



Modelling correlation patterns in mode choice models estimated on multiday travel data



Elisabetta Cherchi ^{a,*}, Cinzia Cirillo ^b, Juan de Dios Ortúzar ^c

^a *Transportation Operations Research Group (TORG), School of Civil Engineering and Geosciences, Newcastle University, Cassie Building, NE1 7RU, UK*

^b *Department of Civil and Environmental Engineering, University of Maryland, College Park, MD 20742, USA*

^c *Department of Transport Engineering and Logistics, Centre for Sustainable Urban Development (CEDEUS), Pontificia Universidad Católica de Chile, Casilla 306, Cod. 105, Santiago 22, Chile*

ARTICLE INFO

Article history:

Received 15 May 2016

Received in revised form 21 October 2016

Accepted 18 November 2016

Available online 6 January 2017

ABSTRACT

Understanding individual choices over time and measuring day-to-day variability in travel behaviour is important to capture the full range of travel behaviour. Although not very common, to date several multi-day travel surveys have been conducted and panel data is available to model different transport choices. However, determining the length of a panel that allows revealing variability in travel behaviour remains an open question. Also, no final agreement has been reached about modelling the various dimensions of correlation over the repeated observations. In this paper, we use the six-week panel data from the Mobidrive survey to estimate a mode choice model that accounts for correlation across individual observations over two time periods: all days of a single week and different days of the week (e.g. all Mondays) in the wave. We first analyse these effects separately, estimating different models for each type of correlation; then we try to disentangle the relative effects of each type of correlation, estimating both types jointly. We found that both types of correlation appeared highly significant when estimated alone, while only the correlation across a given day over the six-week period remained significant, when both types were estimated jointly. This implies that for the Mobidrive panel there is much less variability in mode choice across weeks than across the days of each week. It also suggests that one week could be an appropriate length for a panel to estimate modal choice and to correctly reveal day-to-day variability.

© 2016 Elsevier Ltd. All rights reserved.

1. Introduction

Transport planners and modellers have often questioned if a one-day data set is able to capture the full range of travel undertaken by an individual, as there are many activities that are not necessarily performed on a daily basis. In these days, they are still confronted with the problem of measuring variability (day-to-day, week-to-week) in travel behaviour. However, although several multi-day (or panel data¹) travel surveys and more advanced modelling techniques are now available, analyses based on these kind of data are still not a standard tool for travel behaviour modelling and demand forecasting.

* Corresponding author at: Transportation Operations Research Group (TORG), School of Civil Engineering and Geosciences, Newcastle University, Cassie Building, NE1 7RU, UK.

E-mail addresses: Elisabetta.Cherchi@newcastle.ac.uk (E. Cherchi), ccirillo@umd.edu (C. Cirillo), jos@ing.puc.cl (J.D. Ortúzar).

¹ In this paper we use the term “panel data” to indicate a typical revealed preference survey repeated several times (strictly speaking, repeated data from a stated preference survey could also be considered a panel).

Panel data can be classified into two categories: “long survey panels” and “short survey panels”. The first consists of repeating the same survey (i.e. with the same methodology and design) at “separate” times; for example, once or twice a year, for a certain number of years. Long survey panels have been extensively used in the past to gain insights into activity scheduling and planning, to study dynamic effects, and for their ability to identify temporal variation in travel behaviour (Hanson and Huff, 1988; Jones and Clarke, 1988; Pas, 1988; Pas and Sundar, 1995).

The “short survey panels” are, in general, more recent applications and consist of multi-day data where repeated measurements are collected on the same sample of units over a “continuous” period of time (e.g. seven or more successive days), but the survey is not necessarily repeated in subsequent years. In 1999 Axhausen and his group collected a six-week travel survey (called Mobidrive) for the German cities of Karlsruhe and Halle. A similar survey design, improved with the experience gained in Germany, was transferred in 2003 to Thurgau (Switzerland). These panels have been used to detect rhythms of daily life (Axhausen et al., 2002), to compare different indices that measure similarities of travel behaviour (Schlich and Axhausen, 2003), to draw evidence on the parametric assumptions behind the value of time distribution (Cirillo and Axhausen, 2006), to examine the length between successive participations in several activity purposes (Bhat et al., 2005), to explain current behaviour on the individuals’ history and experience (Cirillo and Axhausen, 2010), to study the intrinsic variability in individual preferences for mode choice, the effect of long period plans and habitual behaviour in daily mode choices (Cherchi and Cirillo, 2014), to account for the dynamic effect of inertia over time in the mode choice (Cherchi et al., 2013) and to study the effect of intra-individual variation in preferences in the estimation of equity effects of congestion charges (Börjesson et al., 2013).

On the other hand, data from Mobidrive and a three-day activity diary collected in Santiago in 2003, were used by Jara-Díaz et al. (2007) to estimate discrete-continuous mode choice and activity duration models. Furthermore, a two-day time-use diary, extracted from the 2002 US National Panel Study of Income Dynamics, was used by Copperman and Bhat (2007) to examine time-use in children’s activities. Using a seven-day diary from the Flemish Time-Use Surveys, Minnen et al. (2015) studied transport habits under different delineations of the time-use data and for different assumptions about regularity (i.e. varying its *tempo* and *timing*), while Minnen and Glorieux (2011) discuss the length of the time-use survey in order to capture the organization of individuals along cycles of work and leisure.

