# ACCOUNTING FOR STOCHASTIC VARIABLES IN DISCRETE CHOICE MODELS

**Federico Díaz,** *Transmetro S.A.S. fediaz@transmetro.gov.co*

**Víctor Cantillo, Julian Arellana.** *Department of Civil and Environmental Engineering, Universidad del Norte vcantill@uninorte.edu.co; jarellana@uninorte.edu.co*

**Juan de Dios Ortúzar**. *Department of Transport Engineering and Logistics, Pontificia Universidad Católica de Chile. jos@ing.puc.cl*

## ABSTRACT

The estimation of discrete choice models requires collecting data about the socioeconomic characteristics of individuals and measuring the attributes describing the alternatives within each individual's choice set. Even though some attributes are intrinsically stochastic (e.g. travel times) or are subject to non-negligible measurement errors (e.g. waiting times), they are usually assumed fixed and deterministic. Indeed, even an accurate measurement can be biased as it might differ from the original (experienced) value perceived by the individual.

Experimental evidence suggests that discrepancies between the values measured by the modeller and experienced by the individuals can lead to incorrect parameter estimates. On the other hand, there is, as usual, an important trade-off between data quality and collection costs. This paper explores the inclusion of stochastic variables in discrete choice models through an econometric analysis that allows identifying the most suitable specifications. Various model specifications were experimentally tested using Monte Carlo simulation. Comparisons included tests for unbiased parameter estimation, computation of marginal rates of substitution and demand forecasts (response analysis) under the implementation of different transport policies.

Results show that error components models can effectively deal with stochastic variables. Also, as in previous misspecification tests reported in the literature, the Multinomial Logit model proved to be quite robust for estimating marginal rates of substitution and forecasting demand for realistic policies, especially when it was estimated with a large number of observations.

*Keywords: stochastic variables, errors in variables, discrete choice models, mixed logit)*

## INTRODUCTION

The estimation of discrete choice models requires data such as socioeconomic characteristics of individuals and attributes of the alternatives within their choice sets. These explanatory variables are usually assumed to be inherently deterministic, that is, that they would yield the same values if measured repeatedly. The problem is that some variables are actually intrinsically stochastic (e.g. travel time under congested conditions[1]) and thus assuming that an accurate measurement made by the modeller is equal to the value originally perceived by the individual can be fairly heroic. In fact even an accurate measure can be biased if it is different from the value perceived by the decision maker.

Furthermore, variables which are intrinsically non-stochastic can still be measured inaccurately producing measurement errors. These errors induce a particular kind of randomness from the modeller's point of view. For instance, in strategic planning applications it is common practice to use zone-based network models to obtain level of service attributes, such as travel time, instead of measuring this key attribute at an individual level due to the high data collection costs involved. Also, trips with different levels of service are usually temporally and spatially aggregated (e.g. the set of trips between two specific zones at a peak hour) and a single level of service value (e.g. an "average" value) is assigned to them, which is evidently different from the true values experienced by the users (Train, 1978). Measurement errors also occur when values are directly provided by the individual in a revealed preference (RP) survey (e.g. waiting time to board a bus, income, or preferred departure time). In this case the difference between the reported value and the real one can be significant due to cognitive issues or even policy bias (Daly and Ortuzar, 1990).

When a discrepancy between the "true" value and the value measured by the modeller exists, an estimation bias arises as shown by Ortuzar and Willumsen (2011, section 9.2). Let us consider a simple Multinomial Logit (MNL) model with a typical utility function $U = \beta x + \varepsilon$ (where $\beta$ are parameters to be estimated, $x$ are measured attributes and $\varepsilon$ is an independent and identically distributed Gumbel error term with mean zero and standard deviation $\sigma_\varepsilon$).

Assume there is a difference between the attribute values as perceived by the modeller $(x^*)$ and the true values $(x)$, such that: $x = x^* + \eta$, where $\eta$ distributes with mean zero and standard deviation $\sigma_\eta$. In this case the utility function is transformed to: $U = \beta(x^*+\eta) + \varepsilon = \beta x^* + (\varepsilon + \beta\eta) = \beta x^* + \delta$. The outcome of this is that in the original model, the estimated parameter $\beta'$ would be:

$$\beta' = \frac{\pi}{\sqrt{6} \cdot \sigma_\varepsilon} \beta \tag{1}$$

---

[1] Related problems arising from the inherent variability of some level of service attributes such as travel time are reliability and risk aversion (Jackson and Jucker, 1982). In this research we will only address the difference between the true value and the values measured by the modeller as a result of this variability.

whilst in the second model the estimated parameter $\beta''$ would be:

$$\beta'' = \frac{\pi}{\sqrt{6} \cdot \sigma_\delta} \beta \qquad (2)$$

Where the standard deviation of the distribution function of the new error component $\delta$ is:

$$\sigma_\delta = \sqrt{\sigma_\varepsilon^2 + \beta^2 \cdot \sigma_\eta^2} \qquad (3)$$

Hence $\beta'' < \beta'$ and this estimation bias may affect the model forecasts.

There is also experimental evidence of bias estimation and miscalculation of marginal rates of substitution when measurement errors occur. For instance, Train (1978) explores the use of more accurate data in the estimation of mode choice models concluding that it is sometimes advisable to carry out an additional effort to correct for the measurement bias of some attributes, such as transit transfer time, when analysing transport policies. Ortuzar and Ivelic (1986) showed that using very precise real data measured at the individual level when estimating mode choice models resulted in better fit and clearly different subjective values of time in comparison with models estimated with aggregate data. More recently, Bhatta and Larsen (2011) show, using synthetic data, how measurement biases may induce biased parameter estimates on a MNL model, besides miscalculation of marginal rates of substitution. Therefore the use of more accurate (but more expensive) data results in better parameter estimates and this clearly establishes a trade-off between data quality and data collection costs (Daly and Ortuzar, 1990).

In this paper we will deal with the problem of working with incorrigibly biased data due to the stochastic nature (inherent or not) of some variables. After a brief review of relevant literature in section2, we will carry out an econometric analysis to identify appropriate specifications to account for stochastic variables in discrete choice modelling (section 3). Then, in section 4 the performance of some specifications arising from the econometric analysis will be tested and compared in terms of parameter estimate bias, computation of marginal rates of substitution and forecasting ability. Finally, section 0 presents our main conclusions.

## THE PROBLEM OF ERRORS IN VARIABLES (EIV)

Much of the effort to specify stochastic variables when estimating econometric models has arisen from the need to solve the EIV problem. In this sense, although there is a vast literature in the case of regression models, research underlying EIV within discrete choice models is scarce, but has shown lately some significant progress. For instance, in the fields of biology and medicine, the EIV problem has been explored in the case of binary models, proposing adaptations of maximum likelihood estimators for specific circumstances (Carroll *et al*., 1984; Stefansky and Carroll, 1985; 1987; 1990; Schnell and Kao, 1987). More generally, Steinmetz and Brownstone (2005) presented a model that considers EIV using multiple imputations that can be used when there is accurate information for a subsample of observations. More recently, Yamamoto and

Komori (2010) estimated a latent class model for handling errors when measuring access distances to public transport.

Walker *et al*. (2010) proposed a hybrid choice model, which includes a latent variable to account for travel time measurement errors given that they were obtained from a network zone based model. In their specification, the true travel time is treated as a latent variable observed through the modeller measured travel time (i.e. the modeller measured travel time is taken as an indicator of the true travel time). Although the theoretical specification of the model is consistent, it is not easy to justify this formulation in practice, especially when high dispersion exists in travel time values, which is the most frequent case. Even more, future indicators are no longer needed for forecasting scenarios.

Other examples of hybrid choice models which can be used to deal with the EIV problem can be found in Bolduc and Alvarez-Daziano (2009) and Brey and Walker (2011), where latent variables are used to account for measurement errors in variables such as the income in a vehicle choice experiment, or the preferred departure time in an airline itinerary choice context. In both cases, the structural equations of the latent variable component are a function of individual characteristics; this has the added benefit of being favourable for forecasting scenarios because the structural equations are used to predict the latent variables values. Nevertheless, the specification of a structural equation associated with individual characteristics in the case of an exogenous variable such as travel time is less natural.

Wansbeek and Meijer (2000) proposed the use of latent variables for treating EIV simply by adding directly the measurement equations to the utility function, without necessarily specifying an additional structural function. The econometric analysis developed in this research builds upon the above idea, keeping in mind variables such as travel time.

## ECONOMETRIC ANALYSIS

### The EIV Problem

Let us consider the following additive and linear in parameters specification for the utility function in a discrete choice context based on random utility theory:

$$U_{in} = \sum_{k=1}^{K} \beta_{ink} x_{ink} + \varepsilon_{in} \tag{4}$$

where $U_{in}$ represents the utility of alternative $A_i$ perceived by individual $n$, $x_{ink}$ refers to the value of the $k$th explanatory attribute of alternative $A_i$ for individual $n$, $\beta_{ink}$ is an unknown parameter to be estimated (referring to alternative $A_i$, for individual $n$ and the $k$th explanatory attribute). $\varepsilon_{in}$ is an error term that distributes independently and identically (IID) Gumbel. Individual $n$ chooses alternative $A_i$ if and only if the utility of that alternative is the maximum among the utilities of

alternatives $A_j$ within her choice set (i.e. $U_{in} > U_{jn}$ for all $A_j \neq A_i$). This formulation corresponds to the classical MNL model (Ortuzar and Willumsen, 2011, Chapter 7).

