

An integrated framework for forecasting travel behavior using Markov Chain Monte Carlo simulation and profile Hidden Markov Models

Ismaïl Saadi^{a,*}, Ahmed Mustafa^a, Jacques Teller^a, Mario Cools^a

^aUniversity of Liège, ArGEnCo, Local Environment Management & Analysis (LEMA) Quartier Polytech 1, Allée de la Découverte 9, BE-4000 Liège, Belgium

Abstract

Recent advances in agent-based micro-simulation modeling have further highlighted the importance of a thorough full synthetic population procedure for guaranteeing the correct characterization of real-world populations and underlying travel demands. In this regard, we propose an integrated approach including Markov Chain Monte Carlo (MCMC) simulation and profiling-based methods to capture the behavioral complexity and the great heterogeneity of agents of the true population through representative micro-samples. The population synthesis method is capable of building the joint distribution of a given population with its corresponding marginal distributions using either full or partial conditional probabilities or both of them simultaneously. In particular, the estimation of socio-demographic or transport-related variables and the characterization of daily activity-travel patterns are included within the framework. The fully probabilistic structure based on Markov Chains characterizing this framework makes it innovative compared to standard activity-based models. Moreover, data stemming from the 2010 Belgian Household Daily Travel Survey (BELDAM) are used to calibrate the modeling framework. We illustrate that this framework effectively captures the behavioral heterogeneity of travelers. Furthermore, we demonstrate that the proposed framework is adequately adapted to meeting the demand for large-scale micro-simulation scenarios of transportation and urban systems.

Keywords: Travel behavior analysis, travel demand modeling, activity sequences, daily activity-travel patterns, simulation-based population synthesis, profiling analysis

1. Background

Activity-based analyses of travel behavior within large-scale micro-simulation models are particularly adapted to understanding the dynamics and the transitional patterns of travel behavior. In this context, most activity-based models require a full (synthetic) population to obtain greater behavioral realism within such complex transport and urban systems. However, the lack of highly disaggregate data suggests the use of more efficient methods for synthesizing individual/household socio-demographic attributes as well as their daily activity information (Beckman et al., 1996). Obtaining accurate individual-level information for a large population is a great challenge especially in the context of restrictive data availability. Typically, such detailed data can be derived from national censuses. However, in practice, only aggregate information is available to researchers and practitioners (Anderson et al., 2014).

*Corresponding author

Email address: ismail.saadi@ulg.ac.be (Ismaïl Saadi)

19 Furthermore, micro-samples may also include incomplete observations, making the true population
20 identification process through a virtual process more complex. To overcome this problem, several tech-
21 niques have been developed to address multiple data sources and incomplete datasets. Generally, detailed
22 individual-level micro-data and aggregate large-scale datasets are the two components used as input for
23 population synthesizing purposes. In the literature, the classical techniques identified are Iterative Proportional
24 Fitting (IPF) (Beckman et al., 1996; Mohammadian et al., 2010), Iterative Proportional Updating
25 (IPU) (Ye et al., 2009), Combinatorial Optimization (CO) (Williamson et al., 1998; Voas and Williamson,
26 2000) and probabilistic models (Farooq et al., 2013; Sun and Erath, 2015; Saadi et al., 2016b).

27 In this paper, we opt for a simulation-based approach because of its ability to address incomplete
28 datasets using partial conditional distributions (Farooq et al., 2013). More specifically, this operational
29 approach based on a Gibbs sampler can generate agents using partial views of the true joint distribution as
30 if the synthesized agents were drawn from the real population. In this context, the method enables the gen-
31 eration of populations of any given size. A comparative study established by Farooq et al. (2013) between
32 IPF and simulation-based approaches revealed that the latest technique clearly outperforms IPF using, for
33 both methods, almost the same amount of data. We build a representation of the true population by im-
34 plicitly merging different information provided by multiple micro-data sources using the simulation-based
35 approach. Then, from this representative/synthesized population, we can group the population into ho-
36 mogeny clusters characterized by similar attributes to understand their related travel behavior and to enable
37 a comparative attribute-based study based on the activity sequences.

38 For several decades, characterizing people’s activity-travel behavior has been an important issue to re-
39 searchers (Pendyala and Goulias, 2002; Auld et al., 2015; Saadi et al., 2016a). To include the sequential
40 dependencies of daily activities, some studies have suggested the use of Sequence Alignment Methods
41 (SAMs) (Joh and Timmermans, 2011; Joh et al., 2002; Wilson, 1998), whose inputs representing the daily
42 activity behavior in the form of activity sequences are extracted from activity-travel diaries (Bhat and Singh,
43 2000; Spissu et al., 2009). Typically, the application of SAMs includes pairwise sequence alignment algo-
44 rithms for scoring and comparing activity sequences in-between them. In this way, the main activity patterns
45 can be identified quantitatively and qualitatively within their respective subset of activity chains (Joh et al.,
46 2006).

47 When a relationship can be established between the derived patterns and the variables of interest, general
48 activity-travel behavior trends can be analyzed to investigate more deeply the determinants influencing those
49 travel patterns (Wilson, 1998). Although SAMs have been extensively used for characterizing activity-
50 travel patterns, such methods clearly suffer from serious limitations. Indeed, SAMs are unable to identify
51 the complete nature of a cluster. Even if the information regarding the frequent activity patterns is extracted,
52 irregular activities are not considered; thus, only parts of the clusters are characterized (Liu et al., 2015). In
53 this context, a number of questions have been raised regarding the consistency of such sequence alignment
54 techniques.

55 In an attempt to obtain improvements, Liu et al. (2015) proposed a profiling method called profile Hid-
56 den Markov Models (pHMMs) for characterizing the complete set of activity sequences including irregular
57 activities. pHMMs belong to the family of sequence profile methods that are essentially used to charac-
58 terize protein sequences. Furthermore, pHMMs may be described as position-specific scoring parameters
59 emanating from a multiple alignment of a group of protein sequences (Durbin et al., 1998).

60 From a methodological point of view, the positions of alignment are categorized according to three
61 possible states: match, insertion and deletion. Fundamentally, the configuration of the model is a successive
62 layout of match states representing the conserved activities that have been identified within their respective
63 cluster. This successive layout forms the reference activity sequences or a base template characterized by
64 the most frequent activities. In parallel, insertion states model the introduction of new conserved residues,
65 which can be interpreted as additional activities incorporated into the previously defined base template.

66 Conversely, deletion states represent the omission of particular conserved activities from the base template.
 67 In this regard, random differences between the activity sequences within the overall characterization process
 68 are implicitly included throughout these two states. Subsequently, estimated pHMMs can generate new
 69 activity sequences so that the relationships with their corresponding cluster are preserved (Liu et al., 2015).

70 Furthermore, pHMMs can be considered as quantitative descriptors that assign weights - probabilities
 71 of occurrence - to each activity at each corresponding match state characterizing both regular and irregular
 72 activities. For instance, considering daily activity-travel sequences, some travelers might see their behavior
 73 differing from the general behavior of a cluster (e.g., work cluster). This difference is captured by pHMMs
 74 by introducing new activities or, inversely, canceling or ignoring others. Subsequently, the scored activity
 75 sequences can be assessed to measure how similar or dissimilar these are from the cluster (Liu et al., 2015).

76 In this paper, we propose an integrated framework including both a population synthesis approach (Fa-
 77 rooq et al., 2013) and a profiling method (Liu et al., 2015) capable of describing and assigning the activity
 78 sequences for each individual of the synthetic population. The model is capable of considering, in an effi-
 79 cient manner, the behavioral effects of different influencing factors, which might vary between clusters. In
 80 this regard, the main trends in terms of travel behavior can emerge from the characterization achieved using
 81 the pHMMs.

82 The remainder of this paper is organized as follows. First, the data, obtained from the Belgian National
 83 Household Travel Survey, are described. Subsequently, in Section 3, the details of the integrated framework
 84 are provided to enable implementation in different contexts. Then, the main results and the validation of
 85 the integrated framework are discussed (Section 4). Finally, the concluding remarks are formulated in
 86 Section 5.

87 2. Data

88 To investigate individuals' daily activity-travel behaviors, we use data from the Belgian National House-
 89 hold Travel Survey of 2010 (Cornelis et al., 2012). The data include 37,680 recorded trips with respect to
 90 15,821 individuals distributed across the country. For each respondent, the data include a detailed sequenc-
 91 ing of trips with their respective characteristics (e.g., time expenditure, start-end locations, trip purpose,
 92 and mode preferences). With respect to the variables of interest, age, gender, socio-professional status
 93 and working time expenditure are considered in the modeling framework (Table 1). Furthermore, public
 94 transport subscription and driving license ownership are also synthesized as transport-related variables.