A Computerized Household Activity-Scheduling Survey (CHASE) was designed in Toronto (Canada) over a period of seven days (Doherty and Miller, 2000) and used to estimate a demand model system for daily activity programming (Khandker and Miller, 2007). Stopher et al. (2008) pioneered the use of Global Positioning Systems (GPS) as a means of measuring personal travel; they used a 28-day GPS survey of 50 households to analyse the variability in daily travel of individuals and the proportions of variance due to intrapersonal and interpersonal variability. Finally, Vanhulsel et al. (2007) introduce an extended reinforcement learning approach to produce weekly activity patterns in Belgium.

The *Santiago Panel* (2006–2010) combines both “short” and “long” survey panel approaches (Yañez et al., 2010), as it is a five-day pseudo diary that has been repeated five different times, just before and four times after the implementation of the radically new Santiago’s public transport system (Muñoz et al., 2009). Data from this panel have been used to study the effect of shock and inertia in individual behaviour and to examine the effect of repeated observations (Yañez et al., 2008) and, more recently, to study the evolution and persistence of modality styles and travel mode choice behaviour in a dynamic context (Vij, 2013).

In a developing country context, the National Planning and Development Agency of the Republic of Indonesia and the Japan International Cooperation Agency collected SITRAMP 2004. It contains an activity diary survey for respondents from households within the Jakarta Metropolitan Area for two consecutive weekdays, Thursday and Friday, and two consecutive weekend days, Saturday and Sunday; this dataset was used to study day-to-day variability in travellers’ activity-travel patterns (Dharmowijoyo et al., 2016). A multi-dimensional three-week household time-use and activity diary was collected in the Bandung Metropolitan Area, Indonesia. Preliminary analyses have been conducted to examine the use of motorised modes, activity participation and multi-tasking, and the relation between transport choices and physical health (Dharmowijoyo et al., 2016).

Although this literature review is not exhaustive by any means, it clearly shows that the use of panel data, especially short survey panels, has recently increased. Nevertheless, the applications mainly refer to time-use, activity choice and activity duration, while estimation of mode choice models with panel data has received less attention. Another topical element in the recent literature is the importance of day-to-day variability and the correct length of the panel to allow revealing or avoiding such variability. A number of questions related to modal choice are still open. Do individual choices differ among days of the same week? Do they differ among weeks? Do we really need to have information on several weeks to capture individual behaviour dynamics?

There are several ways these problems can be addressed. Properly accounting for the various dimensions of correlation over the repeated observations provided by a given individual is crucial to understand the long-term structure of individual choices and should give insights into the proper length of the survey. It should also allow to correctly estimating models and properly use them in prediction. In this context, the use of short survey panels appears fundamental as several dimensions of correlation across responses can be studied, i.e. over trips made the same day, the same week, a given day-of-week in a longer period, and/or over individuals and households.

In this paper, we use the six-week panel data from the Mobidrive survey to estimate a mode choice model that accounts for correlation across individuals over two time periods: a single week and a day of the week (i.e. all Mondays) in the whole panel. We first analyse these effects separately, estimating different models for each type of correlation. However, in that

case it is not possible to ascertain their relative effects over different time periods. To overcome this problem, we then made the important step forward of trying to disentangle the relative effects of each type of correlation. A flexible model that jointly accounts for both dimensions of correlation was estimated, allowing to directly compare their relative effects.

The rest of the paper is organized as follows. In Section 2 we discuss the issue of modelling with panel data and how to account for correlation across the observations provided by a given individual. We first discuss the general theory and then how we extend it to disentangle the relative effects of the two levels of correlation considered in the paper: (1) over a single week and (2) over days of the week. Section 3 provides a brief description of the dataset used and a descriptive analysis of some characteristics of the panel data sample. Section 4 reports our modelling estimation results, and Section 5 summarizes our conclusions.

2. Mixed Logit model on panel data: accounting for correlation across individuals

The Mixed Logit (ML) is any model the choice probability of which can be expressed as an integral of standard logit probabilities, evaluated at parameters μ , over a density of parameters (Train, 2009). As well known, the vector μ can assume any desired distribution. The ML model, in fact, allows for a very rich decomposition of the vector of unobserved components (Greene and Hensher, 2007); in particular, it allows for correlation within individuals, which is the typical effect that occurs in panel data where each individual provides more than one piece of information. In these cases, the utility (U_{qjt}) that individual q associates with alternative j on choice situation t ; has the following structure²:

$$U_{qjt} = (b_k + v_{qk}\sigma_k)X_{kqjt} + [\dots] + \varepsilon_{qjt} \tag{1}$$

where X_{kqjt} is the k_{th} attribute in choice situation t (X might also be an alternative specific constant³) the parameter of which shows random heterogeneity; $\mu_{qk} = v_{qk}\sigma_k$; b_k and σ_k are the mean and standard deviation of the k_{th} random parameter, while v_{qk} is a random variable that varies among individuals but is fixed across the choices of the sequence \mathbf{j} , and accounts for correlation across the T choices of each individual. For simplicity we only show one random parameter but, as well known, Eq. (1) can be extended to account for almost any number of random parameters as well as for correlation among them. Finally ε_{qjt} is the typical additive extreme value type 1 (EV1) error, which generates the standard logit probability.