The stochastic nature of the explanatory variables can be expressed as following:

$$x_{ink} = \overline{x_{ink}} + \eta_{ink} \tag{5}$$

where $\overline{x_{ink}}$ is the mean measured value of the variable (i.e. the value that the modeller would typically use) and $\eta_{ink}$ is the discrepancy between this mean measured value and the true value, perceived by the individual. Equation (5) can be seen as a measurement equation in the context of a hybrid choice model which includes latent variables into the discrete choice model (Bolduc *et al.*, 2008). The variation or discrepancy $\eta_{ink}$ is a stochastic component that follows a certain distribution. Replacing (5) in (4) we get:

$$U_{in} = \sum_{k=1}^{K} \beta_{ink} \left( \overline{x_{ink}} + \eta_{ink} \right) + \varepsilon_{in} \tag{6}$$

Equation (6) is equivalent to a Mixed Logit (ML) formulation because the random error component is a mix of a Gumbel distribution and some other distributions contained in $\eta_{ink}$ (McFadden and Train, 2000). The above suggests that the EIV problem can be approximated through a particular specification of the ML model, and depends on the definition of $\eta_{ink}$. Note that the use of MNL models when stochastic variables are present would be a wrong approach because the IID Gumbel error term of the model cannot represent the full error structure given by the data. Below we will discuss two particular ML formulations for dealing with the EIV problem.

## Stochastic Variables Model

The stochastic variations $\eta_{ink}$ can be specified as follows:

$$\eta_{ink} = \sigma_{ik} \cdot u_{ink} \tag{7}$$

where $u_{ink}$, the components of a vector $\boldsymbol{u}$, are independently distributed variables with zero mean and unitary standard deviation; $\sigma_{ik}$ is a fixed real number representing the standard deviation of the probability function related with the stochastic variation; this value is alternative specific but constant among individuals. Replacing (7) in (6), and assuming generic tastes in the population, the Stochastic Variables (SV) model can be written as:

$$U_{in} = \sum_{k=1}^{K} \beta_{ik} \left( \overline{x_{ink}} + \sigma_{ik} \cdot u_{ink} \right) + \varepsilon_{in} \tag{8}$$

and the probability $P_n$ that individual $n$ chooses alternative $A_i$ for given values $u_{ink}^d$ can be computed as a class of *logit* model:

$$P_n\left(A_i \middle| \boldsymbol{\beta}, \boldsymbol{\sigma}, \boldsymbol{u}^d\right) = \frac{e^{\mu \sum_{k=1}^{K} \beta_{ik}\left(\overline{x_{nik}} + \sigma_{ik} \cdot u_{nik}^d\right)}}{\displaystyle\sum_{A_j \in C_{(n)}} e^{\mu \sum_{k=1}^{K} \beta_{jk}\left(\overline{x_{njk}} + \sigma_{jk} \cdot u_{njk}^d\right)}} \qquad (9)$$

where $\boldsymbol{\beta}$ and $\boldsymbol{\sigma}$ are vectors with elements $\beta_{ik}$ and $\sigma_{ik}$, $\mu$ is the scale factor and $C_{(n)}$ the individual's choice set. The final choice probability can be obtained by integrating (11) over the range of $\boldsymbol{u}$ values. For example, if $\boldsymbol{u}$ is assumed to distribute standard normal the probability that individual $n$ chooses alternative $A_i$ can be expressed as follows:

$$P_{in} = \int_{\boldsymbol{u}} P_n\left(A_i \middle| \boldsymbol{\beta}, \boldsymbol{\sigma}, \boldsymbol{u}\right) \cdot \phi(\boldsymbol{u}) \, d\boldsymbol{u} \qquad (10)$$

which has a typical Mixed Logit (ML) form and thus can be computed using simulated maximum log-likelihood techniques (Train, 2009) as in (13):

$$\overline{P_{in}} = \frac{\sum_{d=1}^{D} P_n\left(A_i \middle| \boldsymbol{\beta}, \boldsymbol{\sigma}, \boldsymbol{u}^d\right)}{D} \qquad (11)$$

for a set of $D$ draws from the distribution of $\boldsymbol{u}$.

It can be shown that in this case due to identifiability issues, for $I$ alternatives it is only possible to estimate $I - 1$ covariance matrix parameters. As a consequence, the modeller has to choose an alternative where all the correspondent covariance matrix parameters must be normalized; this normalization is not simple as it must guarantee that these parameters are fixed at a sufficiently large value in relation to $\boldsymbol{\sigma}$. In the most general case, these values are unknown, and a trial and error procedure needs to be undertaken.

Note that equation (8) can be arranged and set also as follows:

$$U_{in} = \sum_{k=1}^{K} \beta_{ik} \cdot \overline{x_{ink}} + \nu_i + \varepsilon_{in} \qquad (12)$$

where:

$$\nu_i = \sum_{k=1}^{K} \beta_{ik} \cdot \sigma_{ik} \cdot u_{ink} \qquad (13)$$

Equation (12) has an error components ML (ECML) structure, frequently used to account for heteroscedasticity among alternatives (Ortúzar and Willumsen, 2011, section 7.6). Under some circumstances and specially when treating the EIV problem, the SV and ECML structures are mathematically equivalent (i.e. in practice, this equivalence may be interpreted as meaning that the stochastic variables cause heteroscedasticity among alternatives). Furthermore, the effect of the stochastic variables can also be confounded with other sources of heteroscedasticity or heterogeneity in the estimation process.

Due to the equivalence between the SV and ECML structures, instead of estimating the $\sigma_{ik}$ in the SV model we can estimate an alternative specific parameter $\Phi_i$ for each alternative in the ECML. These parameters correspond to the standard deviations of the error term components:

$$v_i = \Phi_i \cdot u_i \tag{14}$$

where $u_i$ distributes Normal with zero mean and unitary standard deviation. The $\Phi_i$ parameters can be expressed as a combination of the attribute variations per alternative:

$$\Phi_i^2 = \sum_{k=1}^{K} \beta_{ik}^2 \cdot \sigma_{ik}^2 \tag{15}$$

Similarly to the SV model, in the ECML structure only $I - 1$ error component variances can be estimated and one of those needs to be normalized. The usual normalization methodology involves estimating a non-identifiable model (i.e. with $I$ variances), and later estimate a new model but fixing to zero the lowest variance in the preliminary estimation (Walker, 2001). As this normalization is easier in the ECML model than in the SV model, it would appear that estimating the former structure is preferable.

## Random Parameters Model

Stochastic variations in the SV model have been assumed to be independent of the corresponding attribute values. Furthermore, the standard deviations $\sigma_{ik}$ are equal for all individuals implying homoscedasticity across observations. However, under some circumstances it is reasonable to expect that the higher an attribute value the larger should be its level of randomness. If such is the case, a proportional direct relationship may be specified:

$$\eta_{ink} = \lambda_{ink} \cdot \overline{x_{ink}} \tag{16}$$

where $\lambda_{ink}$ follows a probability distribution function (pdf) with zero mean and unknown standard deviation $\theta_{ik}$:

$$\lambda_{ink} = \theta_{ik} \cdot u_{ink} \tag{17}$$

In this case then, heteroscedasticity across both respondents and alternatives are found in the model. Replacing (16) in (6), and assuming taste homogeneity, we can write:

$$U_{in} = \sum_{k=1}^{K} \alpha_{ink} \cdot \overline{x_{ink}} + \varepsilon_{in} \tag{18}$$

where:

$$\alpha_{ink} = \beta_{ik} \left(1 + \theta_{ik} \cdot u_{ink}\right) \tag{19}$$

Even though the taste variations across respondents are fixed equation (18) represents a random coefficients (RC) model. Due to formulation equivalence, it is possible that RC model estimates could be confounded with (apparent) random taste heterogeneity in the population if stochastic variables are included in the formulation, and randomness is directly proportional to variable size. However, confounded effects in ML estimates are possible (see Cherchi and Ortúzar, 2008;

Swait and Bernardino, 2000), when stochastic variables in the model are interpreted as taste heterogeneity, even when stochastic variations are independent from the variable size.

# EXPERIMENTAL ANALYSIS

We adopted the classic procedure of Williams and Ortúzar (1982) to generate synthetic databases from known parameters; then we estimated different discrete choice models using this data. The estimated parameters were compared with the known model parameters and WTP measures were computed. In addition, models were tested in terms of their response properties.