Variable	Basic Statistics
Age	Mean: 46.54, Std. Dev.: 21.08
Gender	Male: 47.87%, Female: 52.13%
Socio-Professional status	Not schooled children: 0.08% - Student: 17.23% - Housewife (husband): 4.30% Job seeker: 5.54% - Pensioner: 28.14% - Disabled person: 2.23% Blue-collar worker: 7.72% White-collar worker (executive): 3.51% White-collar worker (non-executive): 21.72% Self-employed person: 3.96% - Liberal profession: 1.11% - Teacher: 3.76% Farmer: 0.23% - Other: 0.48%
Work time expenditure	Mean: 36.68 h/week, Std. Dev.: 11.36 h/week
Public transport subscription	No: 77.33%, Yes: 22.67%
Driving license ownership	In progress: 4.90%, No: 24.87%, Yes: 70.23%

Table 1: Data description of the arbitrary explanatory variables

95 Table 2 presents the percentage of the recorded trips where, the outcome is the trip purpose and the
 96 explanatory variables are age, gender, socio-professional status of the travelers, working time expenditure,
 97 public transport subscription and driving license ownership. The values of the different explanatory vari-
 98 ables are expressed in terms of proportion (%) except for working time expenditure, which is expressed in

Trip purpose (*)	1	2	3	4	5	6	7	8	9	10	11	12
<i>Age</i>												
6-31 years	4.89	40.03	9.05	1.30	14.51	2.09	7.19	1.54	7.51	2.13	7.06	2.68
31-45 years	13.46	37.20	16.77	2.96	0.63	2.41	10.62	2.42	4.99	1.99	4.25	2.30
45-59 years	7.25	38.83	16.73	3.06	0.49	1.92	13.06	3.59	5.74	3.15	3.82	2.37
59+	5.36	40.36	3.10	0.72	0.32	1.88	20.27	5.94	7.94	5.30	4.71	4.10
<i>Gender</i>												
Male	6.64	39.32	12.79	3.10	4.09	2.32	11.14	2.95	6.31	3.05	5.44	2.85
Female	8.77	38.96	10.14	1.04	4.25	1.79	14.04	3.70	6.78	3.15	4.54	2.85
<i>Socio-professional Status</i>												
Student	3.51	41.72	1.49	0.19	23.78	1.75	5.72	1.04	7.21	2.08	8.58	2.92
Housewife (husband)	13.79	39.41	0.74	0.00	0.99	1.48	19.21	4.93	7.88	3.69	4.19	3.69
Job seeker	9.30	39.47	2.09	0.76	3.04	2.28	15.94	4.17	9.11	5.69	4.17	3.98
Pensioner	5.82	40.62	0.99	0.23	0.33	1.83	20.92	6.29	8.11	5.53	5.11	4.22
Disabled person	9.55	40.76	0.64	0.00	1.27	2.55	17.20	7.01	7.64	5.10	5.10	3.18
Blue-collar worker	7.08	40.05	24.52	4.63	0.27	1.09	8.72	1.91	5.45	1.77	3.00	1.50
White-collar worker (executive)	8.64	35.95	21.41	3.93	0.59	3.73	9.63	2.55	4.13	1.77	5.11	2.55
White-collar worker (non-executive)	9.43	37.45	20.48	2.91	0.39	2.26	10.72	2.51	5.52	2.08	4.23	2.01
Self-employed person	8.82	36.65	21.49	6.79	0.23	2.26	8.82	2.71	4.52	1.81	3.17	2.71
Liberal profession	7.78	34.73	20.36	9.58	0.60	2.99	9.58	2.99	4.79	1.80	2.99	1.80
Teacher	10.71	38.16	16.35	1.50	0.56	2.07	12.22	3.01	6.58	2.07	4.32	2.44
Farmer	10.53	42.11	10.53	5.26	0.00	0.00	10.53	0.00	10.53	5.26	5.26	0.00
Other	8.00	38.00	8.00	2.00	2.00	6.00	12.00	4.00	4.00	4.00	8.00	4.00
<i>Working time expenditure</i>												
Mean (h/week)				42.19								
Std. Dev. (h/week)				11.12								
<i>Public transport subscription</i>												
No	8.48	38.92	11.71	2.28	2.92	1.99	13.04	3.32	6.62	3.07	4.77	2.89
Yes	5.04	40.28	10.40	1.24	8.55	2.30	10.98	3.33	6.25	3.20	5.71	2.72
<i>Driving license ownership</i>												
In progress	5.69	40.66	6.95	0.63	11.52	2.67	10.18	2.95	7.23	2.74	5.62	3.16
No	3.49	42.39	4.05	0.31	15.87	1.60	10.36	2.69	6.84	3.14	6.61	2.66
Yes	8.81	38.39	13.39	2.53	1.05	2.14	13.23	3.49	6.43	3.11	4.56	2.88

(*) 1-bring/get, 2-home, 3-work, 4-for work, 5-education, 6-meal, 7-daily-shopping, 8-service, 9-visit, 10-tour, 11-entertainment, 12-other

Table 2: Cross-classification of recorded trips by purpose within each demographic segment (in %)

99 hours/week. Regarding the general distribution of the trips, one can clearly observe that the trips toward
100 home are the most important in terms of proportions.

101 Furthermore, it should be emphasized that commuting patterns also account for a relatively significant
102 share of trips and are mainly represented by professionally active people (see socio-professional status).
103 Young people essentially undertake trips (14.51%) with the objective of attending a school or university to
104 study. In parallel, they are active for visiting friends and family (7.51%) and also participating in extra-
105 activities (7.06%) (e.g., sports and entertainment). Individuals belonging to the oldest age category spend a
106 significant amount of time on daily shopping (20.27%) and visiting (7.94%).

107 Regarding working time expenditure, we can observe that the average number of hours is relatively high
108 (42.19 hours/week). This variable is an indicator of activity time expenditure (work). As mentioned later
109 in the paper, the pHMMs do not represent exact temporal information. Thus, it is technically possible to
110 synthesize a variable that can be used afterward to classify the population and obtain some trends in terms
111 of activity durations.

112 **3. Methodology**

113 *3.1. Framework*

114 To develop the integrated travel demand modeling framework presented in Figure 1, two types of
115 datasets are used: [A] socio-demographic/transport-related variables derived from the individuals' file of
116 the Belgian Household Travel Survey and [B] the trips file including the activity-travel diaries. As detailed
117 before, we propose to select a set of pertinent socio-demographic explanatory variables as the basis for com-
118 paring the different clusters in terms of activity-travel behavior. Subsequently, activity sequences will be
119 derived from the trip diaries. In this context, we will obtain a set of individuals with socio-demographic in-
120 formation (e.g., age, gender, socio-professional status and working time expenditure) and transport-related
121 variables (i.e., public transport subscription and driving license ownership) associated with their respective
122 activity sequences.

123 At this level, the data processing will be performed in parallel and in a completely independent manner.
124 On one side, a synthetic population [2] is built using a simulation-based approach [1]. The population
125 synthesis procedure plays a key role within the modeling framework because it provides a better estimate
126 of the heterogeneity of the population in comparison to standard population synthesis techniques such as
127 IPF (Farooq et al., 2013). In this regard, the synthesized population represents a better approximation and
128 is true regardless of the selected attribute. A direct implication is that this approach enables the estimation
129 of the precise proportion of clusters regarding the whole population.

130 In parallel with the population synthesis, we derive the activity sequences from the activity-travel diaries
131 [B] so that every respondent is associated with a daily activity-travel pattern/plan. At the end of step [3],
132 we will have a full synthetic population as well as a detailed list of individuals characterized by socio-
133 demographic attributes and transport-related variables as well as their respective activity sequences.

134 In the following step [4], it is possible to establish completely homogeneous clusters to isolate the effects
135 of the various explanatory variables to achieve more accurate analysis. Moreover, it is also important to
136 measure the coupled effects of mixed factors for investigating potential interactions.

137 Regarding the characterization of the activity sequences with profile Hidden Markov Models, a proce-
138 durally less complicated version of Liu et al. (2015) is implemented to gain computational efficiency in
139 the calibration phase. In the approach proposed by Liu et al. (2015), before calibrating the Markov Chain
140 profiles, it is necessary to classify the activity sequences according to their longest activities. Then, an
141 identification process of the most recurring activities, based on the definition of the regularity (Liu et al.,
142 2015), is recommended along with the determination of the most frequent activity transitions. Finally, for
143 every subdivision, templates are defined so that the activity sequences related to the clusters are aligned by
144 employing multiple sequence alignment methods based on their respective template. It is only after these
145 three steps that performed that the calibration of the Markov profiles becomes possible. We can clearly
146 notice that the method of Liu et al. (2015) is relatively heavy to implement in its entirety. In this context,
147 we suggest a less heavy methodology to implement while maintaining the key component of the model-
148 ing chain represented by the pHMM. This can be realized by estimating the effects of a combination of
149 explanatory variables of the activity sequences.