Following Train (2009), if $\mathbf{j} = \{j_1, \dots, j_t, \dots, j_T\}$ is the sequence of choices made by each individual q , the probability (P_{qj}) of person q making this sequence of choices is the product of the conditional logit formulae, integrated over the density of the parameters; as well known, the logit formula itself is the exact integral of the density function of iid EV1 variables. In particular, from Eq. (1) the density function of ε_{qjt} is conditional upon the realization of v_{qk} ; hence for a given person q , the probability of making the sequence of choices $\mathbf{j} = \{j_1, \dots, j_t, \dots, j_T\}$ can be written as:

$$P_{qj} = \int_{R_v} \int_{R_\varepsilon} f(\boldsymbol{\varepsilon}|v_q)f(v_q)dv_qd\boldsymbol{\varepsilon} \tag{2}$$

where the domains of integration are:

$$R_v = [-\infty, +\infty]; \quad R_\varepsilon = \begin{cases} \varepsilon_{qjt} \leq \varepsilon_{qjt} + (b_k + v_{qk}\sigma_k)(X_{kqjt} - X_{kqjt}) + [\dots] & \forall j \in I_q \\ U_{qjt} \geq 0 \end{cases} \tag{3}$$

Finally, as $f(\boldsymbol{\varepsilon}|v_q) = f(\varepsilon_{qj_1}|v_q) \dots f(\varepsilon_{qj_t}|v_q) \dots f(\varepsilon_{qj_T}|v_q)$ the probability (P_{qj}) takes the typical expression:

$$P_{qj} = \int_{R_v} \prod_{t=1}^T \frac{e^{v_{qjt}(v_q)}}{\sum_{j \in I_q} e^{v_{qjt}(v_q)}} f(v_q)dv_q \tag{4}$$

The six-week panel data from the Mobidrive survey used to estimate our models include several dimensions of potential correlation across (1) tour mode choices made during the same day, (2) the same week, or (3) by the same individual. In this paper we focus on the correlation over a given day in the whole panel and over all days in a single week, with the specific aimed at disentangling their relative effects.

Figs. 1 and 2 show schematically all the information available for each individual q and the two different ways in which we chose to investigate correlation. In particular, Fig. 1 shows correlation over a given day throughout the weeks, while Fig. 2 shows correlation inside each week. Of course, correlation over all observations for each individual in a wave includes all tours reported by that individual.

In our case, the sequence of choices made by each individual q can be written as $\mathbf{j} = \{j_{Monday,w1}, \dots, j_{d,w}, \dots, j_{Friday,w6}\}$. If only correlation over each week was considered, the choice sequence for each week $w_i = \{w_1, \dots, w_6\}$ would be $\mathbf{j}_{w_i} = \{j_{Monday,w_i}, \dots, j_{d,w_i}, \dots, j_{Friday,w_i}\}$. Analogously, the sequence of choices for each day (e.g. $d = Monday, Tuesday$, and so

² In Eq. (1) we make the assumption that the random term is distributed Normal (0,1). This assumption does not impose any restriction on the methodology proposed and several other distributions can be used, from simple distributions such as the triangular to more flexible, but also more complex, forms such as the S_B distribution (Train and Sonnier, 2005). Some important advantages of using the Normal distribution are discussed by Sillano and Ortúzar (2005).

³ Although both versions of the ML (random parameter and error components) can be used to account for differences between long and short-term variations, the error components specification carries some more risks in complex structures due to theoretical identification problems (Walker, 2002).

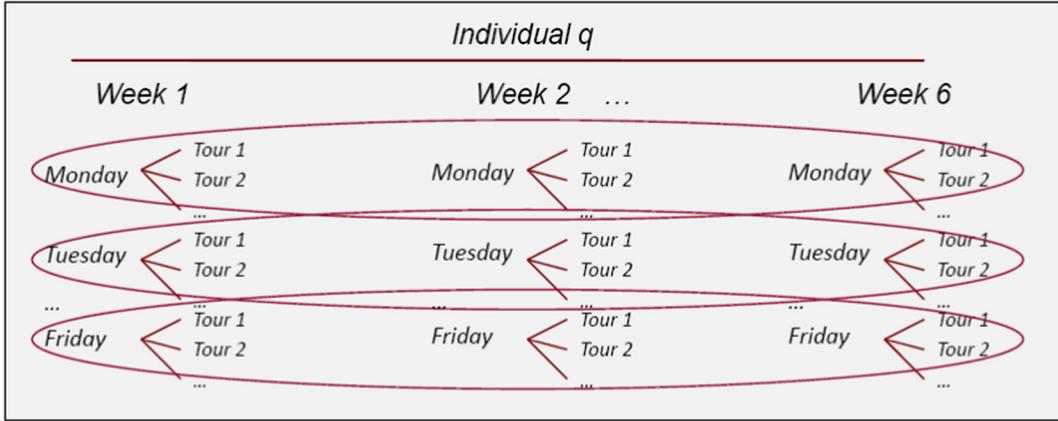


Fig. 1. Correlation over given days of the panel for each individual q .