## Synthetic Population Generation

A collection of datasets was generated in which pseudo-observed individuals behaved according to a known choice rule, with a defined error structure, and had to choose among three options ($A_i$) labelled: Taxi ($i = 1$), Bus ($i = 2$) and Metro ($i = 3$). Attributes included in the choice set description were: Cost ($c$), travel time ($t$) and access time ($a$). The simulated (i.e. pseudo observed) choices for each individual $n$, represented the alternative with higher associated utility $U_{in}$, which was computed as:

$$U_{in} = \beta_{cost} \cdot c_{in} + \beta_{time} \cdot t_{in} + \beta_{access} \cdot a_{in} + \varepsilon_{in} \tag{20}$$

where $\beta_{cost}$, $\beta_{time}$ and $\beta_{access}$ are fixed cost, travel time and access time parameters respectively. The random error terms, $\varepsilon_{in}$, were generated from an *iid* standard Gumbel pdf. In turn, attribute randomness was included as follows:

$$c_{in} = \overline{c_{in}} + \eta_{inc} \tag{21}$$

$$t_{in} = \overline{t_{in}} + \eta_{int} \tag{22}$$

$$a_{in} = \overline{a_{in}} + \eta_{ina} \tag{23}$$

where the first terms represent base values (i.e. modeller measured values), and the second ones correspond to Normal distributed error terms. The base values were, in turn, generated from a truncated Normal distributed pdf to avoid negative attribute values. The final attribute values ($c_{in}$, $t_{in}$ and $a_{in}$) were also truncated according to certain minimum allowed thresholds. Table 1 shows the base attribute values and taste parameters used in the synthetic sample generation process.

Synthetic samples with 500, 2000 and 5000 observations were generated for each of three levels of randomness (i.e. nine synthetic samples). Six out of the nine attributes were of a stochastic nature in the first level (I). In the second level (II), two attributes were stochastic. Finally, in the third level (III) only one attribute was considered stochastic. Stochastic variations were obtained

from a normal pdf with zero mean and a standard deviation computed as a percentage of the corresponding mean base value. The standard deviations used are presented in Table 2.

Table 1 Attribute base values and parameters used in synthetic sample generation

| Attribute | | Taxi | Bus | Metro | Parameter |
|---|---|---|---|---|---|
| *Cost* | *Mean* | 50 | 20 | 22 | |
| | *Standard deviation* | 5 | 4 | 3 | -0.08 |
| | *Minimum value* | 25 | 15 | 12 | |
| *Travel Time* | *Mean* | 15 | 30 | 16 | |
| | *Standard deviation* | 4 | 10 | 3 | -0.12 |
| | *Minimum value* | 7 | 10 | 4 | |
| *Access Time* | *Mean* | 5 | 10 | 18 | |
| | *Standard deviation* | 2 | 3 | 4 | -0.16 |
| | *Minimum value* | 0.5 | 2 | 2 | |

The taxi costs and travel times by bus are the attributes with higher variation in levels of randomness I and II. In terms of global randomness in level I, bus and taxi were the alternatives with higher variability, with standard deviations of $1.18^2$ and 0.99, respectively; on the other hand, the Metro alternative had the lowest variation (standard deviation of 0.69). In level II, the bus and taxi utilities have the higher global standard deviations (1.08 and 0.80), while the metro alternative has no variation. Finally, in level III only the bus alternative has variation with a standard deviation of 1.08.

Table 2 Standard deviations for attribute stochastic variations

| Attribute | Mode | *Level of Randomness* | | |
|---|---|---|---|---|
| | | *I* | *II* | *III* |
| **Cost** | *Taxi* | **10.0** (20%) | **10.0** (20%) | - |
| | *Bus* | - | - | - |
| | *Metro* | - | - | - |
| **Travel Time** | *Taxi* | **4.5** (30%) | - | - |
| | *Bus* | **9.0** (30%) | **9.0** (30%) | **9.0** (30%) |
| | *Metro* | **3.2** (20%) | - | - |
| **Access Time** | *Taxi* | **1.5** (30%) | - | - |
| | *Bus* | **3.0** (30%) | - | - |
| | *Metro* | **3.6** (20%) | - | - |

## Specification of the Estimated Models

Some discrete choice models were estimated using the synthetic database generated by the procedure described in the previous section. The specifications of the model structures tested are shown in equations (26) –(30).

---

[2] $SD_{Ubus} = \sqrt{9^2(-0.12)^2 + 3^2(-0.16)^2} = 1.18$; similar computations were performed for the other alternatives.

*MNL model*

$$U_i = \beta_c \cdot c_i + \beta_t \cdot t_i + \beta_a \cdot a_i + \varepsilon_i \qquad (26)$$

*RCg    (Random    coefficients    model    with    generic    coefficients)*

$$U_i = \left(\beta_c + \sigma_{\cos t} \cdot u_{\cos t}\right) \cdot c_i + \left(\beta_t + \sigma_{time} \cdot u_{time}\right) \cdot t_i + \left(\beta_a + \sigma_{acces} \cdot u_{acces}\right) a_i + \varepsilon_i \qquad (27)$$

*RCs    (Random    coefficients    model    with    specific    coefficients)*

$$U_i = \left(\beta_c + \phi_{c,i} \cdot u_{c,i}\right) \cdot c_i + \left(\beta_t + \phi_{t,i} \cdot u_{t,i}\right) \cdot t_i + \left(\beta_a + \phi_{a,i} \cdot u_{a,i}\right) \cdot a_i + \varepsilon_i \qquad (28)$$

*SV (Stochastic Variables model)*

$$U_i = \beta_c \cdot \left(c_i + \phi_{c,i} \cdot u_{c,i}\right) + \beta_t \cdot \left(t_i + \phi_{t,i} \cdot u_{t,i}\right) + \beta_a \cdot \left(a_i + \phi_{a,i} \cdot u_{a,i}\right) + \varepsilon_i \qquad (29)$$

*ECML (Error Components Mixed Logit model)*

$$U_i = \beta_c \cdot c_i + \beta_t \cdot t_i + \beta_a \cdot a_i + \xi_i \cdot u_i + \varepsilon_i \qquad (30)$$

All random terms *u* in the utility functions above were specified as standard Normal, whilst the error terms $\varepsilon$ were considered to distribute iid Gumbel. For the SV model estimation, the Taxi Cost and Bus Travel Time error terms were estimated. All remaining error terms were fixed to their original values for identification issues of the SV model. For the ECML model, all the error component terms where estimated, but for identification issues the Metro error components were fixed at zero in all cases. On the other hand, the Bus and Metro Cost parameters were fixed at zero in the RCs specification because they were built without any randomness.

## 4.3 PARAMETER ESTIMATION

Estimation results for the levels of randomness I, II and III are given in tables 3 to 5; the following null hypotheses were statistically tested through the estimation analysis:

*Hypothesis I: $\beta_{estimated} = 0$*, in order to evaluate parameter significance; the corresponding *t* ratios are shown in parentheses within the tables.
*Hypothesis II: $(\beta_{estimated} - \beta_{true}) = 0$*, to evaluate each model's capacity to recover the population parameter values; the corresponding *t* ratios are shown in square brackets. The critical value for comparison is 1.96 for a 5% level.

The ratios between estimated and true parameters ($\beta_{estimated}/\beta_{true}$) are presented in curly brackets. The log-likelihood (LL) values are shown as model fit measures and the likelihood ratio (LR) test was used for model comparison (Ortuzar and Willumsen, 2011, page 279), considering the MNL as the restricted model; the computed LR value was contrasted with the critical $\chi^2$ values at

the 5% significance level and degrees of freedom defined by the difference in the number of parameters between both models (critical $\chi^2$ values are shown with an asterisk within the tables). Parameters with two asterisks were not estimated and therefore fixed a-priori.

The MNL estimation results differ significantly from the original base parameter values, for the three variation levels and the three generated sample sizes. On the other hand, the RCg model can neither recover the true parameters except for the sample of only 500 observations where the travel time and access time parameters were unbiased for levels II and III. In this case also we found significant taste differences for the travel time attribute, especially in samples with lower levels of randomness, suggesting that confounding effects could be an issue when attribute variability is considered in random coefficients model estimation (the travel time parameter absorbs part of that randomness). However, the RCg model only has a significantly better fit than the MNL for level III and higher sample sizes (2000 and 5000 observations).

The RCs model was able to recover the original parameter values with little success; it presents the best result so far for the lowest level of randomness (III). In addition, the RCs model has a better fit to the data than the MNL and RCg models in level of randomness III and III when larger sample sizes are used (i.e. 2000 and 5000 observations). As the travel time parameters estimated with the largest sample size (5000 observations) were found to be significantly different from zero, these model results suggest that random effects can indeed be captured using random taste coefficients. On the other hand, few significant access time ($\beta_{access}$) and travel time ($\beta_{time}$) parameters were found in the RCs estimations with 500 observations.

Clearly good results are obtained for the SV model specification, which is in line with the econometric analysis, when estimation is performed with the more reasonable sample sizes (2000 and 5000). In those cases, the SV model can recover the original parameters and have a better model fit than the MNL model, especially for levels of randomness with smaller number of stochastic variables.

The Normal distributed parameters in the SV specification resulted to be significantly different from zero and not different from the originally assumed values. In contrast with this result, at least one of the parameters was found to be significantly different from the original parameter value in the ECML model estimations, except when the model was estimated with 500 observations or 2000 observations and the level of randomness I.

