probit (MDCP) model (Bhat et al., 2013), the logit-based models continue to be the model forms of choice for travel demand forecasting. However, the theory of maximum likelihood estimation indicates that the consistency and efficiency of maximum likelihood estimators depend on the validity of the distributional assumption made on the random error components. If the distributional assumption is violated, then the maximum likelihood estimators are neither consistent nor efficient.

In prediction settings, the MNL model ensures that predicted market shares match the observed shares in the sample (Ben-Akiva and Lerman, 1985). In the case of the MNL model, violations of the distributional assumption will therefore not adversely affect the predicted aggregate market shares. In the case of the MDCEV model, however, such a property does not hold. Jäggi et al. (2013) found that predictions from MDCEV models of vehicle fleet composition and usage are quite sensitive to model specification. As the unobserved but significant factors affecting vehicle fleet composition and usage are absorbed into the random error components, they are bound to influence the nature of the distribution of the random error terms. If the model specification results in a situation where there is violation of the standard Gumbel distributional assumption on the random error terms of the MDCEV model, it is reasonable to expect gross inaccuracies in model predictions depending on the severity of the violation.

It is therefore important to validate the assumed distributions on the random error terms prior to applying the MLE method to estimate model coefficients of either discrete or discrete-continuous travel choice models. The objective of this paper is to propose a practical but strict statistical method to test whether the error terms in random utility functions of MNL or MDCEV models follow the assumed Gumbel distribution.

2. LITERATURE REVIEW

Econometricians have been questioning the validity of the distributional assumption on random error components of utility functions ever since McFadden (1974) first proposed the multinomial logit model formulation (e.g., Manski, 1975). Concerns about error distribution violations motivated the development of semi-parametric and semi-nonparametric choice models. The semi-parametric choice

1 model employs the kernel density method to estimate the distribution of the random errors, and
 2 therefore does not rely on any parametric distributional assumptions (e.g., Klein and Spady, 1993; Lee,
 3 1995). The semi-nonparametric (SNP) choice model is based on a polynomial approximation of a
 4 probability density function (PDF) that takes a flexible form (Gallant and Nychka, 1987). Because the
 5 likelihood function has an explicit analytical expression, the SNP choice modeling method appears to
 6 be applied more widely than the semi-parametric approach in practice (e.g., Chen and Randall, 1997;
 7 Creel and Loomis, 1997; Crooker and Herriges, 2007). In this paper, the SNP approach is used to
 8 derive a statistical test of the validity of the Gumbel distribution in logit models of discrete choice. It
 9 would therefore be prudent to first review the SNP binary choice model.

10 Similar to the binary probit model, the SNP binary choice model is also based on a random utility
 11 (U) function, which can be expressed as $U = V + \varepsilon$, where "V" is the systematic or deterministic
 12 component of the utility function and " ε " is the random component. If a dummy variable "y"
 13 indicates whether an alternative is chosen or not, then $P(y = 1) = P(U > 0) = P(V + \varepsilon > 0) =$
 14 $P(\varepsilon > -V)$. The probability density function takes the following form:

$$15 \quad f(\varepsilon) = \frac{(\sum_{i=0}^K a_i \varepsilon^i)^2 \varphi(\varepsilon)}{\int_{-\infty}^{+\infty} (\sum_{i=0}^K a_i \varepsilon^i)^2 \varphi(\varepsilon) d\varepsilon} \quad (1)$$

16 In Equation (1), $\varphi(\varepsilon)$ represents the PDF of the standard normal distribution. The denominator
 17 ensures that $\int_{-\infty}^{+\infty} f(\varepsilon) d\varepsilon = 1$. Equation (1) can be extended and written in the following form:

$$18 \quad f(\varepsilon) = \frac{(\sum_{i=0}^K \sum_{j=0}^K a_i a_j \varepsilon^{i+j}) \varphi(\varepsilon)}{\int_{-\infty}^{+\infty} (\sum_{i=0}^K \sum_{j=0}^K a_i a_j \varepsilon^{i+j}) \varphi(\varepsilon) d\varepsilon} \quad (2)$$

$$19 \quad \text{Then, } P(y = 1) = P(\varepsilon > -V) = \frac{\int_{-V}^{+\infty} (\sum_{i=0}^K \sum_{j=0}^K a_i a_j \varepsilon^{i+j}) \varphi(\varepsilon) d\varepsilon}{\int_{-\infty}^{+\infty} (\sum_{i=0}^K \sum_{j=0}^K a_i a_j \varepsilon^{i+j}) \varphi(\varepsilon) d\varepsilon} = \frac{\sum_{i=0}^K \sum_{j=0}^K a_i a_j \int_{-V}^{+\infty} \varepsilon^{i+j} \varphi(\varepsilon) d\varepsilon}{\sum_{i=0}^K \sum_{j=0}^K a_i a_j \int_{-\infty}^{+\infty} \varepsilon^{i+j} \varphi(\varepsilon) d\varepsilon} \quad (3)$$

20 For the probability value above, recursion formulas may be applied to compute the integral:

$$21 \quad \int \varepsilon^{i+j} \varphi(\varepsilon) d\varepsilon. \quad (4)$$

22 When $K = 0$, the SNP model will reduce to a binary probit model. Thus, the SNP binary choice
 23 model nests the binary probit model as a special case, and can be used to validate the assumption of
 24 normality for the random error component in the binary probit model based on the likelihood ratio
 25 test. In the case of the logit model, it is possible to replace $\varphi(\varepsilon)$ in Equation (1) with the PDF of the
 26 logistic distribution to be tested. However, in that case, Equation (4) will not have a closed form
 27 expression, leading to considerable computational complexity. Another issue is that the SNP choice
 28 model is usually limited to a binary choice context rather than a multinomial choice context, possibly
 29 due to its computational complexity. It is therefore challenging to extend the original SNP approach to
 30 a multinomial choice modeling situation.

31 Bierens (2008) proposed a new polynomial, called the orthonormal Legendre polynomial, for
 32 estimating distributions semi-nonparametrically on the unit interval. In the transportation literature,
 33 this approach has been used to test normal and log-normal distributions of random coefficients in
 34 mixed logit model (Fosgerau and Bierlaire, 2007). In the method proposed in this paper, the
 35 orthonormal Legendre polynomial will be used to test the validity of the Gumbel distribution in both
 36 MNL and MDCEV models. It will be shown in this paper that the resultant likelihood function based
 37 on the new polynomial will have a closed form expression and nest those of the MNL and MDCEV
 38 models. A standard likelihood ratio test can be invoked to verify the validity of the underlying
 39 distributional assumption.

40

1 **3. METHODOLOGY**

2 **3.1 The Semi-Nonparametric (SNP) Distribution Nesting the Standard Gumbel Distribution**

3 According to Fosgerau and Bierlaire (2007) and Bierens (2008), the orthonormal Legendre polynomial
4 can be recursively defined as:

5 $L_0 = 1, L_1 = \sqrt{3}(2x - 1)$ (5)

6 $L_n = \alpha(2x - 1)L_{n-1} + \beta L_{n-2}, n \geq 2$ (6)

7 In Equation (6), $\alpha = \frac{\sqrt{4n^2-1}}{n}$ and $\beta = -\frac{(n-1)\sqrt{2n+1}}{n\sqrt{2n-3}}$.

8 The advantage of using this polynomial is to ensure that:

9 $\int_0^1 L_m(x)L_n(x)dx = \begin{cases} 0 & \text{if } m \neq n \\ 1 & \text{if } m = n \end{cases}$ (7)

10 Then, the polynomial can be used to construct a semi-nonparametric probability density function that
11 extends and nests the PDF of a standard Gumbel distribution as:

12 $f(x) = \frac{\{1+\sum_{k=1}^K \delta_k L_k[G(x)]\}^2}{1+\sum_{k=1}^K \delta_k^2} g(x)$, (8)

13 where $g(x) = \exp(-e^{-x}) \cdot \exp(-x)$, $G(x) = \exp(-e^{-x})$, and δ_k are parameters. Note that the
14 functional expression "exp(x)" is equivalent to "e^x". Based on Equation (7), it is easy to show that

15 $\int_{-\infty}^{+\infty} f(x) = 1$. To test the standard Gumbel distribution, only consider the situation where $K = 1$ and

16 simplify the formula as:

17 $f(x) = \frac{\{1+\delta_1[\gamma_1+\gamma_2 G(x)]\}^2}{1+\delta_1^2} g(x) = \left[\frac{(1+\delta_1\gamma_1)^2}{1+\delta_1^2} + \frac{2(1+\delta_1\gamma_1)\delta_1\gamma_2}{1+\delta_1^2} G(x) + \frac{(\delta_1\gamma_2)^2}{1+\delta_1^2} G(x)^2 \right] g(x)$
18 $= [\xi_0 + \xi_1 G(x) + \xi_2 G(x)^2]g(x) = \{\sum_{m=0}^2 \xi_m [G(x)]^m\}g(x)$ (9)

19 In the formula above, $\xi_0 = \frac{(1+\delta_1\gamma_1)^2}{1+\delta_1^2}$, $\xi_1 = \frac{2(1+\delta_1\gamma_1)\delta_1\gamma_2}{1+\delta_1^2}$, $\xi_2 = \frac{(\delta_1\gamma_2)^2}{1+\delta_1^2}$, $\gamma_1 = -\sqrt{3}$, and $\gamma_2 = 2\sqrt{3}$.