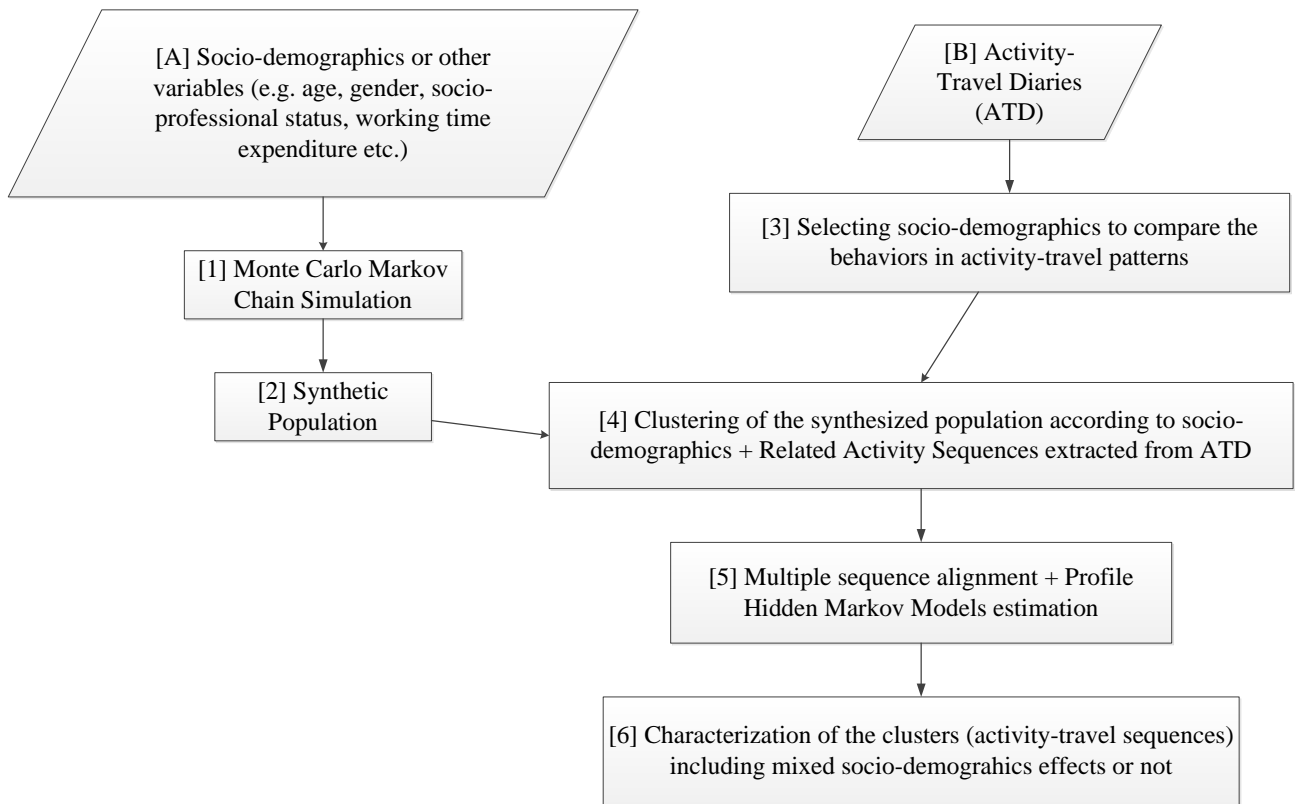


Figure 1: Overall modeling framework

150 With respect to the estimation of the profile Hidden Markov Models [5], the procedure consists of
 151 aligning the activity sequences within each cluster using an existing multiple sequence alignment approach
 152 (SAM). Subsequently, the transition and emission probabilities of the different pHMMs can be estimated
 153 so that every group is characterized by its dedicated model. In this context, the daily activities of the
 154 individuals are qualitatively and quantitatively simulated. In addition, both regular and irregular activities
 155 are implicitly included in the estimation conducted by the pHMM.

156 The strength of this framework lies in its capacity to characterize daily activity-travel sequences of very
 157 heterogeneous groups; in contrast, a purely descriptive statistical analysis would become quasi-incomprehensible
 158 and difficult to dress.

159 3.2. Synthetic population

160 Posterior samples are generated by a Gibbs sampling procedure using partial views of the true joint
 161 probability, i.e., the conditional probabilities. Theoretically, full conditional distributions should be defined
 162 and included within the algorithm; however, full conditionals are rarely available in practice. Thus, partial
 163 or even marginal distributions are used as substitutes (Farooq et al., 2013). Let $q(\mathbf{x})$ be an initial set of
 164 random attributes, $\mathbf{x} = (x_1, x_2, \dots, x_i, x_n)$ the set of attributes and π the conditional probability. Then, the
 165 algorithm is structured as follows:

166 **Step 1: initialize** $\pi(\mathbf{x}) \sim q(\mathbf{x})$

167 **Step 2: samples from CD** q

168 $x_1 \sim \pi(x_1 \mid x_2, x_3, \dots, x_i, \dots, x_n)$

169 $x_2 \sim \pi(x_2 \mid x_1, x_2, \dots, x_i, \dots, x_n)$

170 $x_3 \sim \pi(x_3 \mid x_1, x_2, \dots, x_i, \dots, x_n)$

171 ...

172 $x_i \sim \pi(x_i \mid x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$

173 ...

174 $x_n \sim \pi(x_n \mid x_1, x_2, \dots, x_i, \dots, x_{n-1})$

175 **Step 3: repeat until reaching the population size**

176 As mentioned previously, four socio-demographic and two transport-related variables are included in the
 177 form of conditionals as input for the Gibbs sampler. In this regard, data preparation represents an advantage
 178 compared to the amount of data necessary for the IPF (Farooq et al., 2013). In particular, we attempt to
 179 draw agents from the partial views of the true joint probability using Gibbs sampling. This technique is
 180 particularly flexible for handling multiple data sources with different spatial scales (Farooq et al., 2013). A
 181 simulation-based approach typically needs a set of conditional probabilities $\pi(x_i \mid x_j)$, where i and j are
 182 respectively referring to the dependent and independent variables.

183 To reach steady state, it is advisable to warm the Gibbs sampler (i.e., using approximately 1000 itera-
 184 tions in our case). To reduce computation time, we save the generated population in a .csv file. In this way,
 185 it is possible to extend the size of the population by starting with this file. We simply have to extract the
 186 characteristics of the last agent to use it as the initial condition for the Gibbs sampler. In addition, in the
 187 context of this study, we intend to generate a population of 250,000 agents to serve as a basis. Then, from

Dependent variable	Independent variables
Socio-professional status (full)	Age - Gender - WT expenditure (**) - PT subscription (*) - Driving license ownership
Gender (full)	Age - Socio-professional status - WT expenditure - PT subscription - Driving license ownership
Age (full)	Gender - Socio-professional status - WT expenditure - PT subscription - Driving license ownership
WT expenditure (full)	Age - Gender - Socio-professional status - PT subscription - Driving license ownership
PT subscription (full)	Age - Socio-professional status - WT expenditure - Gender - Driving license ownership
Driving license ownership (full)	Age - Socio-professional status - WT expenditure - Gender - PT subscription

(*) Public Transport subscription

(**) Working Time expenditure

Table 3: List of Conditional Distributions

188 this basis, a sub-sample of 110,000 agents is randomly extracted. This approach allows for a reduction in
189 the possible correlations between successive draws. The sub-sample represents approximately 1% of the
190 population in Belgium. Because we are only interested in the proportions of specific groups with respect
191 to the full population, it is not necessary to generate a larger number of agents. Indeed, the proportions
192 will remain asymptotically identical. In addition, if we were addressing an agent-based micro-simulation
193 problem, it would also have been necessary to approach the problem with a similar population size.

194 As outlined by Farooq et al. (2013), although the method provides a more accurate representation of
195 the population, the simulation-based approach is not yet capable of synthesizing a full household synthetic
196 population with multiple individuals. In this particular situation, an alternative technique, i.e., IPU, could be
197 adopted. Moreover, the synthetic population generation procedure can be easily replaced by IPU. Indeed,
198 the integrated framework presented in this paper is designed to ensure a sufficient level of modularity. Al-
199 ternatively, when the dataset includes the associations in-between households and their related individuals,
200 i.e., through a referencing system, it is possible to include mixed household and individual effects. First, the
201 target cluster of households is isolated. Then, the activity sequences associated with the individuals within
202 the selected households are processed. In this way, household effects can also be considered.