Both the SV and ECML specifications could have some problems when used for estimation with not too many observations (500). Firstly, for levels of randomness I and II, the population taste parameters were found to be non-significant due to their high standard errors in estimation. Furthermore, most of the covariance matrix error component parameters were found to have low significance or different from the original parameter values. Finally for level III, the SV model was unable to find any randomness collapsing to a MNL structure.

Table 3 Parameter estimates for level of randomness I

| Parameter | 500 Observations | | | | | 2000 Observations | | | | | 5000 Observations | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | MNL | RCg | RCs | SV | ECML | MNL | RCg | RCs | SV | ECML | MNL | RCg | RCs | SV | ECML |
| *Cost (-0.08)* | **-0.056** | **-0.058** | **-0.85** | **-0.341** | **-0.093** | **-0.058** | **-0.060** | **-0.0989** | **-0.077** | **-0.065** | **-0.062** | **-0.062** | **-0.0981** | **-0.089** | **-0.072** |
| | (-8.53) | (-8.20) | (-0.59) | (-0.32) | (-1.62) | (-17.56) | (-16.83) | (-2.39) | (-5.34) | (-6.96) | (-28.91) | (-27.59) | (-3.76) | (-6.97) | (-10.69) |
| | [3.66] | [3.15] | [-0.53] | [-0.24] | [-0.23] | [6.49] | [5.76] | [0.46] | [0.21] | [1.59] | [8.46] | [7.91] | [-0.69] | [-0.71] | [1.21] |
| | {0.70} | {0.72} | {10.60} | {4.26} | {1.17} | {0.73} | {0.75} | {1.24} | {0.96} | {0.81} | {0.77} | {0.78} | {1.23} | {1.11} | {0.90} |
| *Travel time (-0.12)* | **-0.084** | **-0.093** | **-1.39** | **-0.592** | **-0.163** | **-0.088** | **-0.094** | **-0.152** | **-0.122** | **-0.103** | **-0.091** | **-0.092** | **-0.146** | **-0.136** | **-0.109** |
| | (-9.23) | (-6.59) | (-0.60) | (-0.32) | (-1.60) | (-19.19) | (-13.18) | (-2.4) | (-4.96) | (-6.32) | (-30.80) | (-21.47) | (-3.73) | (-6.59) | (-9.91) |
| | [3.91] | [1.91] | [-0.55] | [-0.26] | [-0.43] | [6.95] | [3.59] | [-0.51] | [-0.08] | [1.04] | [10.03] | [6.54] | [-0.66] | [-0.78] | [1.00] |
| | {0.70} | {0.78} | {11.58} | {4.93} | {1.36} | {0.73} | {0.79} | {1.27} | {1.02} | {0.86} | {0.75} | {0.77} | {1.22} | {1.13} | {0.91} |
| *Access time (-0.16)* | **-0.120** | **-0.124** | **-1.70** | **-0.648** | **-0.183** | **-0.120** | **-0.122** | **-0.199** | **-0.153** | **-0.130** | **-0.126** | **-0.126** | **-0.191** | **-0.173** | **-0.141** |
| | (-8.56) | (-8.25) | (-0.59) | (-0.32) | (-1.86) | (-16.98) | (-16.45) | (-2.45) | (-6.33) | (-8.75) | (-27.74) | (-26.84) | (-3.81) | (-8.13) | (-13.27) |
| | [2.86] | [2.40] | [-0.53] | [-0.24] | [-0.23] | [5.67] | [5.13] | [-0.48] | [0.29] | [2.01] | [7.49] | [7.22] | [0.62] | [-0.61] | [1.79] |
| | {0.75} | {0.78} | {10.63} | {4.05} | {1.14} | {0.75} | {0.76} | {1.24} | {0.96} | {0.81} | {0.79} | {0.79} | {1.19} | {1.08} | {0.88} |
| $\sigma_{cost}$: *RCg* | | **-2.58E-17** | | | | | **-1.29E-17** | | | | | **2.59E-17** | | | |
| | | (0.00) | | | | | (0.00) | | | | | (0.00) | | | |
| $\sigma_{time}$: *RCg* | | **0.0377** | | | | | **0.0306** | | | | | **-0.0145** | | | |
| | | (1.47) | | | | | (2.12) | | | | | (-0.93) | | | |
| $\sigma_{access}$: *RCg* | | **0.0023** | | | | | **0.0033** | | | | | **0.0022** | | | |
| | | (0.05) | | | | | (0.13) | | | | | (0.11) | | | |
| $\phi_{c,1}$: *RCs & SV (10)* | | | **-0.39** | **23.1** | | | | **-0.0337** | **11.8** | | | | **-0.0273** | **14** | |
| | | | (-0.59) | (4.57) | | | | (-1.45) | (1.80) | | | | (-2.04) | (4.11) | |
| $\phi_{c,2}$: *RCs & SV* | | | **0.72** | **0\*\*** | | | | **-0.0493** | **0\*\*** | | | | **-0.0513** | **0\*\*** | |
| | | | (0.57) | | | | | (-1.1) | | | | | (-1.7) | | |
| $\phi_{c,3}$: *RCs & SV* | | | **-0.10** | **0\*\*** | | | | **0.0694** | **0\*\*** | | | | **0.0489** | **0\*\*** | |
| | | | (-0.31) | | | | | (1.47) | | | | | (1.72) | | |
| $\phi_{t,1}$: *RCs & SV* | | | **0.09** | **4.50\*\*** | | | | **0.0000216** | **4.50\*\*** | | | | **0.0547** | **4.50\*\*** | |
| | | | (0.32) | | | | | (0.00) | | | | | (1.69) | | |
| $\phi_{t,2}$: *RCs & SV (9)* | | | **-0.28** | **18.10** | | | | **-0.0372** | **11.00** | | | | **-0.0321** | **11.30** | |
| | | | (-0.61) | (6.66) | | | | (-1.36) | (3.42) | | | | (-2.09) | (5.87) | |
| $\phi_{t,3}$: *RCs & SV* | | | **-0.11** | **3.20\*\*** | | | | **-0.00817** | **3.20\*\*** | | | | **0.0296** | **3.20\*\*** | |
| | | | (-0.44) | | | | | (-0.14) | | | | | (0.85) | | |
| $\phi_{a,1}$: *RCs & SV* | | | **-0.14** | **1.50\*\*** | | | | **-0.0228** | **1.50\*\*** | | | | **0.041** | **1.50\*\*** | |
| | | | (-0.29) | | | | | (-0.27) | | | | | (0.61) | | |
| $\phi_{a,2}$: *RCs & SV* | | | **1.63** | **3.00\*\*** | | | | **0.0914** | **3.00\*\*** | | | | **0.0721** | **3.00\*\*** | |
| | | | (0.62) | | | | | (1.15) | | | | | (1.56) | | |
| $\phi_{a,3}$: *RCs & SV* | | | **0.65** | **3.60\*\*** | | | | **0.00471** | **3.60\*\*** | | | | **-0.00117** | **3.60\*\*** | |
| | | | (0.57) | | | | | (0.11) | | | | | (-0.04) | | |
| $\xi_1$: *Error component* | | | | | **1.84** | | | | | **-0.62** | | | | | **0.83** |
| | | | | | (0.89) | | | | | (-1.03) | | | | | (2.42) |
| $\xi_2$: *Error component* | | | | | **2.61** | | | | | **0.96** | | | | | **-1.00** |
| | | | | | (1.19) | | | | | (1.82) | | | | | (-3.01) |
| $\xi_3$: *Error component* | | | | | **0.00** | | | | | **0.00** | | | | | **0.00** |
| *Log-likelihood* | **-489.26** | **-488.78** | **-486.21** | **-486.75** | **-487.40** | **-1941.25** | **-1940.40** | **-1938.34** | **-1939.71** | **-1940.43** | **-4813.25** | **-4813.13** | **-4807.95** | **-4808.69** | **-4810.81** |
| *Log-likelihood ratio test* | | 0.97 | 6.10 | 5.03 | 3.72 | | 1.70 | 5.82 | 3.07 | 1.65 | | 0.24 | 10.6 | 9.13 | 4.89 |
| *Critical $\chi^2$ at 5%* | | 7.81* | 14.07* | 5.99* | 5.99* | | 7.81* | 14.07* | 5.99* | 5.99* | | 7.81* | 14.07* | 5.99* | 5.99* |