20 The cumulative distribution function (CDF) of the extended distribution is given as $F(x) =$

21 $\int_{-\infty}^x \{\sum_{m=0}^2 \xi_m [G(\epsilon)]^m\}g(\epsilon)d\epsilon$. Since $\int_{-\infty}^x [G(\epsilon)]^m g(\epsilon) d\epsilon = \int_{-\infty}^x [G(\epsilon)]^m dG(\epsilon)$, and letting $z = G(\epsilon)$,

22 $\int_{-\infty}^x [G(\epsilon)]^m g(\epsilon) d\epsilon = \int_0^{G(x)} z^m dz = \frac{[G(x)]^{m+1}}{m+1}$. Thus, one should have that:

23 $F(x) = \sum_{m=0}^2 \left\{ \frac{\xi_m [G(x)]^{m+1}}{m+1} \right\}$. (10)

24

25 **3.2 Test for Validity of Distributional Assumption in the Multinomial Logit (MNL) Model**

26 In a discrete choice model, it is assumed that $U_j = V_j + \epsilon_j$, where the index of alternatives, $j = 1, 2, \dots, J$.

27 In the interest of brevity, the index "i", corresponding to the individual decision-maker, is suppressed

28 in the equation above. In an MNL model, ϵ_j is independently and identically distributed (i.i.d.) as a

29 standard Gumbel distribution. The proposed method can be used to test whether any random error

30 term, ϵ_j , follows the standard Gumbel distribution or not. Without any loss of generality, if the error

31 term of the first alternative needs to be tested, it is possible to calculate:

32 $q_1(m) = \frac{e^{V_1}}{m \cdot e^{V_1} + \sum_{j=1}^J e^{V_j}}$. (11)

33 For the k^{th} alternative not being tested ($k > 1$), proceed to calculate:

$$1 \quad q_k(m) = \frac{e^{V_k}}{(1+m) \cdot (m \cdot e^{V_1} + \sum_{j=1}^J e^{V_j})}. \quad (12)$$

2 Then, calculate the choice probability for each alternative as:

$$3 \quad P_j = \sum_{m=0}^2 \xi_m q_j(m), \quad (13)$$

4 where $\xi_0 = \frac{(1+\delta_1\gamma_1)^2}{1+\delta_1^2}$, $\xi_1 = \frac{2(1+\delta_1\gamma_1)\delta_1\gamma_2}{1+\delta_1^2}$, $\xi_2 = \frac{(\delta_1\gamma_2)^2}{1+\delta_1^2}$, $\gamma_1 = -\sqrt{3}$, and $\gamma_2 = 2\sqrt{3}$.

5 The log-likelihood function over the entire sample can then be formulated as:

$$6 \quad LL = \sum_{i=1}^N \sum_{j=1}^J y_{ij} \ln(P_{ij}), \quad (14)$$

7 where y_{ij} are dummy variables indicating whether the j^{th} alternative is chosen by the individual "i". If

8 the coefficient δ_1 is fixed at 0, $P_{ij} = \frac{e^{V_{ij}}}{\sum_{j=1}^J e^{V_{ij}}}$ and the model reduces to an MNL model. Thus, the

9 likelihood ratio test can be applied to test the null hypothesis that the random error in the utility
10 function of the first alternative follows the standard Gumbel distribution (a complete mathematical
11 derivation is given in Appendix A).

12

13 **3.3 Test of Validity for the Multiple Discrete-Continuous Extreme Value (MDCEV) Model**

14 As per Bhat (2008), the utility function of an MDCEV model takes the following form:

$$15 \quad U_j = \frac{\gamma_j}{\alpha_j} \psi_j \left[\left(\frac{t_j}{\gamma_j} + 1 \right)^{\alpha_j} - 1 \right], \quad (15)$$

16 where α_j is a satiation parameter that accounts for diminishing marginal utility and γ_j is a
17 translation parameter that accommodates corner solutions (zero consumption of certain alternatives).

18 The baseline utility $\psi_j = e^{V_j + \varepsilon_j}$ and t_j represents the continuous resource being allocated to the
19 alternative "j". The index of alternatives, $j = 1, 2, \dots, K$, where "K" is the total number of alternatives
20 and "M" represents the total number of alternatives that are allocated resources ($M \leq K$). In the
21 interest of brevity, the index denoting the individual "i" is dropped from the equation above.

22 Suppose it is of interest to test the distributional assumption on random error ε_1 of the utility
23 function of the first alternative. If $t_1 > 0$, calculate:

$$24 \quad q(m) = \left(\prod_{j=1}^M c_j \right) \cdot \left(\sum_{j=1}^M \frac{1}{c_j} \right) \cdot (M-1)! \cdot \frac{\prod_{j=1}^M e^{V_j}}{[m \cdot e^{V_1} + \sum_{k=1}^K e^{V_k}]^M}. \quad (16)$$

$$25 \quad \text{If } t_1 = 0, \text{ calculate } q(m) = \left(\prod_{j=1}^M c_j \right) \cdot \left(\sum_{j=1}^M \frac{1}{c_j} \right) \cdot (M-1)! \cdot \frac{\prod_{j=2}^{M+1} e^{V_j}}{(1+m)[m \cdot e^{V_1} + \sum_{k=1}^K e^{V_k}]^M}, \quad (17)$$

26 where $c_j = \frac{1-\alpha_j}{t_j + \gamma_j}$. Then, one can compute the likelihood value for each individual as:

$$27 \quad P = \sum_{m=0}^2 \xi_m q(m), \quad (18)$$

28 where $\xi_0 = \frac{(1+\delta_1\gamma_1)^2}{1+\delta_1^2}$, $\xi_1 = \frac{2(1+\delta_1\gamma_1)\delta_1\gamma_2}{1+\delta_1^2}$, $\xi_2 = \frac{(\delta_1\gamma_2)^2}{1+\delta_1^2}$, $\gamma_1 = -\sqrt{3}$, and $\gamma_2 = 2\sqrt{3}$.

29 The log-likelihood function for the entire sample can be formulated as $LL = \sum_{i=1}^N \ln(P_i)$, which can
30 be maximized to estimate model coefficients as well as the parameter δ_1 . Similar to the case of the
31 MNL model, if the parameter δ_1 is fixed at 0, the model will reduce to the MDCEV model. Therefore,
32 the likelihood ratio test can be applied to test the null hypothesis that the random error term of the
33 utility function of the first alternative follows the standard Gumbel distribution (a complete
34 mathematical derivation is given in Appendix B).

35

4. Simulation Experiments

Before applying the methods for models estimated on travel survey data sets in an empirical context, the suitability of the methods was confirmed through the use of simulation experiments. This section describes results of the simulation experiments that were conducted to demonstrate the applicability of the proposed testing methods. The second objective of the simulation experiments is to determine the sample sizes required for controlling the probability of Type-I and Type-II errors in statistical testing. Travel survey sample sizes are often limited and it would be of value for modelers to be aware of the sample sizes required to apply the test proposed in this paper.

4.1 Simulation Experiments for Test of MNL Model

The experiment is designed with four alternatives; four random utility values are computed as:

$$U_1 = 0.4 - 0.5 \times x_1 + \varepsilon_1,$$

$$U_2 = -0.5 - 0.4 \times x_2 + \varepsilon_2,$$

$$U_3 = -0.6 - 0.3 \times x_3 + \varepsilon_3,$$

$$U_4 = -0.5 \times x_4 + \varepsilon_4.$$

In the equations above, x_1 , x_2 , x_3 , and x_4 follow an independent uniform distribution between 0 and 10. ε_2 , ε_3 , and ε_4 follow an independent standard Gumbel distribution [i.e. $G(0,1)$], while ε_1 follows either $G(0,1)$ or a distribution other than $G(0,1)$. Then, four dummy choice variables, say y_1 , y_2 , y_3 , and y_4 , are generated in accordance with the utility maximization principle:

$$y_j = [U_j \geq \max(U_1, U_2, U_3, U_4)].$$

Both MNL and SNP models are first estimated and then the χ^2 statistic is computed as $2 \times (LL_{SNP} - LL_{MNL})$, where LL_{SNP} and LL_{MNL} are the log-likelihood values of the SNP and MNL models at convergence. At one degree of freedom, the critical χ^2 value is 3.84 at a 0.05 significance level. If the value of the χ^2 statistic is greater than 3.84, the null hypothesis that ε_1 follows the standard Gumbel distribution is rejected at a 95 percent confidence level. Two types of simulation experiments are conducted to investigate the probability of making Type-I and Type-II errors in the hypothesis tests. The Type-I error refers to the case where the null hypothesis is rejected when it is true, while the Type-II error refers to the case where the null hypothesis is not rejected when it is false.

To determine the probability of making a Type-I error, the true standard Gumbel distribution is generated for ε_1 . Then simulation experiments are repeated 100 times and the frequency of rejection of the null hypothesis (even though it is true) is recorded. This frequency of erroneously rejecting the null hypothesis is used to estimate the probability of making a Type-I error. The first two rows of Table 1 show the estimated probabilities of making a Type-I error when the sample size is 200 (small) and 4000 (large) respectively. It is found that the Type-I error probability is less than 0.05 for both sample sizes. Results of the simulation experiment thus demonstrate that the probability of making a Type-I error is very small for the proposed testing method.

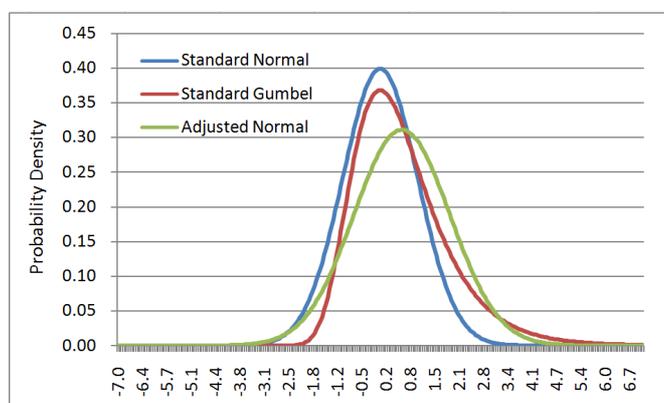
For examining the probability of making Type-II errors, ε_1 should follow a distribution other than the standard Gumbel distribution, $G(0,1)$. Two types of normal distributions are chosen for ε_1 in this simulation experiment. One is the standard normal distribution, i.e., $N(0,1)$, and the other is a normal distribution adjusted to have the same expectation and standard deviation as that of a standard Gumbel distribution. Figure 1 compares the PDFs of the three distributions. Based on a visual examination, the standard normal distribution seems closer to the standard Gumbel distribution than the adjusted normal distribution. Distributions that are very similar to one another are chosen for this experiment to examine the statistical power of the proposed test. The statistical power can

1 be defined as the probability of correctly rejecting the null hypothesis, i.e., $P(\text{Reject } H_0 | H_0 \text{ is wrong})$,
 2 which is equal to $[1 - P(\text{Type-II error})]$. Through simulation experiments, it is found that the
 3 probability of making Type-II errors (or statistical power of the test) depends on both the sample size
 4 and how close the erroneous distribution is to the standard Gumbel distribution with respect to
 5 expectation and standard deviation.