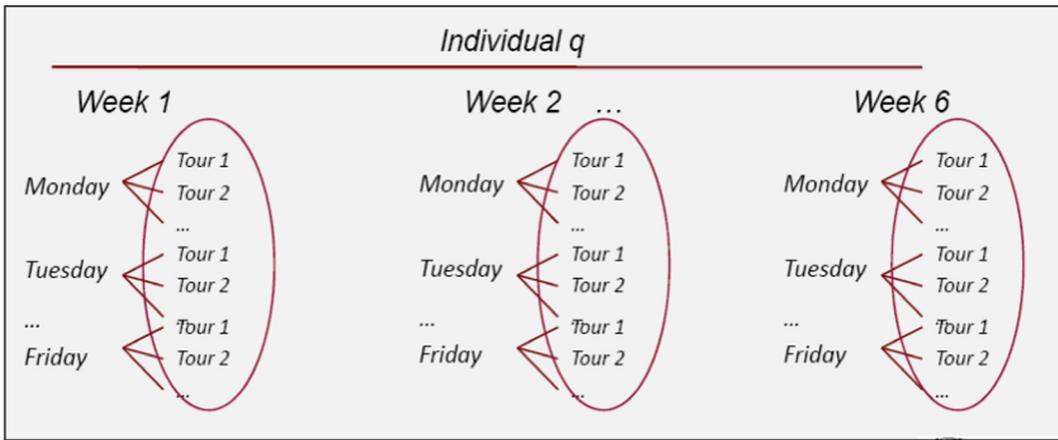


Fig. 2. Correlation over each week for each individual q .

on) over the six weeks would be $\mathbf{j}_d = \{j_{d,w1}, \dots, j_{d,w2}, \dots, j_{d,w6}\}$. It is important to note that correlation over days and over weeks can be thought of as “subsets”⁴ of the correlation over individuals. In fact, as illustrated in Fig. 1, the sequence of choices over days is a subset of the sequence for each individual ($\mathbf{j}_d \subset \mathbf{j}$). Analogously, as illustrated in Fig. 2, the sequence of choices over each week is a subset of the sequence for each individual ($\mathbf{j}_{w_i} \subset \mathbf{j}$).

When only one type of correlation is accounted for (i.e. only the correlation for each day over the six weeks) the ML model in Eq. (4) is the standard formulation for repeated observations (Train, 2009). Instead, when both types of correlation are estimated jointly, new assumptions need to be made because the sequences of choices over days and over each week represent two overlapping but disjoint sets ($\mathbf{j}_d \not\subset \mathbf{j}_{w_i}$, and $\mathbf{j}_{w_i} \not\subset \mathbf{j}_d$).

To jointly estimate these two correlation effects we assumed that attribute tastes were distributed across days and weeks with the same mean (b_k) but with different standard deviations across days (σ_{k1}) and across weeks (σ_{k2}). The specification of the random terms in Eq. (1) then can be rewritten as follows:

$$U_{qj_{d,w}} = (b_k + v_{qk,d}\sigma_{k1} + v_{qk,w}\sigma_{k2})X_{kqj_{d,w}} + [\dots] + \varepsilon_{qj_{d,w}} \tag{5}$$

where $U_{qj_{d,w}}$ is the utility that individual q associates to alternative j on the day d of week w ; $X_{kqj_{d,w}}$ is the k_{th} attribute with a parameter showing random heterogeneity; $v_{qk,d}$ is a random variable distributed $N(0,1)$ that varies among individuals and across days but is fixed for the same day across the weeks of the panel (i.e. all Mondays and so on). Analogously v_{qk,w_i} distributes $N(0,1)$ and varies among individuals and weeks but is fixed for all days inside each week.

For simplicity, Eq. (5) only reports the specification of the random terms. The final specification adopted to estimate the ML model with our panel data includes also several socio-economic variables (SE) with linear and parameterized effects to account for systematic heterogeneity around the preferences for specific alternatives and around the responses to level-of-

⁴ This sentence does not imply that the models are nested from a modelling perspective and should not be confused with that.

service attributes (Ortúzar and Willumsen, 2011, p. 279). Moreover, Eq. (5) refers to only one random parameter but, as we discuss in the results, random heterogeneity was tested for both travel time and cost. Hence in our case $k = 1, 2$. We also tested for correlation between the time and cost parameters but it came out as insignificant.

Finally, in our specification the random term $\varepsilon_{qj,d,w}$ has, as usual, the typical additive extreme value type 1 (EV1) error, with conditional density function equal to:

$$f(\boldsymbol{\varepsilon} | v_{qk,d}, v_{qk,w}) = \prod_{d=Mon,\dots,Fri} \prod_{w_i=w_1,\dots,w_6} f(\varepsilon_{qj,d,w} | v_{qk,d}, v_{qk,w_i}) \quad (6)$$

Then, the sequence of choices made by individual q is $\mathbf{j} = \{j_{Monday,w1}, \dots, j_{d,w}, \dots, j_{Friday,w6}\}$, and the probability (P_{qj}) of person q making this choice sequence is:

$$P_{qj} = \int_{R^+} \prod_{d=Mon,\dots,Fri} \prod_{w_i=w_1,\dots,w_6} \frac{e^{V_{ej_k}(v_{k,d}, v_{k,w})}}{\sum_{j \in I_q} e^{V_{ej_k}(v_{k,d}, v_{k,w})}} f(v_{k,d}) f(v_{k,w}) dv_{k,d} dv_{k,w} \quad (7)$$

The vector of unknown parameters can be estimated by maximizing the simulated log-likelihood function, i.e. by solving⁵:

$$\max_{b_k, \sigma_{k1}, \sigma_{k2}} SLL(b_k, \sigma_{k1}, \sigma_{k2}) = \max_{b_k, \sigma_{k1}, \sigma_{k2}} \sum_{q=1}^Q \ln \frac{1}{R} \sum_{r=1}^R P_{qj,d,w}^r(b_k, \sigma_{k1}, \sigma_{k2}) \quad (8)$$

3. Dataset analysis

The panel dataset used in this paper was derived from a six-week travel diary held in Karlsruhe (Germany) in 1999, part of the *Mobidrive* survey, which involved 160 households and 360 individuals in the main survey. Details about the data and their collection techniques can be found in Axhausen et al. (2002) and in Fell et al. (2000)⁶. The recorded days were structured according to activity chains on a daily basis and all trips were grouped into tours. Each tour had associated a main activity (i.e. work/education) for working days and a principal activity, which was defined as the longest duration out-of-home activity, for non-working days. All daily activity chains were represented in relation to these pivotal activities; the sequence of chains for a given day (week/individual) was called daily (weekly/individual) schedule. A detailed description of the framework applied to the original *Mobidrive* survey to define tours can be found in Cirillo and Axhausen (2006).