Table 4 Parameter estimates for level of randomness II

| Parameter | 500 Observations | | | | | 2000 Observations | | | | | 5000 Observations | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | MNL | RCg | RCs | SV | ECML | MNL | RCg | RCs | SV | ECML | MNL | RCg | RCs | SV | ECML |
| Cost (-0.08) | -0.059 | -0.0625 | -0.641 | -0.123 | -0.114 | -0.062 | -0.064 | -0.794 | -0.073 | -0.067 | -0.064 | -0.0648 | -0.113 | -0.077 | -0.071 |
| | (-8.82) | (-8.46) | (-0.71) | (-1.41) | (-1.43) | (-18.25) | (-17.39) | (-0.64) | (-6.78) | (-14.60) | (-29.49) | (-28.04) | (-3.2) | (-9.89) | (-10.93) |
| | [3.14] | [2.37] | [0.62] | [-0.49] | [-0.43] | [5.46] | [4.35] | [-0.58] | [0.70] | [2.97] | [7.59] | [6.58] | [-0.93] | [0.40] | [1.35] |
| | {0.74} | {0.78} | {8.01} | {1.54} | {1.43} | {0.77} | {0.80} | {9.93} | {0.91} | {0.83} | {0.80} | {0.81} | {1.41} | {0.96} | {0.89} |
| Travel time (-0.12) | -0.097 | -0.114 | -1.18 | -0.226 | -0.213 | -0.092 | -0.104 | -1.31 | -0.119 | -0.109 | -0.093 | -0.098 | -0.176 | -0.121 | -0.112 |
| | (-10.05) | (-7.00) | (-1.40) | (-1.40) | (-1.43) | (-19.74) | (-13.39) | (-0.57) | (-6.29) | (-10.03) | (-31.27) | (-21.47) | (-3.19) | (-9.22) | (-10.19) |
| | [2.38] | [0.37] | [-0.65] | [-0.66] | [-0.62] | [5.99] | [2.06] | [-0.57] | [0.05] | [1.01] | [9.16] | [4.82] | [-1.02] | [-0.08] | [0.73] |
| | {0.81} | {0.95} | {9.83} | {1.88} | {1.78} | {0.77} | {0.87} | {10.92} | {0.99} | {0.91} | {0.77} | {0.82} | {1.47} | {1.01} | {0.93} |
| Access time (-0.16) | -0.128 | -0.135 | -1.29 | -0.239 | -0.226 | -0.129 | -0.133 | -1.68 | -0.147 | -0.138 | -0.131 | -0.133 | -0.229 | -0.153 | -0.145 |
| | (-8.93) | (-8.58) | (-0.71) | (-1.54) | (-1.57) | (-17.86) | (-17.23) | (-0.64) | (-8.27) | (-14.77) | (-28.61) | (-27.6) | (-3.25) | (-12.10) | (-13.62) |
| | [2.24] | [1.59] | [-0.62] | [-0.51] | [-0.46] | [4.31] | [3.51] | [-0.58] | [0.73] | [2.35] | [6.32] | [5.60] | [-0.98] | [0.55] | [1.42] |
| | {0.80} | {0.84} | {8.06} | {1.49} | {1.41} | {0.81} | {0.83} | {10.50} | {0.92} | {0.86} | {0.82} | {0.83} | {1.43} | {0.96} | {0.91} |
| $\sigma_{cost}$: RCg | | 1.26E-17 | | | | | 1.29E-17 | | | | | -1.56E-17 | | | |
| | | (0.00) | | | | | (0.00) | | | | | (0.00) | | | |
| $\sigma_{time}$: RCg | | 0.0535 | | | | | 0.043 | | | | | 0.028 | | | |
| | | (2.22) | | | | | (3.35) | | | | | (2.87) | | | |
| $\sigma_{access}$: RCg | | -0.000151 | | | | | -0.000106 | | | | | 0.000955 | | | |
| | | (0.00) | | | | | (0.00) | | | | | (0.05) | | | |
| $\phi_{c,1}$: RCs & SV (10) | | | -0.259 | 21.4 | | | | -0.278 | 9.23 | | | | -0.0306 | 10.8 | |
| | | | (-0.67) | (2.87) | | | | (-0.64) | (1.33) | | | | (-1088) | (2.87) | |
| $\phi_{c,2}$: RCs & SV | | | -0.52 | 0** | | | | -0.712 | 0** | | | | -0.0756 | 0** | |
| | | | (-0.66) | | | | | (-0.63) | | | | | (-2.06) | | |
| $\phi_{c,3}$: RCs & SV | | | -0.33 | 0** | | | | 0.611 | 0** | | | | 0.07 | 0** | |
| | | | (-0.62) | | | | | (0.62) | | | | | (1.99) | | |
| $\phi_{t,1}$: RCs & SV | | | 0.369 | 0** | | | | 0.199 | 0** | | | | 0.0466 | 0** | |
| | | | (0.65) | | | | | (0.49) | | | | | (1.29) | | |
| $\phi_{t,2}$: RCs & SV (9) | | | -0.299 | -14.90 | | | | -0.362 | 11.50 | | | | -0.0399 | 11.20 | |
| | | | (-0.72) | (-4.83) | | | | (-0.59) | (4.49) | | | | (-2.13) | (6.62) | |
| $\phi_{t,3}$: RCs & SV | | | -0.0705 | 0** | | | | 0.278 | 0** | | | | 0.0145 | 0** | |
| | | | (-0.29) | | | | | (0.68) | | | | | (0.18) | | |
| $\phi_{a,1}$: RCs & SV | | | 0.415 | 0** | | | | 0.0369 | 0** | | | | 0.0658 | 0** | |
| | | | (0.46) | | | | | (0.07) | | | | | (0.77) | | |
| $\phi_{a,2}$: RCs & SV | | | 1.05 | 0** | | | | 0.781 | 0** | | | | 0.1 | 0** | |
| | | | (0.77) | | | | | (0.59) | | | | | (1.72) | | |
| $\phi_{a,3}$: RCs & SV | | | -0.251 | 0** | | | | -0.363 | 0** | | | | 0.00505 | 0** | |
| | | | (-0.61) | | | | | (-0.64) | | | | | (0.17) | | |
| $\xi_1$: Error component | | | | | 2.35 | | | | | 0.00 | | | | | 0.49 |
| | | | | | (0.94) | | | | | (0.00) | | | | | (1.07) |
| $\xi_2$: Error component | | | | | 3.14 | | | | | -1.13 | | | | | -1.14 |
| | | | | | (1.14) | | | | | (-3.01) | | | | | (-3.70) |
| $\xi_3$: Error component | | | | | 0.00 | | | | | 0.00 | | | | | 0.00 |
| Log-likelihood | -477.40 | -476.13 | -474.77 | -475.00 | -475.01 | -1918.46 | -1916.12 | -1910.60 | -1915.73 | -1916.68 | -4779.25 | -4777.91 | -4771.52 | -4773.30 | -4775.33 |
| Log-likelihood ratio test | | 2.54 | 5.26 | 4.80 | 4.80 | | 4.66 | 15.72 | 5.45 | 3.56 | | 2.68 | 15.46 | 11.91 | 7.85 |
| Critical $\chi^2$ at 5% | | 5.99* | 5.99* | 5.99* | 5.99* | | 5.99* | 5.99* | 5.99* | 5.99* | | 5.99* | 5.99* | 5.99* | 5.99* |