6
7 **Table 1. Simulation Results for MNL models**

Distribution	Sample Size	Accepted	Rejected
$\varepsilon_1 \sim G(0, 1)$	200	1.00	0.00 (Type-I Error)
$\varepsilon_1 \sim G(0, 1)$	4000	0.96	0.04 (Type-I Error)
$\varepsilon_1 \sim N(0, 1)$	4000	0.03 (Type-II Error)	0.97
$\varepsilon_1 \sim \gamma + N(0, 1) \cdot \frac{\pi}{\sqrt{6}}$	200000	0.02 (Type-II Error)	0.98

8 γ is the Euler constant ≈ 0.577216 ; number of repetitions = 100; 0.05 level of significance is used
 9



10
11 **Figure 1. Comparison of Distributions Being Tested**

12
13 A sample size of 4000 is found to be adequate to provide a statistical power of 0.97 and a Type-II
 14 error probability less than 0.05 in distinguishing the standard normal distribution from the standard
 15 Gumbel distribution. However, when the erroneous distribution is the adjusted normal distribution
 16 with the same expectation and standard deviation as that of a standard Gumbel distribution, the
 17 sample size needs to be as high as 200000 to obtain satisfactory statistical power and Type-II error
 18 probability. Although the standard normal distribution appears more similar (visually) to the
 19 standard Gumbel distribution than the adjusted normal distribution, it is actually much more difficult
 20 to distinguish the adjusted normal distribution (than the standard normal distribution) from the
 21 standard Gumbel distribution presumably because the adjusted normal distribution has the same
 22 expectation and standard deviation as the standard Gumbel distribution.
 23

24 **4.2 Simulation Experiments for Test of MDCEV Model**

25 In the MDCEV model specification, either parameter α_j or γ_j needs to be normalized for identification
 26 purposes. Alternative normalization approaches will result in two different model profiles, namely,
 27 the " α " profile and " γ " profile. Simulation experiments are conducted for both profiles.

28 Similar to previous experiment for the MNL model, the experiment is designed with four
 29 alternatives and corresponding random utility values, computed as:

1 $u_1 = -1.0 + 0.9 \times x_1 + \varepsilon_1;$

2 $u_2 = -0.6 + 0.8 \times x_2 + \varepsilon_2;$

3 $u_3 = -0.7 + 0.4 \times x_3 + \varepsilon_3;$

4 $u_4 = 0.6 \times x_4 + \varepsilon_4.$

5 Once again, $x_1, x_2, x_3,$ and x_4 follow an independent uniform distribution between 0 and 10, i.e., $U(0,10)$.
 6 $\varepsilon_2, \varepsilon_3,$ and ε_4 follow an independent standard Gumbel distribution, i.e., $G(0,1)$, while ε_1 follows either
 7 $G(0,1)$ or a distribution other than $G(0,1)$. The continuous resource budget $T = \text{Trunc} [U(0,1000)]+10$,
 8 where $\text{Trunc} []$ is a function to convert a real number to an integer by truncating its decimal part. The

9 utility function for the α -profile MDCEV model is $U_j = \frac{1}{\alpha_j} \psi_j \left[\left(\frac{t_j}{1} + 1 \right)^{\alpha_j} - 1 \right]$, where $\psi_j = e^{u_j}$, $\alpha_1 =$
 10 $0.5, \alpha_2 = 0.6, \alpha_3 = 0.7$ and $\alpha_4 = 0.8$. Then, four resource allocation variables, $t_1, t_2, t_3,$ and $t_4,$ are
 11 calculated by applying the efficient algorithm proposed in Pinjari and Bhat (2011).

12 Table 2 furnishes the simulation results for α -Profile MDCEV models. When the sample size is
 13 500 (relatively small) or 4000 (relatively large), it is found that the proposed test provides satisfactory
 14 Type-I error probabilities. In the case of Type-II errors, it is found that a sample size of 1000 is
 15 adequate to realize a satisfactorily small Type-II error probability when ε_1 follows the standard normal
 16 distribution. A sample with just 3000 observations is adequate to distinguish the adjusted normal
 17 distribution from the standard Gumbel distribution. In the case of the MNL model, the simulation
 18 experiment results reported in the previous section suggest that the sample size required is 4000 and
 19 200000 to yield satisfactory Type-II error probabilities. The difference in statistical power between
 20 the MNL and MDCEV models is likely due to the different nature of the two types of models. Whereas
 21 the MNL model focuses exclusively on a single discrete choice from among a set of alternatives, the
 22 MDCEV model utilizes information about multiple discrete choices and the allocation of a continuous
 23 budget to each of the chosen alternatives. As the MDCEV model utilizes more information about
 24 choices and resource allocation behavior, a smaller sample is presumably needed to avoid making a
 25 Type-II error.

26
27 **Table 2. Simulation Experiment Results for MDCEV models**

Distribution	Sample Size	Accepted	Rejected
α -Profile MDCEV Model			
$\varepsilon_1 \sim G(0, 1)$	500	0.97	0.03 (Type-I Error)
$\varepsilon_1 \sim G(0, 1)$	4000	0.95	0.05 (Type-I Error)
$\varepsilon_1 \sim N(0, 1)$	1000	0.05 (Type-II Error)	0.95
$\varepsilon_1 \sim \gamma + N(0, 1) \cdot \frac{\pi}{\sqrt{6}}$	3000	0.05 (Type-II Error)	0.95
γ -Profile MDCEV Model			
$\varepsilon_1 \sim G(0, 1)$	500	0.97	0.03 (Type-I Error)
$\varepsilon_1 \sim G(0, 1)$	4000	0.95	0.05 (Type-I Error)
$\varepsilon_1 \sim N(0, 1)$	750	0.02 (Type-II Error)	0.98
$\varepsilon_1 \sim \gamma + N(0, 1) \cdot \frac{\pi}{\sqrt{6}}$	750	0.03 (Type-II Error)	0.97

28 Number of Repetitions = 100; 0.05 significance level is used

29

1 The experiment for γ -profile MDCEV models also considers four alternatives with random utility
2 values computed as:

3 $u_1 = 0.4 - 0.5 \times x_1 + \varepsilon_1;$

4 $u_2 = -0.5 - 0.4 \times x_2 + \varepsilon_2;$

5 $u_3 = -0.6 - 0.3 \times x_3 + \varepsilon_3;$

6 $u_4 = -0.5 \times x_4 + \varepsilon_4.$

7 $x_1, x_2, x_3,$ and x_4 follow an independent uniform distribution between 0 and 10, i.e., $U(0, 10)$. $\varepsilon_2, \varepsilon_3,$ and
8 ε_4 follow an independent standard Gumbel distribution, i.e., $G(0,1)$, while ε_1 follows either $G(0,1)$ or a
9 distribution other than $G(0,1)$. The continuous resource budget $T = \text{Trunc}[U(0,500)] + 10$. The utility

10 function is $U_j = \gamma_j \psi_j \ln \left[\frac{t_j}{\gamma_j} + 1 \right]$, where $\psi_j = e^{u_j}$, $\gamma_1 = 2.0$, $\gamma_2 = 1.0$, $\gamma_3 = 0.5$ and $\gamma_4 = 1.5$. Four resource
11 allocation variables $t_1, t_2, t_3,$ and t_4 are calculated by applying the efficient forecasting algorithm of
12 Pinjari and Bhat (2011). Table 2 provides the simulation experiment results for γ -Profile MDCEV
13 models.

14 Similar to the α -profile MDCEV models, γ -profile MDCEV models have satisfactory Type-I error
15 probabilities when the sample size is 500 (relatively small) and 4000 (relatively large), respectively.
16 As for the statistical power, γ -profile MDCEV models require even fewer observations (750) than the
17 α -profile MDCEV models to distinguish the normal or adjusted normal distribution from the standard
18 Gumbel distribution.

19

20 5. CASE STUDIES

21 The simulation experiments demonstrated the efficacy of the statistical test developed in this research
22 effort. The proposed testing methods were applied to MNL and MDCEV models estimated on real-
23 world travel survey data sets to examine the extent to which violations of the standard Gumbel
24 distribution occur in different empirical contexts.

25

26 5.1 MNL Model of Long-distance Travel Mode Choice

27 The data set used for the MNL model case study is obtained from the "AER" package available in R
28 statistical platform (Greene, 2011). The data set includes 210 observations of mode choice for long-
29 distance travel among four alternative modes: Air, Train, Bus, and Car. The MNL model is estimated
30 using a specification that includes a number of explanatory variables. All of the explanatory variables
31 exhibit behaviorally intuitive and statistically significant coefficients in the MNL model. The model
32 estimation results furnish goodness-of-fit statistics that are consistent with those typically seen for
33 MNL models in research studies and practice. The proposed method is applied to test the validity of
34 the standard Gumbel distribution in each utility function. Estimation results are presented in Table
35 3. It is seen that the test does not reject the null hypothesis that the random error follows the
36 standard Gumbel distribution for the first, third, and fourth alternatives (air, bus and car). However,
37 the test does reject the standard Gumbel distributional assumption in the case of the second utility
38 function (train mode). The χ^2 statistic is 8.933 and the corresponding p-value is 0.003. As per the
39 simulation result in Table 1, the Type-I error probability is almost zero when the sample size is as small
40 as 200. A Type-I error occurs when the standard Gumbel distribution is rejected even though the
41 random error truly follows the standard Gumbel distribution.