203 3.3. Activity sequences characterization

204 After sorting beforehand the activity sequences according to their main activity, the standard charac-
205 terization method of activity sequences is structured according to three different steps (Liu et al., 2015).
206 (a) The first step consists of measuring the regularity of the activities with respect to their related clusters.
207 In the same way, the most probable sequential order of the activities is identified for every group. In this
208 context, a complete template (reference activity chain) characterizing the most frequent activities as well as
209 their sequential order is defined.

210 Subsequently, (b) the following step consists in aligning the activity sequences based on the templates
211 and with respect to every cluster. This approach allows activity sequences that are perfectly aligned, with
212 identical dimensions, to be obtained.

213 Finally, (c) the aligned activity sequences are characterized by calibrating the pHMMs. The characteri-
214 zation implies the estimation of the transition, emission, insertion and deletion probabilities.

215 Note that, within the framework of our study, we do not group the agents according to their main
216 activities. Thus, it is not necessary to perform a cluster analysis to apply step (a) and allow the identification
217 of the template. On the contrary, we only attempt to estimate the effects of selected explanatory variables.
218 In this context, step (a) can be bypassed. Indeed, because the population is extremely heterogeneous, it
219 is not possible to define a template if activity chains of the same group possess different types of main
220 activities. One can refer to the research of Liu et al. (2015) for a more thorough description of the modeling
221 framework.

222 Figure 2 describes the full parameters of a pHMM, including the emission and transition conditional
223 probabilities, as well as the match m_i , insertion i_i and deletion d_i states. Note that both the match and

224 insertion states are capable of emitting an activity type A_i . m_0 is the beginning state, and m_N the final
225 state. N is the length of the chain containing the largest sequence of activities within the cluster. The
226 bold arrows represent an illustration of all the possible transition combinations between the states m_2 and
227 m_3 . Furthermore, an insertion state can evolve toward the same insertion state (symbolized by the loop);
228 otherwise, the following match state is selected. The parameters of the pHMM and the SAM have been
229 estimated using the Bioinformatics Toolbox of MATLAB.

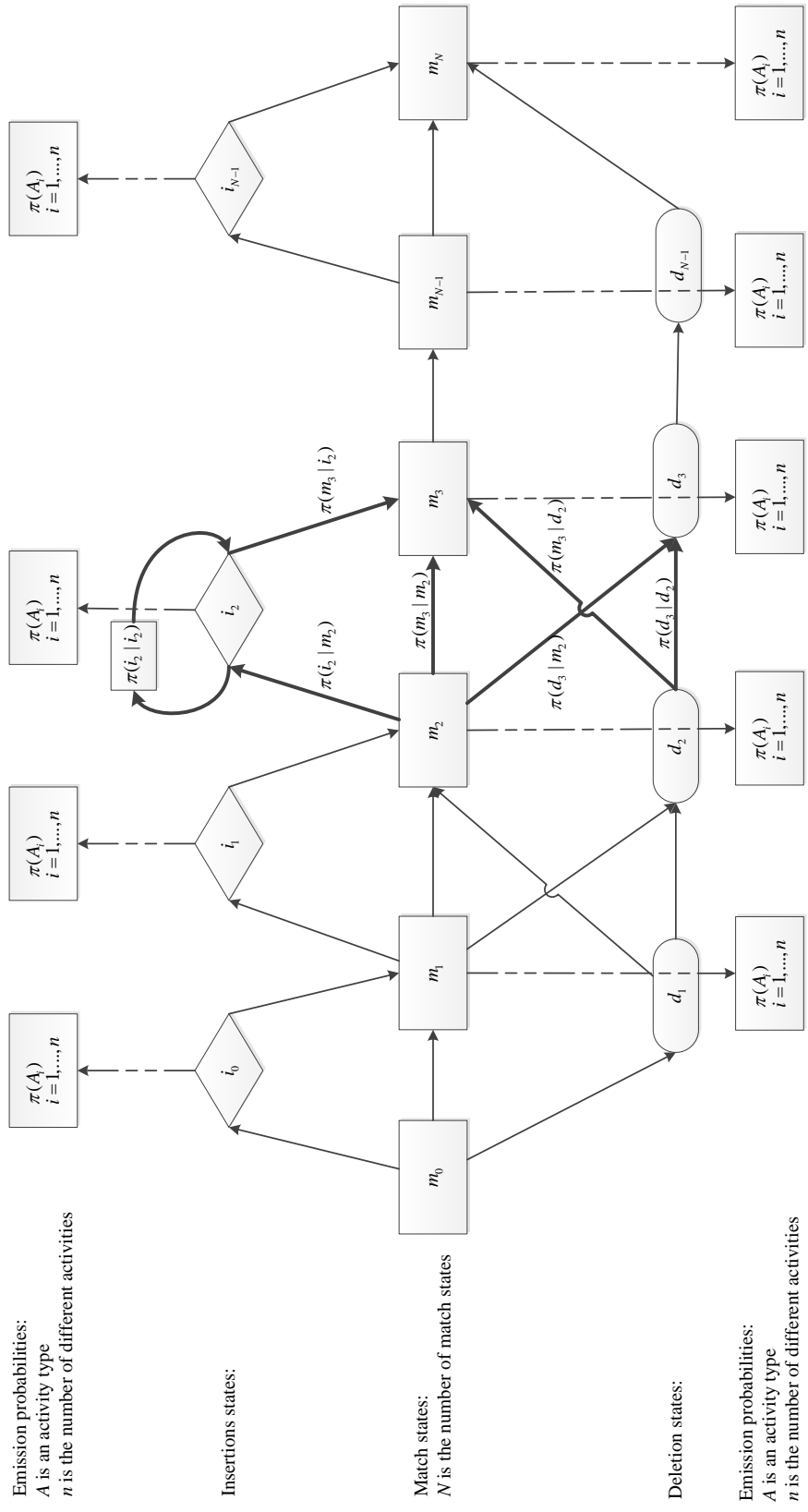


Figure 2: Parameters of the profile Hidden Markov Model

230 **4. Results**

231 *4.1. Synthetic population*

232 Figure 3 presents the comparison between the synthetic population, considering every explanatory vari-
 233 able, and the reference dataset. The results clearly indicate that the generated synthetic population is a suit-
 234 able approximation. Furthermore, the simulation-based approach provides good estimates of the marginal
 235 distributions for the selected attributes. With respect to the joint distributions presented in Fig. 4, the results
 236 demonstrate a good fit between the synthetic population and the reference dataset. Indeed, the slope is close
 237 to 1, with an R^2 value of 0.87. Note that each point represents the proportion of the combination of six at-
 238 tributes. In this study, only full conditions have been implemented in the framework. Therefore, the spread
 239 of the data points cannot be explained by the use of partial conditions. Moreover, the combined effects of
 240 scalability and dimensionality can explain slight deviations in the joint distributions of the simulation-based
 241 approach.