As in Cherchi and Cirillo (2014), other than some further tests on the availability of alternatives, we concentrated only on week days, so our final sample consists of 4089 single tours, 2488 daily schedules, 674 weekly schedules, 129 individual schedules and 56 household schedules. Saturdays and Sundays were also excluded in line with the work of Schlich and Axhausen (2003) where weekend days were treated separately.

Other than the typical level-of-service attributes, several *SE* and other variables were available from the survey and, as we will show in the next section, resulted significant in model estimation. In particular, it is important to mention that we tested not only the linear effects of single *SE* attributes (such as age and so on) but also analysed some categories formed by combining *SE* characteristics (i.e. allowing for systematic taste variations, Ortúzar and Willumsen, 2011). Hereafter, we report a brief description of the percentage of these categories into our sample. In particular, the sample was composed mainly of non-workers (75%); however, among the workers, 23% worked full time and 20% of the part-time workers were females. Tours were equally distributed among work/study (35%), shop/leisure (35%) and other activities (30%), and almost all ages were represented, although the majority (one third of the sample) lies in the range 51–65 years, while youngsters account only for 10% (exactly 5% between 18 and 25, and another 5% between 25 and 35). On the other hand, 27% of the sample was married with children under 18 years old; 45% declared to be the main car user and, in line with this number, 39% of the tours in the sample were made as car driver.

Although public transport is used only for 15% of the tours in our sample, 33% of the individuals declared to have a seasonal ticket. Even more interesting is to note that 70% of the tours are single tours (i.e. without additional stops between the main origin and destination); 16% have one stop and only 14% have more than one stop. Even car drivers (notwithstanding the flexibility offered by the car) mainly do single tours (64%), although, as expected, they tend to do a higher number of stops (even four and five stops) than car passengers, pedestrians and public transport users.

Finally, it is important to note that our sample is evenly distributed over days of the week and over the six weeks of the panel. The characteristics of the tours are quite similar over the days and over the weeks, with few exceptions. As mentioned before, *Car Driver* is by far the most popular travel mode and its use is stable over the days of the week and over the six weeks of the panel.

⁵ A specific code was written in Gauss to estimate this model. We used Paul Ruud's maximization routine and the Halton sequences made available by Prof. Kenneth Train in his web-page.

⁶ See also <http://www.ivt.ethz.ch/vpl/research/mobidrive> for a list of papers employing the *Mobidrive* data.

4. Modelling results

Using the dataset described in the previous section, several models were estimated accounting, as much as possible, for systematic taste variations. Then, random heterogeneity in tastes for travel time and cost was also estimated and correlation among individuals over different time periods explicitly analysed. Table 1 reports the results of several ML models estimated with different assumptions about correlation over observations. A simpler MNL allowing for systematic taste variations (parameterized effects) is also shown for comparison. Firstly, as expected, accounting for random heterogeneity, whatever the type of correlation considered, significantly improves model estimation. The MNL model estimated with an equivalent utility specification to the ML is inferior to any of the estimated ML models (the LR test allows rejecting the MNL at the 99.99% significance level).

Before analysing the effect of correlation, it is interesting to note that the mean values of all parameters do not vary among the different correlation specifications. All parameters are significant (t -test rejected at more than 95% significance) with very few exceptions. We also found that there are seven significant systematic variations around the travel time variable and 16 SE attributes explaining systematic alternative specific preferences. In particular, we found that the marginal utility of travel time is lower for individuals who are married with children, for females working part time, for work trips and for bikers between 35 and 65 years old. Conversely, the marginal utility of travel time increases (is smaller in absolute value) with the number of stops made during the tour (maybe because activities are performed at each stop), for those individuals living in suburban locations and using *Public Transport* (the longer the trips the smaller the disutility of an extra minute travelling), and for educational trips by *Public Transport*, *Walk* or *Bike* (students care less about travel time). Interactions with the cost attribute were also tested but none was found significant.

As for the preferences for each alternative, it is not surprising that *Car Driver* is preferred by those who are mainly car-users and travel many miles. It is also not surprising that *Public Transport* is preferred by season ticket owners, while it is not liked for leisure trips, by older people or by those who have children. It is worth noting that variables such as miles travelled and possession of a season ticket are usual indicators of habit (Gärling and Axhausen, 2003).

It is also interesting to remark that the low level of significance of the mean value of the alternative specific constants (ASC) in all alternatives (except the *Bike*) is due precisely to the inclusion of systematic heterogeneity in preferences for specific alternatives; in the models estimated without the SE attributes the ASC were all extremely significant. As for the *Bike* ASC, we believe that preference for biking depends more on individual attitudes towards this specific mode than on the specific characteristic of the individual or mode, and so it deserves a different analysis that cannot be carried out with the data we had available.

Regarding the effect of correlation over individuals in panel data, as expected, models improved significantly as correlation was extended from the single day, to the week and to all tours performed in the six weeks (these results are not reported in Table 1). In a short survey panel individual correlation might be due to two different effects: correlation across the same day over the entire period (i.e. six weeks here) and correlation over different days of a given week. Separate analyses were performed to try and disentangle these two effects.