Table 3. Test Results of the MNL model (N = 210)

Model	MNL Model		Test (Air Utility)		Test (Train Utility)		Test (Bus Utility)		Test (Car Utility)	
	Coeff.	t-test	Coeff.	t-test	Coeff.	t-test	Coeff.	t-test	Coeff.	t-test
Constant for Air	8.038	5.58	7.822	4.86	7.618	5.39	7.955	5.54	6.425	4.35
Air travel time (min)	-0.030	-4.18	-0.031	-4.20	-0.031	-4.50	-0.030	-4.17	-0.024	-3.40
Party size in mode choice	-0.951	-3.66	-0.962	-3.62	-0.864	-3.45	-0.937	-3.62	-0.837	-3.49
Air waiting time (min)	-0.103	-5.72	-0.105	-5.75	-0.101	-5.69	-0.103	-5.71	-0.094	-5.41
δ_i	--	--	0.133	0.40	--	--	--	--	--	--
Constant for Train	4.409	5.03	4.429	5.04	4.253	6.45	4.349	4.98	3.185	3.61
Train travel time (min)	-0.005	-2.94	-0.005	-2.99	-0.005	-3.80	-0.005	-2.97	-0.004	-2.55
Train cost (\$)	-0.024	-1.84	-0.024	-1.85	-0.019	-2.04	-0.024	-1.83	-0.021	-1.75
Household annual income (K \$)	-0.048	-3.66	-0.048	-3.65	-0.036	-3.96	-0.047	-3.62	-0.042	-3.45
Train waiting time (min)	-0.064	-3.83	-0.064	-3.82	-0.045	-4.17	-0.064	-3.80	-0.058	-3.64
δ_i	--	--	--	--	-0.745	-3.43	--	--	--	--
Constant for Bus	4.905	3.85	4.919	3.86	4.461	3.68	5.033	4.19	3.643	2.88
Bus travel time (min)	-0.006	-3.26	-0.006	-3.30	-0.006	-3.55	-0.006	-3.40	-0.004	-2.29
Bus waiting time (min)	-0.151	-5.17	-0.151	-5.16	-0.141	-5.09	-0.138	-4.13	-0.147	-5.18
δ_i	--	--	--	--	--	--	-0.195	-1.07	--	--
Car travel time (min)	-0.006	-5.13	-0.007	-5.14	-0.007	-5.70	-0.006	-5.15	-0.005	-3.68
δ_i	--	--	--	--	--	--	--	--	-0.588	-1.16
$LL(\beta)$	-160.092	--	-159.963	--	-155.626	--	-159.751	--	-159.339	--
χ^2 Statistic	--	--	0.258	--	8.933	--	0.683	--	1.506	--
p-value	--	--	0.611	--	0.003	--	0.409	--	0.220	--

1 Therefore, it can be concluded that the standard Gumbel distribution is not valid for the random error
 2 term in the utility equation of the train mode. It is also interesting to see that the additional coefficient
 3 " δ_1 " is significant and considerably changes the magnitude of coefficients in the second utility function.

4
 5 **5.2 MDCEV Model of Activity Engagement and Time Allocation**

6 The proposed method is next applied to an MDCEV model of activity (stop) engagement and time
 7 allocation for home-based work tours, where the primary purpose of the tour is to go to the workplace.
 8 The model predicts the secondary activities that will be undertaken during the tour, along with the
 9 time allocated to each activity. The model is estimated on 2009 National Household Travel Survey
 10 (NHTS) data from the Greater Phoenix Metropolitan Area (Garikapati et al. 2014). In the interest of
 11 brevity, model estimation results are suppressed in this paper. Activities within a home-based work
 12 tour are classified into 11 types. Explanatory variables include commuters' demographic and socio-
 13 economic characteristics, commute and work status, built environment attributes associated with
 14 residential and work locations, and accessibility measures. The sample size of the estimation sample
 15 is 1968 and the log-likelihood value at convergence is -17528.658.

16 The proposed method is applied to test the validity of the standard Gumbel distribution for the
 17 random error term in the utility function of each activity type. The test results, including the log-
 18 likelihood value at convergence, χ^2 statistics and corresponding p-values, are listed in Table 4. It is
 19 seen that the likelihood ratio tests reject the distributional assumption for *all* of the utility functions in
 20 the MDCEV model. It is interesting to see that the χ^2 statistic values for two outside goods (going to
 21 work and returning home) are considerably lower than those for other goods, which implies that the
 22 distributional assumption is even more invalid for alternatives that are not outside goods. Outside
 23 goods are those that are consumed (chosen) by all observations in the data set. In this case, going to
 24 work and returning home are activities that must be undertaken in the context of home-based work
 25 tours and are therefore treated as outside goods.

26
 27 **Table 4. Test Results for the MDCEV Model of Home-Based Work Tour**
 28 **Activity Engagement and Time Allocation**

Activity Type	LL(b)	χ^2 Statistics	p-value
Work Epoch (outside good 1)	-17218.635	620.047	0.000
Home Epoch (outside good 2)	-17377.888	301.539	0.000
Other Escort Epoch 1	-14957.450	5142.417	0.000
Other Escort Epoch 2	-14871.097	5315.122	0.000
School Escort Epoch 1	-14987.671	5081.973	0.000
School Escort Epoch 2	-14860.288	5336.741	0.000
Shopping Epoch 1	-15118.892	4819.533	0.000
Maintenance Epoch 1	-15071.536	4914.243	0.000
Meal Epoch 1	-14987.993	5081.330	0.000
Social Visit Epoch 1	-14850.375	5356.567	0.000
Other Discretionary Epoch 1	-14925.381	5206.554	0.000

29
 30 **5.3 MDCEV Model of Household Vehicle Fleet Composition and Utilization**

31 The proposed method is finally applied to a MDCEV model of household vehicle fleet composition and
 32 utilization. The model was estimated on 2009 National Household Travel Survey (NHTS) data from

1 the Greater Phoenix Metropolitan Area (You et al. 2014). Model estimation results are suppressed in
 2 the interest of brevity. Household vehicles are categorized into 14 types, defined by a cross-
 3 classification of body type and vintage. Explanatory variables include household composition and
 4 structure, economic status, and built environment attributes of the household location. The estimation
 5 sample has 4262 observations and the log-likelihood value at convergence is -77020.486.

6 The proposed method is applied to test the validity of the standard Gumbel distribution for the
 7 random error term in the utility function of each vehicle type. Results of the test are shown in Table
 8 5. It is seen once again that the likelihood ratio tests reject the distributional assumption for all utility
 9 functions of the MDCEV model. As in the previous case, the χ^2 statistic for the outside good is much
 10 lower than those for other goods, which implies that the distributional assumption is more invalid for
 11 alternatives that are not treated as outside goods. In this model, non-motorized vehicle is treated as
 12 the outside good as all individuals are assumed to walk for at least some minimal duration over the
 13 course of a day. Thus, the non-motorized vehicle is chosen by all observations in the data set.

14
 15 **Table 5. Test Results for the MDCEV Model of Vehicle Fleet Composition and Utilization**

Vehicle Type	LL(b)	Chi-Sq. Statistics	p-value
Non-motorized vehicle (outside good)	-76839.481	362.018	0.000
Car (years) 0–5	-73185.231	7670.517	0.000
6–11	-72898.926	8243.128	0.000
≥12	-72234.349	9572.282	0.000
Van (years) 0–5	-71472.336	11096.308	0.000
6–11	-71493.089	11054.802	0.000
≥12	-71239.748	11561.484	0.000
SUV (years) 0–5	-72159.229	9722.522	0.000
6–11	-71842.624	10355.732	0.000
≥12	-71363.351	11314.278	0.000
Pickup (years) 0–5	-71736.722	10567.536	0.000
6–11	-71832.517	10375.946	0.000
≥12	-71591.310	10858.360	0.000
Motorbike	-71387.042	11266.896	0.000

16
 17 **6. CONCLUSIONS AND DISCUSSIONS**

18 In this paper, a practical but statistically rigorous method is proposed to test the validity of the standard
 19 Gumbel distribution assumption that is often associated with the random error components in
 20 multinomial travel-related choice models, including both discrete and discrete-continuous models.
 21 The method is based on the use of the orthonormal Legendre polynomial to derive a closed-form
 22 likelihood expression that nests the likelihood functions of the multinomial logit (MNL) and multiple
 23 discrete-continuous extreme value (MDCEV) models. The standard likelihood-ratio test can then be
 24 applied to test the validity of the Gumbel distribution innate to the logit-based choice models. The
 25 efficacy of the proposed method is first examined via simulation experiments. Results of the
 26 simulation experiments show that acceptably low Type-I and Type-II error probabilities in the
 27 application of the test may be realized at typically available travel survey sample sizes, except for the
 28 case of the Type-II error probability for the multinomial logit model (which needs a sample size on the

1 order of 200000 to ensure Type-II error probability less than 0.05). The proposed method is then
2 applied to three real-world case studies, including a multinomial logit model of long-distance travel
3 mode choice, a multiple discrete-continuous choice model of activity-time allocation in home-based
4 work tours, and another multiple discrete-continuous choice model of vehicle fleet composition and
5 utilization. For all three models, the proposed test shows that the assumption of a standard Gumbel
6 distribution for the random error components is rejected at a high level of confidence.