Table 5 Parameter estimates for level of randomness III

| Parameter | 500 Observations | | | | | 2000 Observations | | | | | 5000 Observations | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | *MNL* | *RCg* | *RCs* | *SV* | *ECML* | *MNL* | *RCg* | *RCs* | *SV* | *ECML* | *MNL* | *RCg* | *RCs* | *SV* | *ECML* |
| *Cost (-0.08)* | **-0.060** | **-0.063** | **-0.098** | **-0.060** | **-0.118** | **-0.065** | **-0.0671** | **-0.115** | **-0.070** | **-0.0713** | **-0.068** | **-0.0704** | **-0.0924** | **-0.075** | **-0.0746** |
| | (-9.00) | (-8.61) | (-1.57) | (-9.00) | (-1.5) | (-18.95) | (-17.92) | (-2.33) | (-15.04) | (-9.51) | (-31.00) | (-29.15) | (-4.83) | (-24.40) | (-22.04) |
| | [2.96] | [2.27] | [-0.29] | [2.96] | [-0.48] | [4.52] | [3.44] | [-0.71] | [2.05] | [1.16] | [5.32] | [3.98] | [-0.65] | [1.73] | [1.60] |
| | {0.75} | {0.79} | {1.23} | {0.75} | {-1.48} | {0.81} | {0.84} | | {0.88} | {-0.89} | {0.85} | {0.88} | {1.16} | {0.93} | {-0.93} |
| *Travel time (-0.12)* | **-0.092** | **-0.106** | **-0.169** | **-0.092** | **-0.205** | **-0.091** | **-0.101** | **-0.183** | **-0.109** | **-0.11** | **-0.097** | **-0.106** | **-0.146** | **-0.116** | **-0.116** |
| | (-9.79) | (-6.83) | (-1.54) | (-9.79) | (-1.51) | (-19.63) | (-13.52) | (-2.31) | (-10.55) | (-8.19) | (-32.20) | (-22.13) | (-4.65) | (-17.53) | (-16.41) |
| | [2.96] | [0.90] | [-0.45] | [2.96] | [-0.63] | [6.34] | [2.54] | [-0.80] | [1.07] | [0.74] | [7.61] | [2.92] | [-0.83] | [0.60] | [0.57] |
| | {0.77} | {0.88} | {1.41} | {0.77} | {-1.71} | {0.76} | {0.84} | {1.53} | {0.91} | {-0.92} | {0.81} | {0.88} | {1.22} | {0.97} | {-0.97} |
| *Access time (-0.16)* | **-0.125** | **-0.131** | **-0.19** | **-0.125** | **-0.223** | **-0.131** | **-0.135** | **-0.236** | **-0.143** | **-0.144** | **-0.137** | **-0.141** | **-0.187** | **-0.149** | **-0.149** |
| | (-8.82) | (-8.45) | (-1.68) | (-8.82) | (-1.65) | (-18.15) | (-17.41) | (-2.35) | (-14.86) | (-11.11) | (-29.44) | (-28.13) | (-4.87) | (-23.99) | (-22.58) |
| | [2.46] | [1.87] | [-0.27] | [2.46] | [-0.47] | [4.02] | [3.22] | [-0.76] | [1.77] | [1.23] | [4.95] | [3.79] | [-0.70] | [1.77] | [1.67] |
| | {0.78} | {0.82} | {1.19} | {0.78} | {-1.39} | {0.82} | {0.84} | {1.48} | {0.89} | {-0.90} | {0.86} | {0.88} | {1.17} | {0.93} | {-0.93} |
| *σ$_{cost}$: RCg* | | **-1.97E-17** | | | | | **-3.40E-17** | | | | | **4.96E-18** | | | |
| | | (0.0) | | | | | (0.0) | | | | | (0.0) | | | |
| *σ$_{time}$: RCg* | | **-0.0486** | | | | | **-0.0409** | | | | | **-0.0363** | | | |
| | | (1.98) | | | | | (3.22) | | | | | (4.29) | | | |
| *σ$_{access}$: RCg* | | **0.00141** | | | | | **0.00123** | | | | | **0.0017** | | | |
| | | (0.04) | | | | | (0.05) | | | | | (0.08) | | | |
| *φc,1: RCs & SV (10)* | | | **-0.0309** | **0\*\*** | | | | **-0.0149** | **0\*\*** | | | | **-0.00692** | **0\*\*** | |
| | | | (-0.79) | | | | | (-0.51) | | | | | (-0.48) | | |
| *φc,2: RCs & SV* | | | **-0.0848** | **0\*\*** | | | | **-0.099** | **0\*\*** | | | | **-0.0588** | **0\*\*** | |
| | | | (-0.94) | | | | | (-1.7) | | | | | (-2.42) | | |
| *φc,3: RCs & SV* | | | **-0.00403** | **0\*\*** | | | | **0.0729** | **0\*\*** | | | | **0.0419** | **0\*\*** | |
| | | | (-0.07) | | | | | (1.53) | | | | | (1.91) | | |
| *φt,1: RCs & SV* | | | **0.0171** | **0\*\*** | | | | **0.0698** | **0\*\*** | | | | **0.0442** | **0\*\*** | |
| | | | (0.25) | | | | | (1.35) | | | | | (1.52) | | |
| *φt,2: RCs & SV (9)* | | | **-0.0321** | **0.00** | | | | **-0.0477** | **10.80** | | | | **-0.0339** | **10.30** | |
| | | | (-0.71) | (0.00) | | | | (-1.64) | (4.53) | | | | (-2.6) | (7.51) | |
| *φt,3: RCs & SV* | | | **-0.0197** | **0\*\*** | | | | **0.0421** | **0\*\*** | | | | **0.0297** | **0\*\*** | |
| | | | (-0.22) | | | | | (0.64) | | | | | (-2.6) | | |
| *φu,1: RCs & SV* | | | **0.105** | **0\*\*** | | | | **0.0758** | **0\*\*** | | | | **0.0284** | **0\*\*** | |
| | | | (0.43) | | | | | (0.66) | | | | | (0.41) | | |
| *φu,2: RCs & SV* | | | **-0.137** | **0\*\*** | | | | **0.0656** | **0\*\*** | | | | **0.0586** | **0\*\*** | |
| | | | (-0.91) | | | | | (0.61) | | | | | (1.21) | | |
| *φu,3: RCs & SV* | | | **0.00887** | **0\*\*** | | | | **-0.00725** | **0\*\*** | | | | **-0.00613** | **0\*\*** | |
| | | | (0.10) | | | | | (-0.16) | | | | | (-0.18) | | |
| *ξ1: Error component* | | | | | **3.32** | | | | | **-0.564** | | | | | **0.453** |
| | | | | | (0.97) | | | | | (-0.31) | | | | | (0.33) |
| *ξ2: Error component* | | | | | **4.43** | | | | | **4.18** | | | | | **-6.27** |
| | | | | | (1.22) | | | | | (2.97) | | | | | (-5.11) |
| *ξ3: Error component* | | | | | **0.00** | | | | | **0.00** | | | | | **0.00** |
| *Log-likelihood* | -481.75 | -480.79 | -479.78 | -481.75 | -478.79 | -1916.53 | -1914.48 | -1910.22 | -1914.03 | -1914.28 | -4722.41 | -4719.11 | -4714.06 | -4715.96 | -4717.02 |
| *Log-likelihood ratio test* | | 1.92 | 3.94 | 0.00 | 5.92 | | 4.10 | 12.62 | 5.00 | 4.50 | | 6.60 | 16.70 | 12.90 | 10.78 |
| *Critical χ$^2$ at 5%* | | 3.84* | 3.84* | 3.84* | 3.84* | | 3.84* | 3.84* | 3.84* | 3.84* | | 3.84* | 3.84* | 3.84* | 3.84* |

## 4.4 Willingness to Pay Computations

Willingness to pay (WTP) values for travel time (*WTPTV*) and access time (*WTPTA*) are presented in tables 6 and 7. The ratio between the estimated and true WTP values is defined by $r = WTP_{estimated}/ WTP_{true}$. In the case of the RC models, the reported values represent the mean of the parameter distribution. Even though, the *WTPTA* computed from complex models estimated with the largest sample size are closer to the true WTP values, the WTP estimations with the MNL structure are close enough to true values in most cases. However, wrong *WTPTV* were obtained from the RC and SV models estimated with only 500 observations (for the levels of randomness I and II), suggesting again that complex models may encounter problems when they are estimated with low sample sizes (Cherchi and Ortúzar, 2008).

Table 6 Willingness to pay for travel time (WTPTV) computations

| True (Goal) | Level of randomness | Model | 500 obs. Estimate | r | 2000 obs. Estimate | r | 5000 obs. Estimate | r |
|---|---|---|---|---|---|---|---|---|
| 1.5 | I | MNL | 1.51 | 1.00 | 1.51 | 1.01 | 1.46 | 0.97 |
| | | RCg | 1.61 | 1.07 | 1.58 | 1.05 | 1.48 | 0.99 |
| | | RCs | 1.60 | 1.07 | 1.55 | 1.03 | 1.49 | 0.99 |
| | | SV | 1.74 | 1.16 | 1.58 | 1.06 | 1.53 | 1.02 |
| | | ECML | 1.74 | 1.16 | 1.58 | 1.06 | 1.53 | 1.02 |
| | II | MNL | 1.64 | 1.10 | 1.50 | 1.00 | 1.46 | 0.97 |
| | | RCg | 1.82 | 1.22 | 1.63 | 1.08 | 1.51 | 1.01 |
| | | RCs | 1.80 | 1.20 | 1.66 | 1.11 | 1.54 | 1.03 |
| | | SV | 1.84 | 1.22 | 1.64 | 1.09 | 1.57 | 1.05 |
| | | ECML | 1.84 | 1.22 | 1.64 | 1.09 | 1.57 | 1.05 |
| | III | MNL | 1.53 | 1.02 | 1.40 | 0.94 | 1.42 | 0.95 |
| | | RCg | 1.69 | 1.13 | 1.51 | 1.00 | 1.50 | 1.00 |
| | | RCs | 1.69 | 1.13 | 1.54 | 1.03 | 1.55 | 1.03 |
| | | SV | 1.53 | 1.02 | 1.55 | 1.03 | 1.55 | 1.04 |
| | | ECML | 1.53 | 1.02 | 1.55 | 1.03 | 1.55 | 1.04 |

## 4.5 Response Analysis

The forecasting ability (response properties) of the estimated discrete choice models was examined by contrasting the market shares estimated by the various models under different transport policies against the simulated results for the individual choices under the relevant policy conditions (Williams and Ortúzar, 1982; Cantillo *et al.*, 2006). In addition, model forecasts for "perfect" SV and ECML (i.e. using the true parameters) models were also computed.

The various transport policy scenarios proposed to evaluate the models' forecasting ability are shown in Table 8, where we also show the percentage change associated with every attribute in each policy. The first policy, P0, represents the base situation; while P1, P2 and P3 are policies that have differential impacts on those alternatives with higher global levels of randomness (i.e. Bus and Taxi). Policies P4 and P5 represent higher percentage changes in the

cost attribute of a specific alternative (30% changes), and policies P6 and P7 propose simultaneous changes to some alternatives on selected attributes (Cantillo et al., 2010).