In particular, model ML1 (Panel-week) accounts for correlation across days of the same week but assumes that observations across weeks are independent. The ML2 model (Panel day-of-the-week) accounts for correlation across the same day in different weeks (i.e. across all tours made by the same individual on the six Mondays, the six Tuesdays, and so on) while observations belonging to different days of the same week are considered independent. As reported in Table 1, both models show significant random heterogeneity for the travel time and cost parameters and although not directly comparable, they only differ in the way correlation is accounted for. Thus, comparing their log-likelihoods at convergence suggests that, in our six-week panel data there is more correlation across a given day over the six weeks (ML1) than across different days of the same week (ML2). But we cannot tell if both types of correlation do actually exist or if only one type is really important.

The specification adopted in Eq. (5) helps clarifying this issue. In fact, when both correlations (across days of the same week and across a given day-of-the-week) are estimated jointly, as in model ML3 (Panel week&day), results show that correlation across days of the same week loses significance for both the parameters of travel time and cost.

Finally, it is important to mention that none of our models presented empirical identification problems. Different from the theoretical identification (Walker, 2002), which is inherent to the model specification regardless of the data at hand, the empirical identification depends only on the information used to estimate the model and can be revealed only during the estimation process. Following Chiou and Walker (2007), we checked if our results remained stable when increasing the number of draws in the estimation; this was an expected result as we are using panel data. As reported by Cherchi and Ortúzar (2008), empirical identification problems reduce considerably as the number of observations available for each individual increases. Notwithstanding, although we did not find empirical identification problems, we found that the correlation across days of the same week was highly sensitive to the seed used for the random draws. Instead, the correlation across a given day-of-the-week was always highly significant. This issue deserves further investigation.

5. Conclusions

The use of panel data has recently increased in the area of travel behaviour research. Nevertheless, applications mainly refer to time-use and activity choice and duration, while use of panels in mode choice modelling to better capture, for exam-

Table 1

Model results: Correlation over days and weeks in panel data.

Variable	Alts.	MNL		ML1 Panel-week		ML2 Panel days of week		ML3 Panel-week&day	
		Estimate	t-test	Estimate	t-test	Estimate	t-test	Estimate	t-test
ASC_Car Pass	CP	0.1268	0.51	0.3685	0.74	0.4032	0.84	0.3934	0.80
ASC_PT	PT	-0.8676	-2.80	-0.7840	-0.95	-0.5368	-0.72	-0.5780	-0.76
ASC_Walk	Walk	-0.0306	-0.12	0.7118	1.10	0.7881	1.31	0.7573	1.23
ASC_Bike	Bike	-1.5992	-6.26	-1.5566	-2.30	-1.4557	-2.23	-1.4782	-2.22
Time (mean)	All	-0.0196	-5.20	-0.0637	-5.52	-0.0698	-6.52	-0.0707	-6.54
Time (s.d) Corr. day-weeks	All					-0.0779	-9.50	-0.0732	-5.56
Time (s.d) Corr. weeks	All			-0.0711	-8.69			-0.0055	-0.55
Cost (mean)	All	-0.0866	-5.09	-0.0870	-2.30	-0.1092	-2.57	-0.1127	-2.81
Cost (s.d) Corr. day-weeks	All					0.1197	3.86	0.1027	2.66
Cost (s.d) Corr. weeks	All			0.0946	3.08			0.0669	1.28
Time*Marchild	All	-0.0199	-4.42	-0.0292	-1.63	-0.0316	-0.72	-0.0308	-0.54
Time*Work	All	-0.0218	-0.26	-0.0205	-0.19	-0.0131	-0.96	-0.0135	-0.02
Time*FemPartTime	All	-0.0169	-0.05	0.0013	0.09	0.0001	0.01	-0.0001	-0.01
Time*Nstop	All	0.0056	3.18	0.0137	3.17	0.0150	3.21	0.0152	3.42
Time*Education	All	0.0142	3.69	0.0174	1.69	0.0190	1.98	0.0191	2.00
Time*Sub-urb. loc.	PT	0.0178	4.93	0.0341	3.05	0.0337	3.32	0.0340	3.37
Time*age65	B	0.9246	4.17	1.3667	2.43	1.4067	2.53	1.4112	2.54
Age 26–35	CD	2.2000	4.57	2.6764	2.87	2.6000	2.79	2.5697	2.76
Age 51–65	W	0.7455	5.03	0.9512	1.48	1.0772	2.37	1.0750	2.32
Age 18–25	B	-0.0255	-0.03	-0.0335	-0.09	-0.0353	-0.17	-0.0354	-0.19
Age 26–35	B	1.0874	5.23	1.3184	1.86	1.3150	1.85	1.3104	1.84
Age 51–65	B	1.8072	5.17	2.5028	2.89	2.4002	2.69	2.3832	2.73
Main car user	CD	1.4406	11.8	1.6704	4.27	1.7040	4.65	1.6987	4.55
Annual Mileage	CD	0.0178	5.04	0.0149	1.51	0.0159	1.65	0.0158	1.64
Urban location	PT	0.8664	4.24	1.2712	2.26	1.2766	2.50	1.3012	2.56
Age 51–65	PT	-0.4435	-0.88	-0.1309	-0.31	-0.1556	-0.35	-0.1615	-0.37
Marchild	PT	-0.3499	-0.98	-0.6296	-0.55	-0.6747	-0.57	-0.6859	-0.61
SeasonTicket	PT	2.0320	11.1	1.6597	3.36	1.5764	3.24	1.5821	3.21
Leisure	PT	-0.7429	-0.31	-0.9224	-0.56	-0.8895	-0.57	-0.8611	-0.52
Work	PT	1.3938	7.57	1.5670	2.62	1.4906	2.65	1.5328	2.74
Time Budget	CD	0.0010	3.65	0.0012	2.55	0.0013	2.80	0.0012	2.73
Leisure	CP	1.0798	5.91	1.0951	3.07	1.0842	3.10	1.0893	3.11
Time Budget	B	0.0015	7.10	0.0018	3.28	0.0018	3.22	0.0018	3.19
$l(\theta)$		-765.35		-519.68		-502.78		-501.07	
ρ^2 (C)		0.2055		0.2760		0.2809		0.2814	