7 In theory, the violation of underlying distributional assumptions will lead to the estimation of
8 parameters that are statistically inconsistent and inefficient. However, the extent to which departures
9 from the standard Gumbel distribution affect the magnitudes of model coefficients and model
10 sensitivity to policy scenarios remains unclear. This question could be addressed in future research
11 using simulation experiments in which a robust model based on valid distributional assumptions is
12 developed. Then, the impacts of modifying the distribution on the random error component on
13 parameter estimates in logit-based models can be explicitly quantified.

14 The three real-world models considered in this paper are typical discrete or discrete-continuous
15 choice models with a rich set of explanatory variables and exhibiting goodness-of-fit statistically
16 typically encountered in practice. Given that violations of the standard Gumbel distribution are
17 occurring even in the context of these models, it is prudent to identify approaches that could
18 potentially overcome any ill-effects of the distributional assumption violations. Three strategies are
19 noted below:

20 1. *Adopting an alternative parametric distribution*

21 The normal distribution is usually a preferred distribution for the error term that captures the
22 effects of unobserved factors (as opposed to the Gumbel distribution). Thus, the corresponding
23 multinomial probit (MNP) or multiple discrete-continuous probit (MDCP) models (Bhat et al., 2013)
24 are likely to be better alternatives to MNL and MDCEV models. In addition, heteroskedastic
25 versions of the logit model or discrete-continuous choice model (Bhat, 1995; Sikder and Pinjari,
26 2013) may also prove to be better alternatives. However, tests should be conducted to determine
27 whether violations of the chosen alternative parametric distribution are occurring. There are
28 methods currently available to test for violations of the normal distribution in the econometric
29 literature (e.g., Bera et al., 1984).

30 2. *Adopting mixed logit model with random coefficients of known distribution*

31 Modelers may still use the standard Gumbel distribution for the random error components, but
32 could identify a random coefficient that follows a certain parametric distributional assumption.
33 Then, this coefficient may be scaled up to approximate the random utility values and minimize the
34 impact of the Gumbel random disturbance. Details of this method are described in Train (2009).
35 Note that the existence of such a random coefficient is a necessary condition to apply this method.
36 The test proposed by Fosgerau and Bierlaire (2007) may be used to test the distributional
37 assumption on a random coefficient in a mixed logit model.

38 3. *Developing a robust multinomial choice model free from distributional assumptions*

39 Statistical or econometric models estimated using maximum likelihood methods necessarily
40 involve the making of distributional assumptions. Modelers have longed for the development of
41 robust choice models free from distributional assumptions for several decades (e.g., Manski, 1975;
42 Gallant and Nychka, 1987; Klein and Spady, 1993); however, most practical distribution-free or
43 semi-parametric choice methods have been limited to the analysis and modeling of binary choice
44 variables, rendering their application to multinomial choice contexts computationally challenging.

1 The extension of such distribution-free or semi-parametric approaches to the modeling of
2 multinomial choice variables would constitute a worthy research endeavor.

3 4 **ACKNOWLEDGEMENT**

5 This research is funded by the startup grant of the “Thousand Young Talent” program from the Central
6 Organization Department of China.

7 8 9 **REFERENCES**

- 10
11 Ben-Akiva, M. and S.R. Lerman (1985). Discrete choice analysis: theory and application to travel
12 demand, MIT Press, Cambridge, M.A.
- 13 Bera, A., Jarque, C., and Lee, L. (1984). Testing the Normality Assumption in Limited Dependent
14 Variable Models. *International Economic Review*, 25, pp. 563-578.
- 15 Bhat, C.R. (1995). A heteroscedastic extreme value model of intercity travel mode choice.
16 *Transportation Research Part B: Methodological*, Volume 29, Issue 6, December 1995, pp. 471-483
- 17 Bhat, C.R. (2005). A Multiple Discrete–Continuous Extreme Value Model: Formulation and Application
18 to Discretionary Time-use Decisions. *Transportation Research Part B*, Vol. 39, No. 8, 2005, pp. 679–
19 707.
- 20 Bhat, C.R. (2008). The Multiple Discrete-Continuous Extreme Value (MDCEV) Model: Role of Utility
21 Function Parameters, Identification Considerations, and Model Extensions. *Transportation*
22 *Research Part B*, 2008, Vol. 42, No. 3, pp. 274–303.
- 23 Bhat, C.R., M. Castro and M. Khan (2013). A new estimation approach for the multiple discrete–
24 continuous probit (MDCP) choice model. *Transportation Research Part B*, 55, pp. 1-22.
- 25 Bierens, H.J. (2008). Semi-Nonparametric Interval-Censored Mixed Proportional Hazard Models:
26 Identification and Consistency Results. *Econometric Theory*, 24(3), pp. 749-794.
- 27 Chen, H.Z. and A. Randall (1997). Semi-nonparametric estimation of binary response models with an
28 application to natural resource valuation. *Journal of Econometrics*, 76(1-2), pp. 323-340.
- 29 Creel, M. and J. Loomis (1997). Semi-nonparametric Distribution-Free Dichotomous Choice Contingent
30 Valuation. *Journal of Environmental Economics and Management*, 32(3), pp. 341-358.
- 31 Crooker, J.R. and J.A. Herriges (2007). Parametric and Semi-Nonparametric Estimation of Willingness-
32 to-Pay in the Dichotomous Choice Contingent Valuation Framework, *Environmental and Resource*
33 *Economics*, 27(4), pp. 451-480.
- 34 Fosgerau, M. and M. Bierlaire (2007). A practical test for the choice of mixing distribution in discrete
35 choice models. *Transportation Research Part B: Methodological*, 41, pp. 784–794.
- 36 Gallant, A.R. and D. Nychka (1987). Semi-Non-Parametric Maximum Likelihood Estimation,
37 *Econometrica*, 55, pp. 363–390.
- 38 Garikapati, V.M., D. You, R.M. Pendyala, P.S. Vovsha, V. Livshits and K. Jeon (2014). Multiple Discrete-
39 Continuous Model of Activity Participation and Time Allocation for Home-Based Work Tours.
40 *Transportation Research Record: Journal of the Transportation Research Board*, No. 2429,
41 Transportation Research Board of the National Academies, Washington, D.C., 2014, pp. 90–98.
- 42 Greene, W.H. (2011). *Econometric Analysis* (7th Edition), Prentice Hall. The web linkage of the data
43 Table F18-2: <http://people.stern.nyu.edu/wgreene/Text/Edition7/TableF18-2.csv> or the Table F21-
44 2 from <http://pages.stern.nyu.edu/~wgreene/Text/tables/tablelist5.htm>

1 Jäggi, B., C. Weis and K.W. Axhausen (2013). Stated response and multiple discrete-continuous choice
2 models: Analyses of residuals. *Journal of Choice Modelling*, 6, pp. 44-59.

3 Klein, R.W. and R.H. Spady (1993). An Efficient Semiparametric Estimator for Binary Response Models.
4 *Econometrica*, 61(2), pp. 387-421.

5 Lee, L.F. (1995). Semiparametric maximum likelihood estimation of polychotomous and sequential
6 choice models. *Journal of Econometrics*, 65, pp. 381-428

7 Manski, C.F. (1975). Maximum score estimation of the stochastic utility model of choice. *Journal of*
8 *Econometrics*, 3(3), pp. 205–228.

9 Mcfadden, D. (1974). Conditional logit analysis of qualitative choice behavior. in P. Zarembka (ed.),
10 *Frontiers in Econometrics*, pp. 105-142, Academic Press: New York.

11 Pinjari, A.R. and C.R. Bhat (2011). Computationally efficient forecasting procedures for Kuhn-Tucker
12 consumer demand model systems: application to residential energy consumption analysis.
13 Technical paper, Department of Civil and Environmental Engineering, University of South Florida.

14 Sikder, S. and A.R. Pinjari (2013). The benefits of allowing heteroscedastic stochastic distributions in
15 multiple discrete-continuous choice models. *Journal of Choice Modelling*, 9, pp. 39-56.

16 Train, K.E. (2009). *Discrete choice methods with simulation*. Cambridge University Press.

17 You, D., V.M. Garikapati, R.M. Pendyala, C.R. Bhat, S. Dubey, K. Jeon, and V. Livshits (2014).
18 Development of Vehicle Fleet Composition Model System for Implementation in Activity-Based
19 Travel Model. *Transportation Research Record: Journal of the Transportation Research Board*, No.
20 2430, pp. 145–154.

21
22
23
24
25
26
27
28
29

1 **Appendix A: Mathematical Derivation of the Test for MNL Model**

2 Suppose there are "J" alternatives in the choice set and their random utility functions are U_1, U_2, \dots, U_J . And
 3 the utility U_j is expressed as the sum of the systematic component V_j and the random component ε_j , i.e. U_j
 4 $= V_j + \varepsilon_j$. Suppose the standard Gumbel distributional assumption of ε_1 needs to be tested. All the other
 5 ε_j ($j > 1$) really follow the standard Gumbel distribution. One can derive the probability that the alternative
 6 1 is chosen as:

$$\begin{aligned} 7 \quad P(y = 1) &= P(U_1 > U_2, U_1 > U_3, \dots, U_1 > U_J) \\ 8 \quad &= P(V_1 + \varepsilon_1 > V_2 + \varepsilon_2, V_1 + \varepsilon_1 > V_3 + \varepsilon_3, \dots, V_1 + \varepsilon_1 > V_J + \varepsilon_J) \\ 9 \quad &= P(\varepsilon_2 < V_{12} + \varepsilon_1, \varepsilon_3 < V_{13} + \varepsilon_1, \dots, \varepsilon_J < V_{1J} + \varepsilon_1). \end{aligned}$$

10 In the formulae above, V_{ij} represents $V_i - V_j$. Since ε_j are assumed to be independently distributed in an

11 MNL model, then $P(y = 1) = \int_{-\infty}^{+\infty} \left[\prod_{j=2}^J G(V_{1j} + \varepsilon_1) \right] f(\varepsilon_1) d\varepsilon_1$.