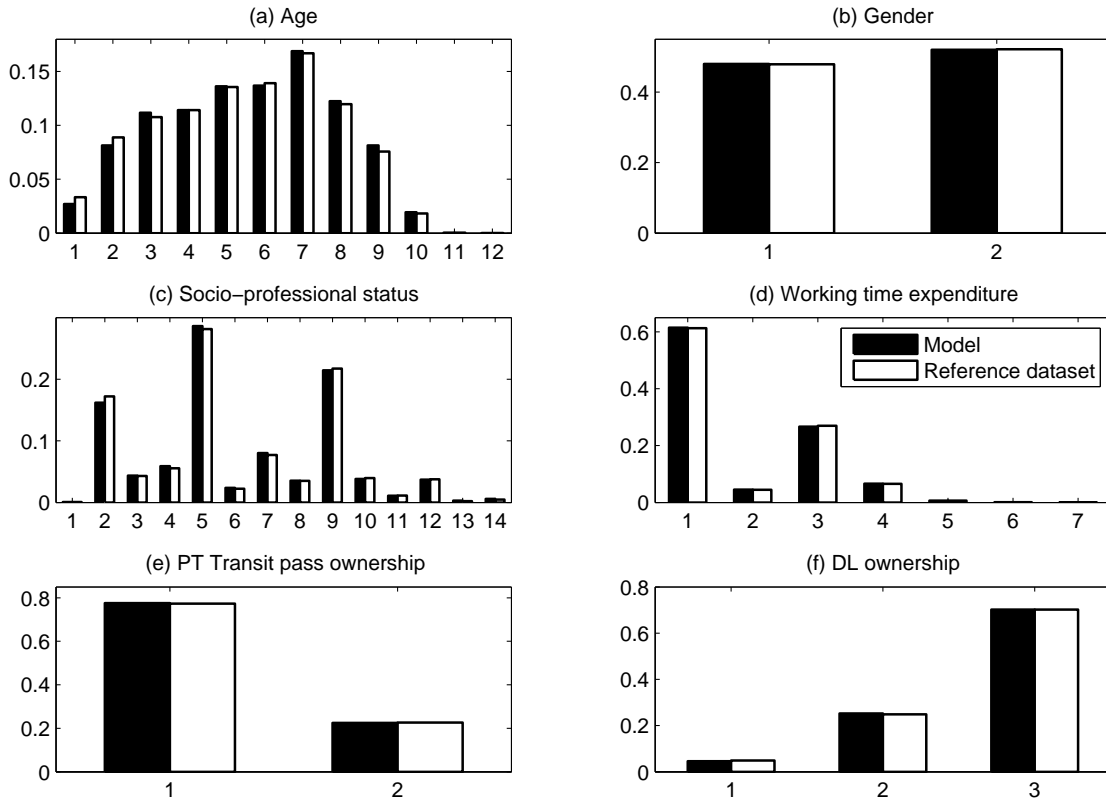


Figure 3: Comparison between the synthetic population and the reference dataset

242 In addition, some slight deviations may be observed within the marginal distributions (Fig. 3) as well as
 243 the joint conditional distributions (Fig. 4) due to the random effects included in the Gibbs sampler, which
 244 are mainly related to the stochastic nature of the model. In this regard, we assume that an increase in the
 245 size of the micro-data can play an effective role in the mitigation of the error.

246 Furthermore, studies have shown that an increase in the sample fraction is conducive to mitigating the
 247 input uncertainty (Rasouli and Timmermans, 2014). To enhance the stability of the forecasts, we propose
 248 to run the Gibbs sampler repeatedly such that the final outcome is the average of multiple model runs. As

249 outlined by Rasouli et al. (2012), this procedure also contributes to decreasing uncertainty but from the
 250 model perspective only.

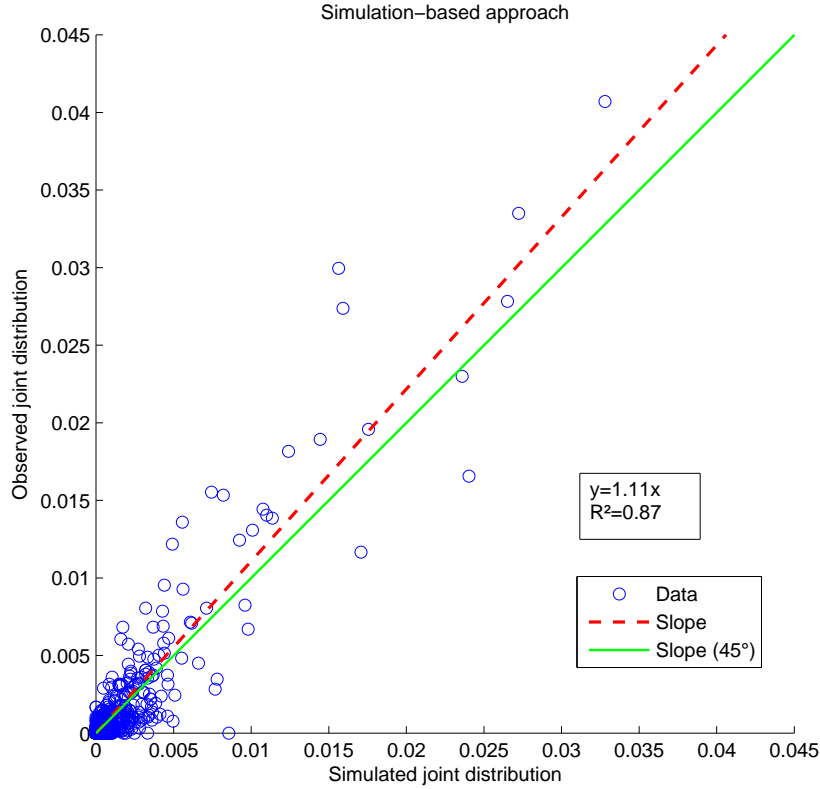


Figure 4: Comparison between the simulated and observed joint distributions

251 To assess the reliability of the population synthesis method, a more disaggregate approach consisting in
 252 a detailed clustering of the simulated and observed populations with respect to all the levels of the attributes
 253 defined in Table 3. Based on this clustering procedure, the deviations between the two populations from
 254 their related joint distributions are estimated. In this regard, the Root Mean Square Error (RMSE) is an
 255 indicator which allows to assess how close the synthesized population is to the observed one. The RMSE
 256 is defined as follows:

$$RMSE = \sqrt{E((\tilde{\theta} - \theta)^2)} = \sqrt{\frac{\sum_{i=1}^n (\tilde{y}_i - y_i)^2}{n}} \quad (1)$$

257 where $\tilde{\theta}$ is the estimator of the population, θ is the observed population, and n is the number of predicted
 258 values. Vovsha et al. (2015) used the RMSE to estimate the goodness-of-fit of the synthesized populations
 259 for different zonal systems. The RMSEs presented in Table 4 show that, even in the context of a finer
 260 analysis, the model is capable of maintaining a satisfactory level of accuracy. Overall, most of the errors are
 261 close to 0, except for the values combined with level 13 of the socio-professional variable (i.e., the farmers).
 262 However, this only represents a small portion of the full population (0.23% (Table 1)). In this regard, when
 263 the number of observations within a specific combination of variables is low, it is more probable to observe
 264 important deviations in terms of RMSE.

Socio-professional status \times Age	1	2	3	4	5	6	7	8	9	10	11	12
1	-	0.027	0.058	-	-	-	-	-	-	-	-	-
2	0.013	0.026	0.017	0.027	0.049	0.074	0.080	-	-	-	-	-
3	0.036	0.059	0.018	0.034	0.055	0.024	0.025	0.032	-	-	-	-
4	0.027	-	0.023	0.018	0.021	0.017	0.011	0.055	-	-	-	-
5	-	-	-	0.056	0.031	0.023	0.016	0.007	0.015	0.041	0.070	0.154
6	0.118	-	0.056	0.031	0.023	0.016	0.014	0.061	-	-	-	-
7	0.019	-	0.018	0.030	0.027	0.032	0.037	0.078	-	-	-	-
8	0.036	-	0.036	0.032	0.035	0.029	0.033	0.071	-	-	-	-
9	0.020	-	0.007	0.007	0.012	0.012	0.011	0.037	-	-	-	-
10	0.038	-	0.043	0.043	0.040	0.039	0.047	0.026	-	-	-	-
11	0.051	-	0.069	0.040	0.052	0.035	0.038	0.094	-	-	-	-
12	0.021	-	0.038	0.023	0.015	0.018	0.023	0.072	-	-	-	-
13	0.154	-	0.126	0.045	0.062	0.065	0.066	0.154	-	-	-	-
14	0.081	-	0.063	0.046	0.043	0.058	0.056	-	-	-	-	-

Table 4: Root Mean Square Error (RMSE) of the Simulated Joint Distributions with respect to different combinations of age and status

265 4.2. Characterization of the activity sequences

266 After building the synthetic population, we define criteria of analysis (e.g., gender and socio-professional
267 status) to extract the activity-travel patterns. The objective of the population synthesis is to estimate the
268 proportion of the categories of the studied people within the overall population. Subsequently, the activity
269 sequences are extracted from the activity-travel diaries describing the succession of the activities of the
270 studied clusters.

271 Regarding the simulation tools, various programs in bioinformatics are available for aligning multiple
272 activity sequences and also estimating the parameters of pHMMs. Thus, it is not necessary to develop a full
273 code for sequence alignment and parameter estimation. This proves that this framework can be applied in a
274 rather fast and effective manner.

275 As advised by Liu et al. (2015), some assumptions have been made regarding all the activity sequences.
276 Indeed, we suppose that all the sequences begin and end with the same activity, in this particular case, home
277 (H). Thus, the first and last positions possess a probability that is approximately equal to 1. In this regard,
278 no deletion states occur in the first position, and no transition probabilities need to be defined in the last
279 position because everything is shifting toward home.