Table 7 Willingness to pay for access time (WTPTA) computations

| True (Goal) | Level of randomness | Model | 500 | | 2000 | | 5000 | |
|---|---|---|---|---|---|---|---|---|
| | | | Estimate | r | Estimate | r | Estimate | r |
| 2 | I | MNL | 2.14 | 1.07 | 2.05 | 1.03 | 2.04 | 1.02 |
| | | RCg | 2.15 | 1.07 | 2.05 | 1.02 | 2.03 | 1.01 |
| | | RCs | 2.02 | 1.01 | 2.00 | 1.00 | 1.95 | 0.98 |
| | | SV | 2.02 | 1.01 | 2.00 | 1.00 | 1.95 | 0.98 |
| | | ECML | 1.90 | 0.95 | 1.99 | 0.99 | 1.94 | 0.97 |
| | II | MNL | 2.17 | 1.08 | 2.09 | 1.05 | 2.06 | 1.03 |
| | | RCg | 2.16 | 1.08 | 2.08 | 1.04 | 2.05 | 1.03 |
| | | RCs | 2.06 | 1.03 | 2.05 | 1.03 | 2.03 | 1.01 |
| | | SV | 2.06 | 1.03 | 2.05 | 1.03 | 2.03 | 1.01 |
| | | ECML | 1.94 | 0.97 | 2.03 | 1.01 | 1.99 | 0.99 |
| | III | MNL | 2.08 | 1.04 | 2.03 | 1.01 | 2.01 | 1.00 |
| | | RCg | 2.07 | 1.03 | 2.01 | 1.01 | 2.00 | 1.00 |
| | | RCs | 2.07 | 1.04 | 2.03 | 1.01 | 1.99 | 1.00 |
| | | SV | 2.07 | 1.04 | 2.03 | 1.01 | 1.99 | 1.00 |
| | | ECML | 2.08 | 1.04 | 2.03 | 1.02 | 1.99 | 1.00 |

Table 8 Transport Policies Examined

| Policy | Travel time | | | Cost | | | Access time | | |
|---|---|---|---|---|---|---|---|---|---|
| | Taxi | Bus | Metro | Taxi | Bus | Metro | Taxi | Bus | Metro |
| P0 | | | | | | | | | |
| P1 | | -10% | | | | | | | |
| P2 | | -20% | | | | | | | |
| P3 | | | | -20% | | | | | |
| P4 | | | | | 30% | | | | |
| P5 | | | | | | 30% | | | |
| P6 | | -20% | -20% | 20% | | | | | |
| P7 | | | | 20% | 20% | -20% | | | |

Forecast comparisons were done using the following $\chi^2$ test:

$$\chi^2 = \sum_{i=1}^{I} \frac{(n_i - N_i)^2}{N_i}$$

(31)

where $n_i$ is the estimated number of individuals choosing alternative i and $N_i$ the true (i.e. simulated) value. The result of this test was contrasted against the critical $\chi^2$ value for 95% confidence and I-1 (i.e. 3-1=2) degrees of freedom (5.99).

Table 9 presents the chi-square test values for the demand forecasts using the estimated models described in the previous sections. In line with econometric analysis on section 3.3, RC is not able to capture stochasticity variations in attributes because the generated stochasticity levels do not depend on the attribute values. For this reason, results show that there are no significant differences among the response results from the random coefficients models (i.e. RCg and RCs) and the MNL.

Results also suggest that the MNL models estimated with the larger samples (i.e. 2000 and 5000 observations) are fairly robust in terms of response, especially for policies that imply minor changes in the attributes. Notwithstanding, in all cases the more flexible models, such as SV and ECML, perform consistently better than the MNL, correctly forecasting the market shares when higher impact policies (i.e. implying major changes in some attributes) are considered. Even though the SV and ECML models estimated with a low sample size are also preferred against the MNL in response under high impact policies, all the models show poor forecasting behaviour in that case (the issue of adequate sample sizes to estimate even a simple model such as MNL was noted more than 30 years ago by Williams and Ortúzar, 1982). The MNL model seems to be superior to the SV and ECML models in forecasting only in the case of level of randomness II under policy P2 (which could just be an oddity).

These forecasting results could be seen as specific for the generated policy scenarios and databases, and cannot be assumed to be general enough for different combination of parameter values and policy scenarios, which in turn could consider different range of variations in the attributes and stochasticity levels. However, what we want to note is that results are highly sensible to sample size. A more detailed analysis about the generalization of our results should be undertaken and can be performed following the approach proposed by Hess and Train (2011), which consist on drawing multiple "versions" of the databases varying the sample size, the forecasting scenarios, the attribute values and the stochasticity levels.

Table 9 Estimated models response properties ($\chi^2$ test values)

| Level | Policy | 500 observations | | | | | 2000 observations | | | | | 5000 observations | | | | | Perfect | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | *MNL* | *RCg* | *RCs* | *SV* | *ECML* | *MNL* | *RCg* | *RCs* | *SV* | *ECML* | *MNL* | *RCg* | *RCs* | *SV* | *ECML* | *SV* | *ECML* |
| I | P0 | 1.4 | 2.0 | 2.9 | 2.5 | 2.3 | 0.3 | 0.3 | 0.1 | 0.1 | 0.1 | 0.1 | 0.1 | 0.1 | 0.1 | 0.1 | 0.1 | 0.1 |
| | P1 | 1.1 | 1.8 | 3.1 | 3.3 | 2.7 | 0.3 | 0.5 | 0.2 | 0.1 | 0.1 | 0.2 | 0.2 | 0.1 | 0.1 | 0.1 | 0.2 | 0.2 |
| | P2 | 1.5 | 2.3 | 4.2 | 5.2 | 4.2 | 0.5 | 0.7 | 0.3 | 0.2 | 0.2 | 0.1 | 0.2 | 0.1 | 0.2 | 0.1 | 0.2 | 0.2 |
| | P3 | 0.2 | 0.0 | 2.7 | 5.7 | 2.7 | 1.3 | 1.1 | 0.9 | 0.5 | 0.6 | 0.7 | 0.8 | 0.5 | 0.6 | 0.5 | 0.1 | 0.1 |
| | P4 | 1.9 | 2.6 | 3.2 | 4.3 | 3.3 | 0.2 | 0.2 | 0.1 | 0.1 | 0.1 | 0.0 | 0.0 | 0.0 | 0.2 | 0.1 | 0.0 | 0.0 |
| | P5 | 0.6 | 0.8 | 2.1 | 2.9 | 2.1 | 1.0 | 1.0 | 0.4 | 0.3 | 0.5 | 0.4 | 0.5 | 0.2 | 0.1 | 0.1 | 0.1 | 0.1 |
| | P6 | 12.0 | 9.8 | 24.2 | 5.8 | 7.6 | 4.4 | 3.5 | 6.9 | 2.5 | 3.0 | 1.4 | 1.0 | 2.8 | 0.5 | 0.7 | 0.5 | 0.5 |
| | P7 | 11.1 | 11.4 | 17.4 | 3.1 | 4.7 | 3.1 | 2.8 | 4.0 | 1.4 | 1.5 | 0.9 | 0.7 | 0.9 | 0.2 | 0.2 | 0.4 | 0.4 |
| II | P0 | 5.1 | 5.1 | 4.6 | 2.7 | 3.9 | 0.2 | 0.2 | 0.1 | 0.1 | 0.0 | 0.5 | 0.5 | 0.4 | 0.5 | 0.3 | 0.7 | 0.7 |
| | P1 | 6.3 | 7.5 | 6.6 | 4.1 | 5.5 | 1.0 | 1.5 | 1.1 | 0.5 | 0.7 | 1.2 | 1.3 | 1.0 | 0.6 | 0.5 | 0.7 | 0.7 |
| | P2 | 5.5 | 7.1 | 6.6 | 5.1 | 6.6 | 0.9 | 1.5 | 1.0 | 0.4 | 0.6 | 1.0 | 1.0 | 0.7 | 0.3 | 0.3 | 0.7 | 0.7 |
| | P3 | 7.1 | 4.6 | 2.8 | 1.2 | 0.6 | 7.2 | 6.2 | 5.6 | 4.9 | 4.8 | 8.1 | 7.6 | 6.8 | 5.8 | 5.1 | 3.4 | 3.4 |
| | P4 | 4.5 | 4.7 | 3.8 | 1.5 | 2.6 | 0.2 | 0.2 | 0.1 | 0.7 | 0.8 | 0.0 | 0.0 | 0.1 | 0.9 | 0.8 | 0.4 | 0.4 |
| | P5 | 10.1 | 8.7 | 6.2 | 3.6 | 4.8 | 5.2 | 4.8 | 4.2 | 3.1 | 3.9 | 5.7 | 5.3 | 4.4 | 2.8 | 3.0 | 2.0 | 2.1 |
| | P6 | 22.4 | 19.0 | 36.9 | 18.7 | 20.6 | 7.9 | 6.0 | 7.1 | 4.3 | 4.2 | 3.8 | 2.9 | 4.1 | 2.0 | 2.1 | 0.4 | 0.4 |
| | P7 | 20.0 | 19.1 | 26.5 | 8.3 | 10.8 | 7.2 | 6.1 | 5.5 | 3.3 | 4.6 | 3.7 | 3.3 | 3.3 | 1.8 | 2.8 | 0.9 | 0.9 |
| III | P0 | 2.0 | 2.4 | 2.0 | 2.0 | 1.2 | 0.2 | 0.2 | 0.2 | 0.3 | 0.4 | 0.4 | 0.3 | 0.5 | 0.5 | 0.3 | 0.9 | 0.9 |
| | P1 | 5.2 | 6.8 | 5.7 | 5.2 | 4.0 | 0.4 | 0.6 | 0.3 | 0.1 | 0.2 | 0.6 | 0.7 | 0.6 | 0.4 | 0.3 | 0.3 | 0.3 |
| | P2 | 5.8 | 7.8 | 6.7 | 5.8 | 5.5 | 0.4 | 0.7 | 0.5 | 0.4 | 0.3 | 0.2 | 0.4 | 0.3 | 0.1 | 0.1 | 0.3 | 0.3 |
| | P3 | 7.7 | 5.5 | 5.0 | 7.7 | 1.8 | 6.9 | 6.1 | 4.4 | 3.8 | 4.5 | 6.1 | 4.9 | 4.4 | 3.3 | 2.9 | 1.4 | 1.4 |
| | P4 | 3.2 | 3.4 | 2.5 | 3.2 | 1.4 | 0.2 | 0.1 | 0.1 | 0.4 | 0.3 | 0.8 | 0.7 | 0.2 | 0.1 | 0.0 | 0.0 | 0.1 |
| | P5 | 5.5 | 4.9 | 5.0 | 5.5 | 4.0 | 3.7 | 3.3 | 2.2 | 1.6 | 2.1 | 3.6 | 2.9 | 2.9 | 2.1 | 1.8 | 1.5 | 1.5 |
| | P6 | 21.8 | 18.5 | 19.5 | 21.8 | 13.2 | 7.9 | 4.8 | 5.3 | 3.9 | 3.4 | 3.5 | 2.1 | 1.6 | 0.9 | 1.2 | 0.1 | 0.1 |
| | P7 | 24.8 | 24.2 | 22.0 | 24.8 | 13.9 | 9.3 | 7.7 | 7.7 | 6.3 | 5.2 | 4.6 | 4.1 | 2.5 | 1.5 | 2.0 | 0.2 | 0.2 |