12 By plugging the PDF of the semi-nonparametric distribution in Equation (9), one can obtain that

13 $P(y = 1) = \int_{-\infty}^{+\infty} \left[\prod_{j=2}^J G(V_{1j} + \varepsilon_1) \right] \left\{ \sum_{m=0}^2 \xi_m [G(\varepsilon_1)]^m \right\} g(\varepsilon_1) d\varepsilon_1$.

14 Then, it can be simplified as $P(y = 1) = \sum_{m=0}^2 \xi_m q_1(m)$, where

$$\begin{aligned} 15 \quad q_1(m) &= \int_{-\infty}^{+\infty} \left[\prod_{j=2}^J G(V_{1j} + \varepsilon_1) \right] [G(\varepsilon_1)]^m g(\varepsilon_1) d\varepsilon_1 \\ 16 \quad &= \int_{-\infty}^{+\infty} \left[\prod_{j=2}^J \exp(-e^{-V_{1j} - \varepsilon_1}) \right] [\exp(-e^{-\varepsilon_1})]^m \exp(-e^{-\varepsilon_1}) \exp(-\varepsilon_1) d\varepsilon_1 \\ 17 \quad &= - \int_{-\infty}^{+\infty} \left[\prod_{j=2}^J \exp(-e^{-V_{1j} - \varepsilon_1}) \right] [\exp(-e^{-\varepsilon_1})]^m \exp(-e^{-\varepsilon_1}) d[\exp(-\varepsilon_1)]. \end{aligned}$$

18 Let $w = \exp(-\varepsilon_1)$. Since ε_1 ranges from $-\infty$ to $+\infty$, "w" should range from $+\infty$ to 0. Then

$$\begin{aligned} 19 \quad q_1(m) &= - \int_{+\infty}^0 \left[\prod_{j=2}^J \exp(-w \cdot e^{-V_{1j}}) \right] [\exp(-w)]^m \exp(-w) dw \\ 20 \quad &= - \int_{+\infty}^0 \exp(-w \sum_{j=2}^J e^{-V_{1j}} - mw - w) dw. \end{aligned}$$

21 Let $\theta = -\sum_{j=2}^J e^{-V_{1j}} - m - 1$, then $q_1(m) = - \int_{+\infty}^0 \exp(\theta w) dw = -\frac{1}{\theta} \int_{+\infty}^0 \exp(\theta w) d(\theta w)$.

22 Let $\eta = \theta w$. Since "w" ranges from $+\infty$ to 0 and $\theta < 0$, η should range from $-\infty$ to 0. Then, $q_1(m) =$

23 $-\frac{1}{\theta} \int_{-\infty}^0 \exp(\eta) d\eta = -\frac{1}{\theta} = -\frac{1}{-\sum_{j=2}^J e^{-V_{1j}} - m - 1} = \frac{1}{\sum_{j=2}^J e^{V_{1j} - V_1 + m + 1}} = \frac{e^{V_1}}{(m+1)e^{V_1 + \sum_{j=2}^J e^{V_j}} e^{V_j}} = \frac{e^{V_1}}{m \cdot e^{V_1 + \sum_{j=1}^J e^{V_j}}}$.

24 For the alternative other than the 1st one (say the 2nd alternative), its choice probability:

$$\begin{aligned} 25 \quad P(y = 2) &= P(U_2 > U_1, U_2 > U_3, \dots, U_2 > U_J) \\ 26 \quad &= P(V_2 + \varepsilon_2 > V_1 + \varepsilon_1, V_2 + \varepsilon_2 > V_3 + \varepsilon_3, \dots, V_2 + \varepsilon_2 > V_J + \varepsilon_J) \\ 27 \quad &= P(\varepsilon_1 < V_{21} + \varepsilon_2, \varepsilon_3 < V_{23} + \varepsilon_2, \dots, \varepsilon_J < V_{2J} + \varepsilon_2). \end{aligned}$$

28 Since $F(\varepsilon_1) = \sum_{m=0}^2 \left\{ \xi_m \frac{[G(\varepsilon_1)]^{m+1}}{m+1} \right\}$ according to Equation (10), $P(y = 2) =$

29 $\int_{-\infty}^{+\infty} \sum_{m=0}^2 \left\{ \xi_m \frac{[G(V_{21} + \varepsilon_2)]^{m+1}}{m+1} \right\} \left[\prod_{j=3}^J G(V_{2j} + \varepsilon_2) \right] g(\varepsilon_2) d\varepsilon_2$. It can be simplified as

30 $P(y = 2) = \sum_{m=0}^2 \xi_m q_2(m)$, where

$$1 \quad q_2(m) = \frac{1}{m+1} \int_{-\infty}^{+\infty} [G(V_{21} + \varepsilon_2)]^{m+1} \left[\prod_{j=3}^J G(V_{2j} + \varepsilon_2) \right] g(\varepsilon_2) d\varepsilon_2$$

$$2 \quad = \frac{1}{m+1} \int_{-\infty}^{+\infty} \exp[-e^{-V_{21}-\varepsilon_2+\ln(m+1)}] \left[\prod_{j=3}^J \exp(-e^{-V_{2j}-\varepsilon_2}) \right] \exp(-e^{-\varepsilon_2}) \cdot e^{-\varepsilon_2} d\varepsilon_2$$

$$3 \quad = -\frac{1}{m+1} \int_{-\infty}^{+\infty} \exp[-e^{-V_{21}-\varepsilon_2+\ln(m+1)}] \left[\prod_{j=3}^J \exp(-e^{-V_{2j}-\varepsilon_2}) \right] \exp(-e^{-\varepsilon_2}) d e^{-\varepsilon_2}$$

4 Let $w = e^{-\varepsilon_2}$. Since ε_2 ranges from $-\infty$ to $+\infty$, "w" should range from $+\infty$ to 0, then

$$5 \quad q_2(m) = -\frac{1}{m+1} \int_{+\infty}^0 \exp[-w \cdot e^{-V_{21}+\ln(m+1)}] \left[\prod_{j=3}^J \exp(-w \cdot e^{-V_{2j}}) \right] \exp(-w) dw$$

$$6 \quad = -\frac{1}{m+1} \int_{+\infty}^0 \exp[-w \cdot e^{-V_{21}+\ln(m+1)}] \left[\exp(-w \cdot \sum_{j=3}^J e^{-V_{2j}}) \right] \exp(-w) dw.$$

7 Let $\theta = -e^{-V_{21}+\ln(m+1)} - \sum_{j=3}^J e^{-V_{2j}} - 1$, then

$$8 \quad q_2(m) = -\frac{1}{m+1} \int_{+\infty}^0 \exp(\theta w) dw = -\frac{1}{\theta(m+1)} \int_{+\infty}^0 \exp(\theta w) d(\theta w).$$

9 Let $\eta = \theta w$. Since "w" ranges from $+\infty$ to 0 and $\theta < 0$, η should range from $-\infty$ to 0. Then,

$$10 \quad q_2(m) = -\frac{1}{\theta(m+1)} \int_{-\infty}^0 \exp(\eta) d(\eta) = -\frac{1}{\theta(m+1)} = \frac{1}{m+1} \cdot \frac{1}{e^{-V_{21}+\ln(m+1)} + \sum_{j=3}^J e^{-V_{2j}+1}}$$

$$11 \quad = \frac{1}{m+1} \cdot \frac{e^{V_2}}{(m+1) \cdot e^{V_1} + \sum_{j=3}^J e^{V_j} + e^{V_2}} = \frac{e^{V_2}}{(m+1)(m \cdot e^{V_1} + \sum_{j=1}^J e^{V_j})}.$$

12 Without loss of generality, $q_k(m) = \frac{e^{V_k}}{(m+1)(m \cdot e^{V_1} + \sum_{j=1}^J e^{V_j})}$, where $k > 1$.

13

14

15 **Appendix B: Mathematical Derivation of the Test for MDCEV Model**

16 Assume the utility for the alternative "j" takes the following form as in Bhat (2008): $U_j =$

17 $\frac{\gamma_j}{\alpha_j} \psi_j \left[\left(\frac{t_j}{\gamma_j} + 1 \right)^{\alpha_j} - 1 \right]$, where $\psi_j = e^{V_j + \varepsilon_j}$ and the alternative index $j = 1, 2, \dots, K$. Suppose the random

18 error ε_1 in the 1st alternative needs to be tested. "K" represents the total number of alternatives and

19 "M" represents the total number of alternatives being allocated resource ($M \leq K$). Each individual is

20 supposed to allocate continuous resource by maximizing the overall utility value subject to the budget:

21 Maximize $\sum_{j=1}^K \frac{\gamma_j}{\alpha_j} \psi_j \left[\left(\frac{t_j}{\gamma_j} + 1 \right)^{\alpha_j} - 1 \right]$ subject to $\sum_{j=1}^K t_j = T$.