CD = Car Driver; CP = Car Passenger; PT = Public Transport; W = Walk; B = Bike.

Time [min]; Cost [DM] including any parking fees.

Age 18–25/age 26–35/age 51–65/age > 65: (Dummy).

FemPartTime = 1, if female and employed part-time.

Marchild = 1, if married with children under 18.

Sub-urb. loc. = 1 if the origin of the tour is in a suburban location.

Nstop is the number of secondary activities observed within each tour (Discrete).

M_CAR_U = 1 if main car user.

P_A_MVT is the Total annual mileage by car (Continuous).

SeasonTicket = Number of season tickets (Discrete).

Education = 1, if the purpose of the main activity of the tour is study.

Leisure = 1, if the purpose of the main activity of the tour is leisure.

Work = 1, if the purpose of the main activity of the tour is work.

Time budget [min] is the time spent on previous activities (home stay included) and previous travel.

ple, individual heterogeneity has received less attention. In particular, short survey panels are characterised by repeated measurements on the same sample of units gathered over a “continuous” period of time. This raises interesting issues, not yet explored, related with the several dimensions of correlation that might appear across responses provided by the same individual: over trips made in the same day, the same week, a given day-of-week, and/or over individuals and households.

In this paper, we specifically addressed the issue of the several dimensions of correlation in a short (six-week) survey panel. We estimated a random parameters Mixed Logit model accounting for correlation across individuals over two time periods: a single week and a given day (e.g. Mondays) in the whole panel. We first analysed these effects separately, and then we tried to disentangle the relative effects of each type of correlation. The models also account for several types of systematic heterogeneity over individual preferences.

As expected, accounting for correlation over individuals in panel data improves enormously the model results. Not surprisingly, results improved significantly as we extended correlation from a single day, to a week, to all tours performed in the six-week period. However, more interestingly, we found that in our six-day panel data, there is more correlation across the

same day over the six weeks than across different days of the same week. In particular, when both types of correlation were estimated jointly, the effect of correlation across different days of the same week disappeared and only correlation across the same day over the six weeks remained significant. When different models were estimated for each type of correlation, it was not possible to draw conclusions about the relative effect of correlation over different time periods; in fact, it was not possible to ascertain whether only one type of correlation is indeed important or if both types of correlation do actually exist.

The specification adopted helps clarifying the true effect that correlation can have in short panel data. This is crucial to understand the long-term structure of individual choices, to correctly estimate models and to properly use them in prediction. Moreover, as we found that in our sample there is no variability across tours made on the same day over six weeks, this suggests that one week could be an appropriate length for estimating mode choice models and to correctly reveal such variability. Having more weeks (i.e. a survey longer than one week) could be more important in terms of increasing the total number of observations and the number of observations for each individual (better statistical results).

Acknowledgements

The authors would like to thank Kay Axhausen for providing the *Mobidrive* data set, Luis I. Rizzi for an earlier discussion on the correlation issues and Massimiliano Bez for his valuable help on the model specification. Finally, the third author is indebted to the *Institute on Complex Engineering Systems* (ICM: P-05-004-F; CONICYT: FB0816), the *Centre for Sustainable Urban Development*, CEDEUS (CONICYT/FONDAP/15110020) and the *Bus Rapid Transit Centre of Excellence* funded by VREF (www.brt.cl), for their support.