280 Note that the number of positions is defined by the activity chain with the highest number of successively
281 different activities, i.e., 13 activities in this paper. If a smaller size is taken, the profile-HMM will have
282 to aggregate some information to be able to estimate the probabilities, thus reducing the quality of the
283 information. On the other hand, if a greater profile size is fixed, more parameters will have to be estimated;
284 however, this would not improve the accuracy and would make the analysis more complex. In this regard,
285 one can understand that determining the number of positions depends exclusively on the longest chain.
286 By referring to the work of Liu et al. (2015), the approach is similar except that the longest activity is
287 called a template and is built from an identification of the regular and irregular activities as well as their
288 sequential order. Note that in Liu et al. (2015), built clusters are only based on the main activities, i.e.,
289 the work cluster. By including socio-demographics in our paper, some additional travel patterns can be
290 revealed, e.g., education for young people. In this regard, the activity sequencing is sensitive to changes in
291 the socio-demographics.

292 To illustrate the methodology, we propose some case studies wherein the differences in behavior are
293 highlighted through the parameter estimates. Table 5 presents the transition and emission probabilities
294 resulting from an estimated pHMM, where the results for the full population of Belgium are included to
295 highlight the main patterns of conduct.

296 For each column of the emission probabilities, the highest values are in bold. In this way, we can
 297 extract some key information. For example, throughout all the positions, the commuting patterns are quite
 298 significant in terms of importance compared to the remainder of the population and specifically in position
 299 4.

300 The trends suggest that some activities, such as leisure, sports and visiting family and friends (see
 301 between positions 9-12), are preferred to be conducted at the end of the day. In contrast, the daily shopping
 302 activity is distributed throughout the day.

Position (k)	1	2	3	4	5	6	7	8	9	10	11	12	13
<i>Transition probabilities</i>													
$\pi(m_{k+1}, m_k)$	0.548	0.583	0.270	0.644	0.503	0.649	0.633	0.747	0.093	0.999	0.977	1.00	–
$\pi(i_k, m_k)$	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.031	0.000	0.000	0.000	–
$\pi(d_{k+1}, m_k)$	0.452	0.417	0.730	0.356	0.497	0.351	0.367	0.253	0.876	0.000	0.023	0.000	–
$\pi(m_{k+1}, i_k)$	0.500	0.500	0.500	0.500	0.500	0.500	0.500	0.500	0.510	0.500	0.500	0.500	–
$\pi(i_k, i_k)$	0.500	0.500	0.500	0.500	0.500	0.500	0.500	0.500	0.490	0.500	0.500	0.500	–
$\pi(m_{k+1}, d_k)$	0.263	0.041	0.475	0.035	0.056	0.076	0.098	0.115	0.018	0.066	0.186	0.648	–
$\pi(d_{k+1}, d_k)$	0.737	0.959	0.525	0.965	0.943	0.924	0.902	0.885	0.982	0.934	0.814	0.352	–
<i>Emission probabilities</i>													
$\pi(H)$ (home)	1	0.002	0.691	0.046	0.457	0.282	0.154	0.533	0.062	0.097	0.489	0.004	1
$\pi(W)$ (work)	–	0.018	0.018	0.651	0.254	0.287	0.237	0.055	0.080	0.065	0.125	0.237	–
$\pi(G)$ (bring/get)	–	0.140	0.036	0.079	0.029	0.052	0.095	0.093	0.204	0.065	0.057	0.091	–
$\pi(E)$ (education)	–	0.206	0.009	0.041	0.004	0.017	0.041	0.005	0.009	0.032	0.011	0.008	–
$\pi(M)$ (meal)	–	0.020	0.021	0.002	0.043	0.012	0.018	0.017	0.049	0.064	0.023	0.008	–
$\pi(S)$ (daily shopping)	–	0.337	0.124	0.072	0.082	0.098	0.219	0.121	0.227	0.097	0.046	0.079	–
$\pi(A)$ (service)	–	0.080	0.033	0.017	0.029	0.046	0.047	0.028	0.009	0.065	0.011	0.012	–
$\pi(F)$ (visit)	–	0.062	0.018	0.022	0.029	0.069	0.095	0.028	0.191	0.065	0.046	0.289	–
$\pi(V)$ (tour)	–	0.066	0.009	0.017	0.014	0.029	0.012	0.022	0.009	0.097	0.023	0.174	–
$\pi(R)$ (leisure, sports etc.)	–	0.051	0.009	0.029	0.025	0.052	0.024	0.044	0.116	0.032	0.057	0.059	–

Table 5: Parameter Estimates for the Full Population of Belgium

303 It has been reported in various studies (Bhat and Singh, 2000) that most professionally active individuals
 304 go shopping after work on the road returning home or even later in the day. Because this category of people
 305 is important in term of size with respect to the full population, it is thus logical to obtain significant values
 306 (i.e., 21.9% in position 7 – 22.7% in position 9) later in the day.

307 Furthermore, the bring/get activity is significantly present throughout the day as well. However, it is
 308 necessary to note that various groups of the population are implicitly included in the results highlighted
 309 in Table 5, which makes the clusters’ specific features more difficult to capture. In this context, a more
 310 detailed cluster analysis is necessary to allow one to distinguish which part of the population conducts
 311 shopping activities and during approximately which period of the day.

312 When characterizing the full population, it is indeed less obvious to consider what are the proportions
 313 as well as the categories of people who perform a given activity at a given moment of the day.

314 Table 6 presents the parameter estimates of the emission and transition probabilities for individuals less
 315 than 31 years of age. This category groups most of the students and also some young workers. After
 316 isolating the target sub-population, we can note the important increase in the education activity; this is
 317 synonymous with important trips toward schools and universities.

318 Note that young people also dedicate a portion of their time to conducting sports or entertainment activ-
 319 ities. Furthermore, they organize, as a general rule, such activities after their courses (see position 5). If we
 320 need any reminder of this, we simply need to observe the sequencing of the activities within the positions.
 321 The great majority of the educational activities are grouped in position 2 (43.3%). Subsequently, it is only
 322 from position 4 that young people perform secondary activities such as shopping. Indeed, an important

323 portion of young people (students/novice workers) live alone during their studies or at the beginning of
 324 their professional lives; they also have to fulfill their vital needs by moving quasi-daily to shop.

325 These results are logical and compatible with the descriptive analysis presented in Table 2. The main
 326 activity-travel patterns (e.g., education, visiting friends and family, entertainment and sports) have been
 327 characterized by the pHMM. Furthermore, this proves that the calibration of the pHMM was correct. This
 328 mode of comparison clearly reveals the added value of the pHMM compared to a classic analysis of de-
 329 scriptive statistics. Not only is the establishment of activity sharing possible throughout the day but the
 330 result also indicates that the sequencing of the activities can be obtained thanks to the positioning system.

Position (k)	1	2	3	4	5	6	7	8	9	10	11	12	13
<i>Transition probabilities</i>													
$\pi(m_{k+1}, m_k)$	0.876	0.529	0.333	0.362	0.462	0.090	0.998	0.352	0.411	0.757	0.988	1.000	–
$\pi(i_k, m_k)$	0.000	0.000	0.000	0.000	0.000	0.000	0.001	0.000	0.000	0.000	0.000	0.000	–
$\pi(d_{k+1}, m_k)$	0.124	0.471	0.667	0.638	0.538	0.910	0.001	0.648	0.589	0.243	0.012	0.000	–
$\pi(m_{k+1}, i_k)$	0.500	0.500	0.500	0.500	0.500	0.500	0.500	0.500	0.500	0.500	0.500	0.500	–
$\pi(i_k, i_k)$	0.500	0.500	0.500	0.500	0.500	0.500	0.500	0.500	0.500	0.500	0.500	0.500	–
$\pi(m_{k+1}, d_k)$	0.256	0.321	0.035	0.066	0.029	0.034	0.070	0.065	0.042	0.083	0.174	0.668	–
$\pi(d_{k+1}, d_k)$	0.744	0.679	0.965	0.934	0.971	0.997	0.930	0.935	0.958	0.917	0.826	0.332	–
<i>Emission probabilities</i>													
$\pi(H)$ (home)	1	0.000	0.598	0.297	0.180	0.395	0.046	0.099	0.287	0.158	0.474	0.007	1
$\pi(W)$ (work)	–	0.255	0.070	0.143	0.132	0.156	0.046	0.077	0.283	0.131	0.099	0.106	–
$\pi(G)$ (bring/get)	–	0.030	0.068	0.078	0.083	0.070	0.046	0.047	0.038	0.113	0.064	0.146	–
$\pi(E)$ (education)	–	0.433	0.033	0.064	0.132	0.091	0.046	0.077	0.007	0.018	0.020	0.024	–
$\pi(M)$ (meal)	–	0.011	0.036	0.023	0.019	0.049	0.046	0.017	0.026	0.045	0.035	0.061	–
$\pi(S)$ (daily shopping)	–	0.082	0.074	0.146	0.145	0.029	0.136	0.090	0.091	0.153	0.062	0.109	–
$\pi(A)$ (service)	–	0.020	0.015	0.026	0.035	0.017	0.046	0.022	0.019	0.018	0.027	0.026	–
$\pi(F)$ (visit)	–	0.070	0.060	0.092	0.038	0.033	0.046	0.249	0.113	0.131	0.072	0.255	–
$\pi(V)$ (tour)	–	0.030	0.008	0.064	0.024	0.008	0.046	0.009	0.030	0.063	0.030	0.042	–
$\pi(R)$ (leisure, sports etc.)	–	0.066	0.033	0.049	0.186	0.111	0.046	0.270	0.068	0.126	0.094	0.213	–