22 It is equivalent to maximizing $\sum_{j=1}^K \frac{\gamma_j}{\alpha_j} \psi_j \left[\left(\frac{t_j}{\gamma_j} + 1 \right)^{\alpha_j} - 1 \right] - \lambda (\sum_{j=1}^K t_j - T)$. For the observation where

23 the resource allocation $t_1 > 0$, according to KT conditions, one should have

$$24 \quad \psi_1 \left(\frac{t_1}{\gamma_1} + 1 \right)^{\alpha_1-1} = \lambda, \quad \psi_2 \left(\frac{t_2}{\gamma_2} + 1 \right)^{\alpha_2-1} = \lambda, \quad \dots, \quad \psi_M \left(\frac{t_M}{\gamma_M} + 1 \right)^{\alpha_M-1} = \lambda; \quad \psi_{M+1} \left(\frac{t_{M+1}}{\gamma_{M+1}} + 1 \right)^{\alpha_{M+1}-1} <$$

$$25 \quad \lambda, \dots, \psi_K \left(\frac{t_K}{\gamma_K} + 1 \right)^{\alpha_K-1} < \lambda.$$

26 In the above formulae, $t_j > 0$ when $j \leq M$ and $t_j = 0$ when $M < j \leq K$. Then, one will have

$$27 \quad \psi_1 \left(\frac{t_1}{\gamma_1} + 1 \right)^{\alpha_1-1} = \psi_j \left(\frac{t_j}{\gamma_j} + 1 \right)^{\alpha_j-1}, \quad j = 2, 3, \dots, M.$$

1 Take a logarithm function on both sides and one should obtain that:

2 $V_1 + \varepsilon_1 = V_j + \varepsilon_j, j = 2, 3, \dots, M;$

3 $V_1 + \varepsilon_1 > V_j + \varepsilon_j, j = (M+1), (M+2), \dots, K,$ where $V_j = x_j \beta_j + (\alpha_j - 1) \ln \left(\frac{t_j}{y_j} + 1 \right).$

4 As per Bhat (2005), the likelihood value $P = |J| \cdot \int_{-\infty}^{\infty} [\prod_{j=2}^M g(V_{1j} + \varepsilon_1)] [\prod_{j=M+1}^K G(V_{1j} + \varepsilon_1)] f(\varepsilon_1) d\varepsilon_1.$

5 To obtain $f(\varepsilon_1),$ one can plug the SNP density function: $f(x) = \{\sum_{m=0}^2 \xi_m [G(x)]^m\} g(x)$ into the
6 equation above and then have that $P = |J| \int_{-\infty}^{\infty} [\prod_{j=2}^M g(V_{1j} + \varepsilon_1)] [\prod_{j=M+1}^K G(V_{1j} +$
7 $\varepsilon_1)] \{\sum_{m=0}^2 \xi_m [G(\varepsilon_1)]^m\} g(\varepsilon_1) d\varepsilon_1$

8 $= \sum_{m=0}^2 \xi_m |J| \int_{-\infty}^{\infty} [\prod_{j=2}^M g(V_{1j} + \varepsilon_1)] [\prod_{j=M+1}^K G(V_{1j} + \varepsilon_1)] [G(\varepsilon_1)]^m g(\varepsilon_1) d\varepsilon_1 = \sum_{m=0}^2 \xi_m q_1(m),$

9 where $q_1(m) = \int_{-\infty}^{+\infty} |J| (\prod_{j=2}^M g(V_{1j} + \varepsilon_1)) (\prod_{j=M+1}^K G(V_{1j} + \varepsilon_1)) [G(\varepsilon_1)]^m g(\varepsilon_1) d\varepsilon_1$

10 $= |J| \int_{-\infty}^{\infty} [\prod_{j=2}^M g(V_{1j} + \varepsilon_1)] [\prod_{j=M+1}^K G(V_{1j} + \varepsilon_1)] [G(\varepsilon_1)]^m g(\varepsilon_1) d\varepsilon_1$

11 $= |J| \int_{-\infty}^{\infty} \{\prod_{j=2}^M [G(V_{1j} + \varepsilon_1) \cdot e^{-V_{1j} - \varepsilon_1}]\} [\prod_{j=M+1}^K G(V_{1j} + \varepsilon_1)] [G(\varepsilon_1)]^m G(\varepsilon_1) e^{-\varepsilon_1} d\varepsilon_1$

12 $= |J| \int_{-\infty}^{\infty} \{\prod_{j=2}^M [G(V_{1j} + \varepsilon_1) \cdot e^{-V_{1j} - \varepsilon_1}]\} [\prod_{j=M+1}^K G(V_{1j} + \varepsilon_1)] [G(\varepsilon_1)]^m G(V_{11} + \varepsilon_1) e^{-\varepsilon_1} d\varepsilon_1$

13 $= |J| \int_{-\infty}^{\infty} \prod_{j=1}^K G(V_{1j} + \varepsilon_1) [\prod_{j=2}^M e^{-V_{1j} - \varepsilon_1}] [G(\varepsilon_1)]^m e^{-\varepsilon_1} d\varepsilon_1$

14 $= |J| \int_{-\infty}^{\infty} [\prod_{j=1}^K \exp(-e^{-V_{1j} - \varepsilon_1})] (\prod_{j=2}^M e^{-V_{1j}}) (e^{-\varepsilon_1})^{M-1} e^{-\varepsilon_1} \cdot [G(\varepsilon_1)]^m d\varepsilon_1$

15 $= |J| (\prod_{j=2}^M e^{-V_{1j}}) \int_{-\infty}^{+\infty} \exp(-\sum_{j=1}^K e^{-V_{1j} - \varepsilon_1}) (e^{-\varepsilon_1})^{M-1} \exp[-e^{-\varepsilon_1 + \ln(m)}] e^{-\varepsilon_1} d\varepsilon_1$

16 Define the integral part in the formula above as "Int":

17 $Int = \int_{-\infty}^{+\infty} \exp(-\sum_{j=1}^K e^{-V_{1j} - \varepsilon_1}) (e^{-\varepsilon_1})^{M-1} \exp[-e^{-\varepsilon_1 + \ln(m)}] e^{-\varepsilon_1} d\varepsilon_1$

18 $= - \int_{-\infty}^{+\infty} \exp(-\sum_{j=1}^K e^{-V_{1j} - \varepsilon_1}) (e^{-\varepsilon_1})^{M-1} \exp[-e^{-\varepsilon_1 + \ln(m)}] d e^{-\varepsilon_1}$

19 Let $w = e^{-\varepsilon_1}.$ Since ε_1 ranges from $-\infty$ to $+\infty,$ w should range from $+\infty$ to 0. Then,

20 $Int = - \int_{+\infty}^0 \exp(-w \sum_{j=1}^K e^{-V_{1j}}) w^{M-1} \exp(-mw) dw$

21 $= - \int_{+\infty}^0 \exp[-w(\sum_{j=1}^K e^{-V_{1j}} + m)] w^{M-1} dw.$

22 Let $a = \sum_{j=1}^K e^{-V_{1j}} + m$ and $Int = - \int_{+\infty}^0 \exp(-w \cdot a) w^{M-1} dw.$ Let $b = -aw,$ then $w = -\frac{b}{a}$ and

23 $dw = -\frac{db}{a}.$ Since "w" ranges from $+\infty$ to 0 and $a > 0,$ "b" should range from $-\infty$ to 0. Thus, $Int =$

24 $-\int_{-\infty}^0 e^b \left(-\frac{b}{a}\right)^{M-1} \left(-\frac{1}{a}\right) db = \frac{(-1)^{M-1}}{a^M} \int_{-\infty}^0 e^b \cdot b^{M-1} db.$

25 As per Bhat (2005), $\int_{-\infty}^0 e^b \cdot b^{M-1} db = (-1)^{M-1} \cdot (M-1)!.$ Then, $Int = \frac{(-1)^{M-1}}{a^M} (-1)^{M-1} \cdot (M-1)!$

$$1) \quad 1) = \frac{(M-1)!}{a^M} = \frac{(M-1)!}{(m+\sum_{j=1}^K e^{-V_{1j}})^M}.$$

$$2) \quad \text{Then, } q_1(m) = |J|(\prod_{j=2}^M e^{-V_{1j}}) \frac{(M-1)!}{(m+\sum_{j=1}^K e^{-V_{1j}})^M} = |J|(\prod_{j=2}^M e^{-V_1+V_j}) \frac{(M-1)!}{(m+\sum_{j=1}^K e^{-V_1+V_j})^M}$$

$$3) \quad = |J|(M-1)! \frac{(\prod_{j=1}^M e^{V_j})}{(m \cdot e^{V_1+\sum_{j=1}^K e^{V_j}})^M}.$$

4) As per page 704 in Bhat (2005), $|J| = (\prod_{j=1}^M c_j) \left(\sum_{j=1}^M \frac{1}{c_j} \right)$, where $c_j = \frac{1-\alpha_j}{\tau_j+\gamma_j}$. Thus, $q_1(m) =$

$$5) \quad (\prod_{j=1}^M c_j) \left(\sum_{j=1}^M \frac{1}{c_j} \right) \left[\frac{\prod_{j=1}^M e^{V_j}}{(m \cdot e^{V_1+\sum_{j=1}^K e^{V_j}})^M} \right] (M-1)!.$$

6) For the observation where the resource allocation $t_1 = 0$, one may pick out another alternative (say "2")
7) to which the resource is allocated and therefore $t_2 > 0$. According to KT conditions,

$$8) \quad V_2 + \varepsilon_2 = V_j + \varepsilon_j, \quad j = 3, \dots, (M+1);$$

$$9) \quad V_2 + \varepsilon_2 > V_1 + \varepsilon_1;$$

$$10) \quad V_2 + \varepsilon_2 > V_j + \varepsilon_j, \quad j = (M+2), (M+3), \dots, K.$$

11) The likelihood value can be computed as:

$$12) \quad P = |J| \cdot \int_{-\infty}^{\infty} [\prod_{j=3}^{M+1} g(V_{2j} + \varepsilon_2)] F(V_{21} + \varepsilon_2) [\prod_{j=M+2}^K G(V_{2j} + \varepsilon_2)] g(\varepsilon_2) d\varepsilon_2.$$