References

- Axhausen, K.W., Zimmermann, A., Schönfelder, S., Rindsfuser, G., Haupt, T., 2002. Observing the rhythms of daily life: a six-week travel diary. *Transportation* 29, 95–124.
- Börjesson, M., Cherchi, E., Bierlaire, M., 2013. Within individual variation in preferences: equity effects of congestion charges. *Transp. Res. Rec.* 2382, 92–101.
- Cherchi, E., Ortúzar, J. de D., 2008. Empirical identification in the mixed logit model: analysing the effect of data richness. *Network Spatial Econ.* 8, 109–124.
- Cherchi, E., Börjesson, M., Bierlaire, M., 2013. A hybrid mode choice model to account for the dynamic effect of inertia over time. In: 3rd International Choice Modelling Conference, Sydney, Australia.
- Cherchi, E., Cirillo, C., 2014. Understanding variability, habit and the effect of long period activity plan in modal choices: a day to day, week to week analysis on panel data. *Transportation* 41, 1245–1262.
- Chiou, L., Walker, J.L., 2007. Masking identification in the discrete choice model under simulation methods. *J. Econom.* 141, 683–703.
- Cirillo, C., Axhausen, K.W., 2006. Evidence on the distribution of values of travel time savings from a six-week travel diary. *Transp. Res.* 40A, 444–457.
- Cirillo, C., Axhausen, K.W., 2010. Dynamic model of activity type choice and scheduling. *Transportation* 37, 5–38.
- Copperman, R.B., Bhat, C.R., 2007. An exploratory analysis of children's daily time-use and activity patterns using the child development supplement (CDS) to the US Panel Study of Income Dynamics (PSID). *Transp. Res. Rec.* 2021, 36–44.
- Dharmowijoyo, D.B.E., Susilo, Y.O., Karlström, A., 2015. Collecting a multi-dimensional three-weeks household time-use and activity diary in the Bandung Metropolitan Area, Indonesia. *Transp. Res.* 80A, 231–246.
- Dharmowijoyo, D.B.E., Susilo, Y.O., Karlström, A., 2016. Day to day variability in travellers' activity travel patterns in the Jakarta Metropolitan Area. *Transportation* 43 (4), 601–621.
- Doherty, S.T., Miller, E.J., 2000. A computerized household activity scheduling survey. *Transportation* 21, 75–97.
- Gärling, T., Axhausen, K.W., 2003. Introduction: habitual travel choice. *Transportation* 30, 1–11.
- Greene, W.H., Hensher, D.A., 2007. Heteroskedastic control for random coefficients and error components in mixed logit. *Transp. Res.* 43E, 610–623.
- Hanson, S., Huff, J.O., 1988. Systematic variability in repetitive travel. *Transportation* 15, 111–135.
- Jara-Díaz, S., Munizaga, M., Greeven, P., Guerra, R., 2007. The unified expanded goods-activities-travel model: theory and results. In: 11th World Conference on Transport Research, Berkeley, California.
- Jones, P., Clarke, M., 1988. The significance and measurement of variability in travel behaviour. *Transportation* 15 (1), 65–87.
- Khandker, M.N.H., Miller, E.J., 2007. Modelling activity program generation considering within day and day-to-day dynamics in activity travel behaviour. In: 86th Conference of the Transportation Research Board, Washington D.C..
- Minnen, J., Glorieux, I., van Tienen, T.P., 2015. Transportation habits: evidence from time diary data. *Transp. Res.* 76A, 25–37.
- Minnen, J., Glorieux, I., 2011. Two Days a week? A comparison of the quality of time-use data from 2-day, 7-day diaries and a weekly work grid. In: Carrasco, J.A., Jara-Díaz, S., Munizaga, M. (Eds.), *Time Use Observatory*. Grafica LOM, Santiago, pp. 105–118.
- Muñoz, J.C., Ortúzar, J. de D., Gschwendler, A., 2009. Transantiago: the fall and rise of a radical public transport intervention. In: Saaleh, W., Sammer, G. (Eds.), *Travel Demand Management and Road User Pricing: Success, Failure and Feasibility*. Ashgate, Farnham, pp. 151–172.
- Ortúzar, J. de D., Willumsen, L.G., 2011. *Modelling Transport*. John Wiley and Sons, Chichester.
- Pas, E.I., 1988. Weekly travel-activity behaviour. *Transportation* 15, 89–109.
- Pas, E.I., Sundar, S., 1995. Intra-personal variability in daily urban travel behaviour: some additional evidence. *Transportation* 22, 135–150.
- Fell, B., Schönfelder, S., Axhausen, K.W., 2000. *Mobidrive questionnaires*. Arbeitsberichte Verkehrs- und Raumplanung 52, Institut für Verkehrsplanung, Transporttechnik, Strassen- und Eisenbahnbau, ETH Zürich.
- Schlich, R., Axhausen, K.W., 2003. Habitual travel behaviour: evidence from a six-week travel diary. *Transportation* 30, 113–136.
- Sillano, M., Ortúzar, J. de D., 2005. Willingness-to-pay estimation with mixed logit models: some new evidence. *Environ. Plann.* 37A, 525–550.
- Stopher, P.R., Clifford, E., Montes, M., 2008. Variability of travel over multiple days. *Transp. Res. Board* 2054, 56–63.
- Train, K.E., 2009. *Discrete Choice Methods with Simulation*. Cambridge University Press, Cambridge.
- Train, K., Sonnier, G., 2005. Mixed logit with bounded distributions of correlated partworths. In: Alberini, A., Scarpa, R. (Eds.), *Applications of Simulations Methods in Environmental Resource Economics*. Kluwer Academics Publisher, Dordrecht, Chapter 7.
- Vanhulsel, M., Janssens, D., Wets, G., 2007. Calibrating a new reinforcement learning mechanism for modelling dynamic activity-travel behaviour and key events. In: 86th Conference of the Transportation Research Board, Washington DC.
- Vij, A., 2013. Incorporating the Influence of Latent Modal Preferences in Travel Demand Models (Ph.D. thesis). Department of Civil and Environmental Engineering, University of California, Berkeley.
- Walker, J., 2002. The mixed logit (or logit kernel) model: dispelling misconceptions of identification. *Transp. Res. Rec.* 1805, 86–98.
- Yañez, M.F., Cherchi, E., Heydecke, B.G., Ortúzar, J. de D., 2008. On the treatment of repeated observations in panel data: efficiency of mixed logit parameter estimates. In: XV Panamerican Conference on Traffic and Transport Engineering, Cartagena De Indias.
- Yañez, M.F., Mansilla, P., Ortúzar, J. de D., 2010. The Santiago Panel: measuring the effects of implementing Transantiago. *Transportation* 37, 125–149.