Table 6: Parameter Estimates for the Population of Belgium below 31 Years of Age

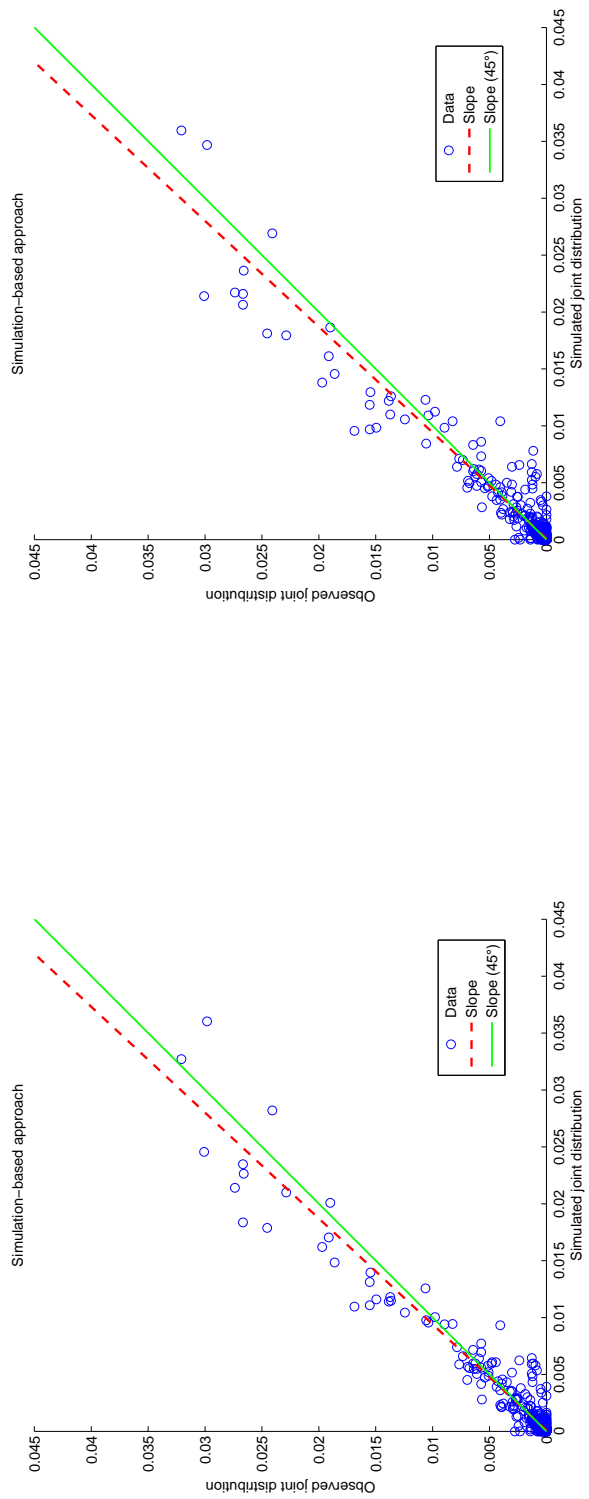
331 4.3. Comparison between the simulation-based approach and the IPU-based approach

332 To demonstrate the efficiency of the integrated framework with respect to standard approaches, we
 333 propose to compare, at each level of the framework, a sub-module with an existing technique. First, the
 334 simulation-based approach is compared with the iterative proportional updating (IPU) algorithm. Then,
 335 we show the added value of the pHMMs over a standard SAM. In this paper, it is important to distinguish
 336 the sequence alignment process (SAM) necessary for aligning all the activity sequences and the pHMM
 337 characterization. This step is fundamental because the activity sequences are structured such that deletion
 338 states are added to ensure equal lengths for each activity sequence. In this context, we can describe this
 339 integrated framework as a step above other previous approaches. Regarding the comparison between the
 340 IPU algorithm and the simulation-based approach, it is quite difficult to determine which method performs
 341 better from an absolute perspective. We must qualify the comparison in light of different aspects involved
 342 within each approach. For example, the IPU algorithm requires a huge amount of data (PUMS and all
 343 the total controls) to determine the weights associated with the corresponding agents (Ye et al., 2009). In
 344 contrast, the simulation-based approach is only based on the calibration of the conditional distributions
 345 (Farooq et al., 2013). In this context, only PUMS are used. Thus, a completely fair comparison is very
 346 complex. In addition, a throughout data preparation needs to be realized to classify all the variables in the
 347 correct format (table of frequencies) so that it can be included within IPU (Ye et al., 2009). In addition, all
 348 the total controls corresponding to all the levels of the set of attributes are needed, and the matching between

349 the format variables of the total controls and the PUMS should be the same. In practice, the implementation
350 is quite heavy and computational intensive. However, the simulation-based approach is characterized by its
351 great flexibility (Farooq et al., 2013). In addition, the PUMS can be used under their original format, i.e., as
352 a set of observations with respect to different attributes. In practice, it is important to propose approaches
353 that are capable of using a minimum amount of data to mitigate the phenomenon of data dependency in
354 such a way that the quality of the results is preserved. To our knowledge, no studies have investigated
355 the comparison between the simulation-based approach and IPU. Note that in Farooq et al. (2013), one
356 can find the comparison with IPF. To ensure a fair comparison between both approaches, we will compare
357 the methods for the synthesis of 4 attributes as advised by Farooq et al. (2013) so that zero-cell and zero-
358 marginal problems, which can lead to non-convergence, can be avoided. In addition, let us consider the
359 BELDAM travel survey because it describes the full population. In this context, the total controls can be
360 derived for IPU as well as the seed. We suppose that the seed that is extracted represents 50% of the survey.
361 In parallel, the same seed is included into the simulation-based approach. We will also consider a full seed
362 for the simulation-based approach to provide it with full information and establish a relatively equal amount
363 of inputs with respect to both methods. Table 7 presents the comparison between the different methods and
364 configurations. Note that the intercepts are not included because they are all approximately zero.

Method	Simulation-based (full PUMS)	Simulation-based (50% PUMS)	IPU-based (50% PUMS + total controls)
R-square	94.6%	93.7%	92.4%
Slope	1.075	1.075	0.932
RMSE	0.00112	0.00121	0.00139

Table 7: Comparison between the simulation-based approach and IPU



(a) 50% PUMS

(b) Full PUMS

Figure 5: Comparison of the joint distributions with respect to the simulation-based approach

365 The results presented in Table 7 clearly show the advantage of the simulation-based approach over IPU.
 366 In this regard, the conclusions are in accordance with those presented by Farooq et al. (2013). We can
 367 see that using the full PUMS for the simulation-based approach does not improve the R-squared value
 368 significantly. A 50% PUMS is largely sufficient to ensure accurate results. Furthermore, although the
 369 amount of data is low compared to IPU, the simulation-based approach is capable of providing a synthetic
 370 population with an RMSE that is reduced by -14.88% and an R-squared that is improved by 1.4%.

371 Figures 5a, 5b and 6 present the comparison between the joint distributions of the simulated populations
 372 and the reference dataset. Each circle is a proportion of a combination of 4 attributes. Using these figures,
 373 we can see that the simulation-based approach is capable of preserving good estimates for the important
 374 proportions, whereas IPU provides a poorer performance. However, IPU maintains very good estimates
 375 for small proportions, which is not the case for the simulation-based approach. However, given that the
 376 absolute differences are more important for high proportions and less important for small proportions, the
 377 RMSE is thus higher under the simulation-based approach than under IPU.