13) Here, the CDF of the SNP distribution, as in Equation (10), is used for $F(x)$. Then, $P = |J| \cdot$

$$14) \quad \int_{-\infty}^{\infty} [\prod_{j=3}^{M+1} g(V_{2j} + \varepsilon_2)] \sum_{m=0}^2 \left\{ \frac{\xi_m [G(V_{21} + \varepsilon_2)]^{m+1}}{m+1} \right\} [\prod_{j=M+2}^K G(V_{2j} + \varepsilon_2)] g(\varepsilon_2) d\varepsilon_2$$

$$15) \quad = \sum_{m=0}^2 \xi_m \frac{|J|}{m+1} \int_{-\infty}^{\infty} [\prod_{j=3}^{M+1} g(V_{2j} + \varepsilon_2)] [G(V_{21} + \varepsilon_2)]^{m+1} [\prod_{j=M+2}^K G(V_{2j} + \varepsilon_2)] g(\varepsilon_2) d\varepsilon_2.$$

$$16) \quad \text{Let } q_2(m) = \frac{|J|}{m+1} \int_{-\infty}^{\infty} [\prod_{j=3}^{M+1} g(V_{2j} + \varepsilon_2)] [G(V_{21} + \varepsilon_2)]^{m+1} [\prod_{j=M+2}^K G(V_{2j} + \varepsilon_2)] g(\varepsilon_2) d\varepsilon_2 \text{ and}$$

$$17) \quad P = \sum_{m=0}^2 \xi_m q_2(m). \text{ Then, let } q_2(m) = \frac{|J|}{m+1} \text{Int1}, \text{ where}$$

$$18) \quad \text{Int1} = \int_{-\infty}^{\infty} [\prod_{j=3}^{M+1} g(V_{2j} + \varepsilon_2)] [G(V_{21} + \varepsilon_2)]^{m+1} [\prod_{j=M+2}^K G(V_{2j} + \varepsilon_2)] g(\varepsilon_2) d\varepsilon_2.$$

$$19) \quad = \int_{-\infty}^{\infty} [\prod_{j=3}^{M+1} G(V_{2j} + \varepsilon_2)] [\prod_{j=3}^{M+1} \exp(-V_{2j} - \varepsilon_2)] [G(V_{21} + \varepsilon_2)]^{m+1} [\prod_{j=M+2}^K G(V_{2j} + \varepsilon_2)] G(V_{22} +$$

$$20) \quad \varepsilon_2) \exp(-\varepsilon_2) d\varepsilon_2$$

$$21) \quad = \int_{-\infty}^{\infty} [\prod_{j=2}^K G(V_{2j} + \varepsilon_2)] [\prod_{j=3}^{M+1} \exp(-V_{2j} - \varepsilon_2)] [G(V_{21} + \varepsilon_2)]^{m+1} \exp(-\varepsilon_2) d\varepsilon_2$$

$$22) \quad = \int_{-\infty}^{\infty} [\prod_{j=2}^K \exp(-e^{-V_{2j}-\varepsilon_2})] [\prod_{j=3}^{M+1} \exp(-V_{2j})] (e^{-\varepsilon_2})^{M-1} \cdot [G(V_{21} + \varepsilon_2)]^{m+1} \cdot \exp(-\varepsilon_2) d\varepsilon_2$$

$$23) \quad = [\prod_{j=3}^{M+1} \exp(-V_{2j})] \int_{-\infty}^{\infty} [\prod_{j=2}^K \exp(-e^{-V_{2j}-\varepsilon_2})] (e^{-\varepsilon_2})^{M-1} \cdot [G(V_{21} + \varepsilon_2)]^m [G(V_{21} + \varepsilon_2)]^1 \cdot$$

$$24) \quad \exp(-\varepsilon_2) d\varepsilon_2$$

$$25) \quad = [\prod_{j=3}^{M+1} \exp(-V_{2j})] \int_{-\infty}^{\infty} [\prod_{j=2}^K \exp(-e^{-V_{2j}-\varepsilon_2})] (e^{-\varepsilon_2})^{M-1} \cdot [G(V_{21} + \varepsilon_2)]^m \exp(-e^{-V_{21}-\varepsilon_2}) \cdot$$

$$26) \quad \exp(-\varepsilon_2) d\varepsilon_2$$

$$1 = [\prod_{j=3}^{M+1} \exp(-V_{2j})] \int_{-\infty}^{\infty} [\prod_{j=1}^K \exp(-e^{-V_{2j}-\varepsilon_2})] (e^{-\varepsilon_2})^{M-1} \cdot [G(V_{21} + \varepsilon_2)]^m \cdot e^{-\varepsilon_2} d\varepsilon_2$$

$$2 = [\prod_{j=3}^{M+1} \exp(-V_{2j})] Int2, \text{ where}$$

$$3 Int2 = \int_{-\infty}^{\infty} [\prod_{j=1}^K \exp(-e^{-V_{2j}-\varepsilon_2})] (e^{-\varepsilon_2})^{M-1} \cdot [G(V_{21} + \varepsilon_2)]^m \cdot e^{-\varepsilon_2} d\varepsilon_2$$

$$4 = \int_{-\infty}^{\infty} [\prod_{j=1}^K \exp(-e^{-V_{2j}-\varepsilon_2})] (e^{-\varepsilon_2})^{M-1} \cdot \exp[-e^{-\varepsilon_2-V_{21}+\ln(m)}] \cdot e^{-\varepsilon_2} d\varepsilon_2$$

$$5 = \int_{-\infty}^{\infty} \exp[-e^{-\varepsilon_2} \sum_{j=1}^K e^{-V_{2j}}] (e^{-\varepsilon_2})^{M-1} \cdot \exp[-e^{-\varepsilon_2} e^{-V_{21}+\ln(m)}] \cdot e^{-\varepsilon_2} d\varepsilon_2$$

$$6 = - \int_{-\infty}^{\infty} \exp[-e^{-\varepsilon_2} \sum_{j=1}^K e^{-V_{2j}}] (e^{-\varepsilon_2})^{M-1} \cdot \exp(-e^{-\varepsilon_2} e^{-V_{21}+\ln(m)}) d e^{-\varepsilon_2}$$

7 Let $w = e^{-\varepsilon_2}$. Since ε_2 ranges from $-\infty$ to $+\infty$, w should range from $+\infty$ to 0. Thus,

$$8 Int2 = - \int_{+\infty}^0 \exp[-w \sum_{j=1}^K e^{-V_{2j}}] w^{M-1} \cdot \exp[-w \cdot e^{-V_{21}+\ln(m)}] dw$$

$$9 = - \int_{+\infty}^0 \exp\{-w[e^{-V_{21}+\ln(m)} + \sum_{j=1}^K e^{-V_{2j}}]\} w^{M-1} dw.$$

10 Let $a = e^{-V_{21}+\ln(m)} + \sum_{j=1}^K e^{-V_{2j}}$ and then $Int2 = - \int_{+\infty}^0 e^{-aw} w^{M-1} dw$. Let $b = -aw$, then $w =$

$$11 -\frac{b}{a}, dw = -\frac{db}{a}. \text{ Let } Int2 = - \int_{-\infty}^0 e^b \left(-\frac{b}{a}\right)^{M-1} \left(-\frac{1}{a}\right) db = \frac{(-1)^{M-1}}{a^M} \int_{-\infty}^0 e^b \cdot b^{M-1} db. \text{ As per Bhat}$$

$$12 (2005), \int_{-\infty}^0 e^b \cdot b^{M-1} db = (-1)^{M-1} \cdot (M-1)!. \text{ Then, } Int2 = \frac{(-1)^{M-1}}{a^M} (-1)^{M-1} \cdot (M-1)! = \frac{(M-1)!}{a^M}.$$

$$13 \text{ Thus, } q_2(m) = \frac{|||}{m+1} [\prod_{j=3}^{M+1} \exp(-V_{2j})] \frac{(M-1)!}{a^M}$$

$$14 = \frac{|||}{m+1} [\prod_{j=3}^{M+1} \exp(-V_{2j})] \frac{(M-1)!}{(e^{-V_{21}+\ln(m)} + \sum_{j=1}^K e^{-V_{2j}})^M} = \frac{|||(M-1)!}{(m+1)} \frac{[\prod_{j=3}^{M+1} e^{-V_{2j}}]}{(me^{-V_{21}+\ln(m)} + \sum_{j=1}^K e^{-V_{2j}})^M}$$

$$15 = \frac{|||(M-1)!}{(m+1)} \frac{[\prod_{j=3}^{M+1} e^{V_j}] (e^{-V_2})^{M-1}}{(me^{V_1+\sum_{j=1}^K e^{V_j}} (e^{-V_2})^M)} = \frac{|||(M-1)!}{(m+1)} \frac{\prod_{j=2}^{M+1} e^{V_j}}{(me^{V_1+\sum_{j=1}^K e^{V_j}})^M}$$

$$16 = (\prod_{j=1}^M c_j) \left(\sum_{j=1}^M \frac{1}{c_j}\right) \frac{(M-1)!}{(m+1)} \frac{\prod_{j=2}^{M+1} e^{V_j}}{(me^{V_1+\sum_{j=1}^K e^{V_j}})^M}.$$

17 Since $t_1 = 0$, e^{V_1} will not be multiplied with the numerator term. However, the numerator term is still the

18 product of exponentials of utilities for "M" alternatives being allocated resource. In summary, the SNP

19 likelihood function can be expressed as $P = \sum_{m=0}^2 \xi_m q(m)$, where

$$20 q(m) = \begin{cases} (\prod_{j=1}^M c_j) \left(\sum_{j=1}^M \frac{1}{c_j}\right) \left[\frac{\prod_{j=1}^M e^{V_j}}{(m \cdot e^{V_1+\sum_{j=1}^K e^{V_j}})^M} \right] (M-1)!, & \text{if } t_1 > 0, \\ (\prod_{j=1}^M c_j) \left(\sum_{j=1}^M \frac{1}{c_j}\right) \left[\frac{\prod_{j=2}^{M+1} e^{V_j}}{(m+1)(me^{V_1+\sum_{j=1}^K e^{V_j}})^M} \right] (M-1)!, & \text{if } t_1 = 0 \end{cases}$$

21