## CONCLUSIONS

We examine the inclusion of stochastic variables in the estimation of discrete choice random utility models. For this we conducted an econometric analysis of the problem and designed an experimental study using simulated data with variables subject to stochastic variations. MNL and different ML model specifications were estimated and their results statistically compared using several measures of fit (i.e. parameter recovery and willingness-to-pay measures) and also their response properties to several transport policies.

The econometric analysis shows that the problem can be approximated using ML formulations with a flexible covariance matrix that allow including heteroscedasticity among alternatives and/or observations, due to the presence of the stochastic variables. The most appropriate specification would depend on the data variation structure.

The simple MNL model cannot solve the problem properly. If the variable randomness does not depend on the attribute magnitude and its variance is constant across observations, the problem can be approximated through a stochastic variable (SV) model, which is a specific kind of ML and can be shown to be equivalent to an error component (ECML) model. The ECML structure has certain estimation benefits in term of parameter identification in comparison with the SV model; furthermore, the structure also allows treating the stochastic variable effect as a particular kind of heteroscedasticity among alternatives.

On the other hand, when stochastic variations are related with the magnitude of the explanatory variables, randomness can be captured more adequately through a random coefficients (RC) model. In this case, however, the presence of stochastic variables, with or without correlation, can be confounded with potential taste heterogeneity in the population.

In line with the econometric analysis, results using simulated data confirm that ML models are more appropriate than MNL models to recover the true population parameters (which are, of course, known in the simulated data) when stochastic variables are considered. Experiments with our simulated dataset suggest that MNL models seem to be fairly robust for computing marginal rates of substitution and for demand forecasting to non-dramatic policy changes. The best results when dealing with the stochastic variables problem were obtained using the ECML structure, but still care should be taken if the model is estimated with small sample size data.

Our findings could be useful to formulate more robust specifications to better explain behaviour and forecast choices when significant randomness in the explanatory variables is suspected. Further researches should be focused on generalizing our forecasting findings using multiple simulated databases varying stochasticy levels, attribute values, sample size, and policy scenarios. Also, it should be relevant to test more complex error structures, especially when stochastic variations are correlated and/or present some dependence among them. Another important issue confirmed here relates to the appropriate estimation of models using low sample sizes; following Sillano and Ortúzar (2005), maybe this should be performed using Bayesian estimation techniques.

## REFERENCES

Bhatta, B., Larsen O., 2011. Errors in variables in multinomial choice modelling: a simulation study applied to a multinomial model of travel mode choice. Transport Policy 18, 326–335.

Brey R., Walker J.L., 2011. Estimating time of day demand with errors in reported preferred times: an application to airline travel. Procedia Social and Behavioural Sciences 17, 150–168.

Bolduc D., Alvarez-Daziano R., 2009. On estimation of hybrid choice models. International Choice Modelling Conference, Harrogate.

Bolduc, D., Boucher N., Alvarez-Daziano R., 2008. Hybrid choice modelling of new technologies for car choice in Canada. Transportation Research Record 2082, 63-71.

Cantillo V., Heydecker B.G., Ortúzar J. de D., 2006. A discrete choice model incorporating thresholds for perception in attribute values. Transportation Research Part B 40, 807–825.

Cantillo, V, Amaya J., Ortúzar J. de D., 2010. Thresholds and indifference in stated choice surveys. Transportation Research Part B 44, 753-763.

Carrol, R.J., Ruppert D., Stefanki L.A., Criniceanu C.M., 2006. Measurement Error in Nonlinear Models. Chapman & Hall, London.

Carroll, R.J., Spiegelman, C.H., Gordon Lan, K.K., Bailey, K.T., Abbott, R.D., 1984. On errors-in-variables for binary regression models. Biometrika 71, 19–25.

Cherchi, E., Ortúzar, J. de D., 2008. Predicting best with mixed logit models: understanding some confounding effects. In: P.O. Inweldi (ed.), Transportation Research Trends, 215–235. Nova Science Publishers, Inc., New York.

Daly, A.J., Ortúzar, J. de D., 1990. Forecasting and data aggregation: theory and practice. Traffic Engineering and Control 31, 632–643.

Diaz, F., 2012. Estudio Econométrico y Experimental de Variables Estocásticas en Modelos de Elección Discreta. M.Sc. Thesis, Universidad del Norte (in Spanish).

Fuller, W.A., 1987. Measurement Error Models. John Wiley and Sons, New York.

Hess, S., Train, K.E., 2011. Recovery of inter- and intra-personal heterogeneity using mixed logit models. Transportation Research Part B 45(7), 973-990.

Jackson W.B., Jucker, J.V., 1982. An empirical study of travel time variability and travel choice behaviour. Transportation Science 16, 460–475.

Kao, C., Schnell, J., 1987. Errors in variables in the multinomial responde model. Economics Letters 25, 249–254.

McFadden, D., Train, K.E., 2000. Mixed MNL models for discrete response. Journal of Applied Econometrics 15, 447–470.

Ortúzar, J de D., Ivelic, A.M., 1986. Efectos de la desagregación temporal de variables de servicio en la especificación y estabilidad de funciones de demanda. In: Actas del IV Congreso Panamericano de Ingeniería de Tránsito y Transporte, Santiago (in Spanish).

Ortúzar, J. de D., Willumsen, L.G., 2011. Modelling Transport. 4th edition, John Wiley and Sons, Chichester.

Stefansky, L.A., Carroll, R.J., 1985. Covariate measurement error in logistic regression. The Anals of Statistics 13, 1335–1351.

Stefansky, L.A., Carroll, R.J., 1987. Trust conditional scores and optimal scores for generalized linear measurement error models. Biometrika 74, 703– 16.

Stefansky, L.A., Carroll, R.J., 1990. Structural logistic regression measurement model. Proceedings Summer Research Conference on Statistical Analysis of Measurement Error Models and Applications.  Providence, Rhode Island.

Swait, J.D., Bernardino, A., 2000. Distinguishing taste variation from error structure in discrete choice data. Transportation Research Part B 34, 1–15.

Steinmetz, S.S.C, Brownstone, D., 2005. Estimating commuters 'value of time' with noisy data: a multiple imputation approach. Transportation Research Part B 36, 865–889.

Train, K.E., 1978. The sensitivity of parameter estimates to data specification in mode choice models. Transportation 7, 301–309.

Train, K.E., 2009. Discrete Choice Methods with Simulation. 3rd edition, Cambridge University Press, Cambridge.

Walker J.L., 2001. Extended Discrete Choice Models: Integrated Framework, Flexible Error Structures and Latent Variables. Ph.D Thesis, MIT.

Walker J.L., Jieping, L., Sirinivasan, S., Bolduc, D., 2010. Travel demand models in the developing world: correcting for measurement errors. Transportation Letters: the International Journal of Transportation Research 4, 231-243.

Wansbeek T., Meijer, E., 2000. Measurement Error and Latent Variables in Econometrics. Elsevier, Amsterdam.

Williams, H.C.W.L., Ortúzar, J. de D., 1982. Behavioural theories of dispersion and the mis-specification of travel demand models. Transportation Research Part B 16, 167–219.

Yamamoto T., Komori, R., 2010. Mode choice analysis with imprecise location information. Transportation 37, 491–503.