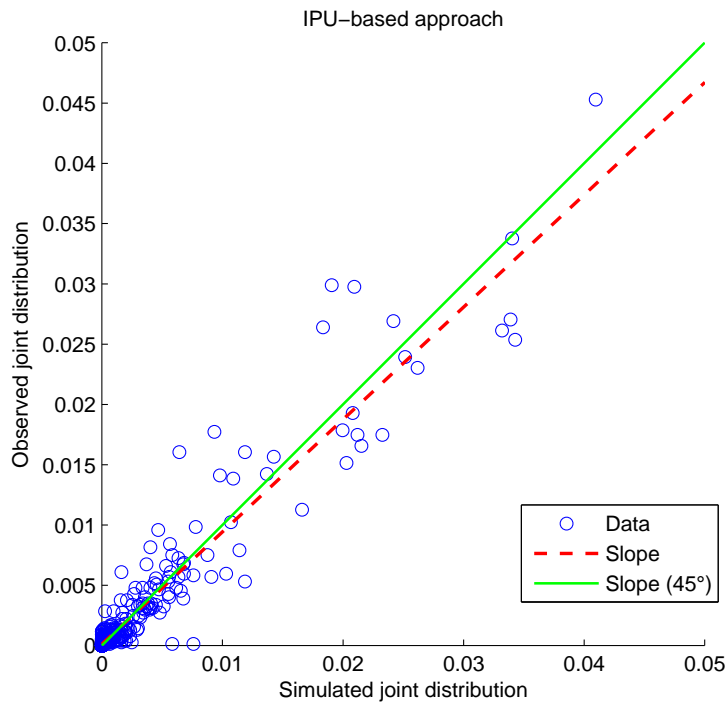


Figure 6: Comparison of the joint distributions with respect to IPU

378 *4.4. Profile-HMM validation results*

379 To ensure that the estimated transition and emission probabilities are accurate, we propose two im-
 380 portant indicators: (i) the proportions of the trip patterns and (ii) the occurrences of the different types
 381 of activities with respect to the full population. First, we generate a set of activity sequences from the
 382 estimated profile-HMM. Then, if the comparison between the indicators of the synthesized and observed
 383 activity sequences demonstrates that they are equivalent, then the parameter estimates have been estimated
 384 properly. As mentioned previously, we focus on the main trip patterns. In this regard, we compare all the
 385 trip patterns starting from home toward any other activity location and vice versa. The results presented in
 386 Figures ?? and 8 reveal that the main trip patterns have been correctly captured by the calibrated profile-
 387 HMM. This result is particularly important in the context of agent-based modeling, i.e., MATSim. Indeed,

388 the estimation of the traffic flows is especially affected by the quantity of trips and their related patterns. An
 389 over- or under-estimation may have a significant impact on the predictions in terms of traffic jams and/or
 390 traffic flows. To complete the comparison, we have also presented the fit of the simulated and observed
 391 proportions of trip patterns. Table 8 presents the main statistical metrics. Although some minor deviations
 392 might be depicted, we can consider that the model is able to produce good estimates of the trip patterns
 393 and their proportions within the overall modeling framework. The R-square values between brackets, pre-
 394 sented in Table 8, correspond to the regression models that include the intercept. Note that these R-square
 395 values are better in comparison to those of the models without intercept. This can be explained by the
 396 under-estimation in the simulated values for uncommon trip patterns (see Fig. 8).

Model	Slope	R-squared	RMSE (in %)
Trip patterns 1 (H to A)	0.80 (0.65)	0.86 (0.96)	0.23
Trip patterns 2 (A to H)	0.87 (0.68)	0.72 (0.82)	0.21

Table 8: Comparison between the synthesized activity sequences and the observed sequences for different statistical metrics

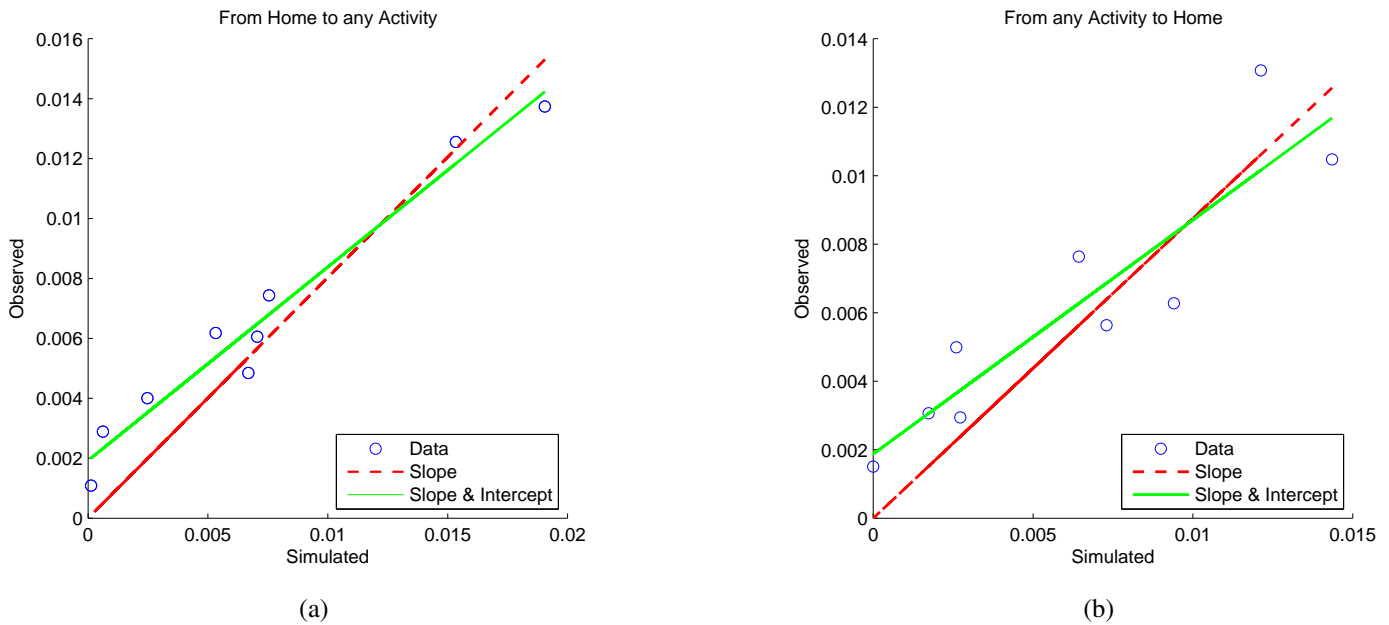


Figure 7: Comparison between the joint distributions with respect to the trip patterns

397 With respect to the comparison of the proportion of activities for the full population, the results presented
 398 in Figure 9 and Figure 10 also indicate a good match between the simulated and observed activity sequences
 399 (R-squared=0.99). In this regard, these results prove that the previously presented emission probabilities
 400 are sufficiently accurate and reliable.

401 5. Conclusion

402 In this paper, we present an integrated framework including a synthetic population approach (Farooq
 403 et al., 2013), together with a profiling method (Liu et al., 2015), for characterizing activity-travel patterns
 404 from both qualitative and quantitative perspectives.

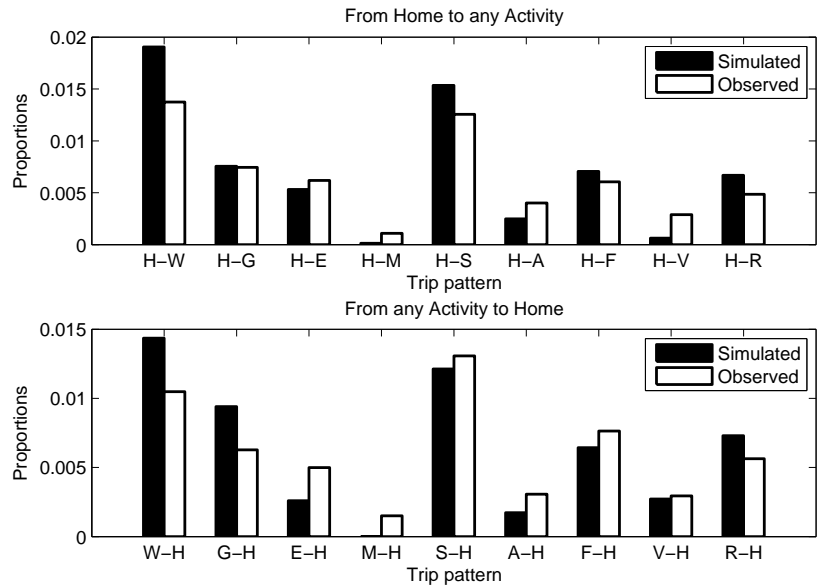


Figure 8: Comparison between the marginal distributions with respect to the trip patterns

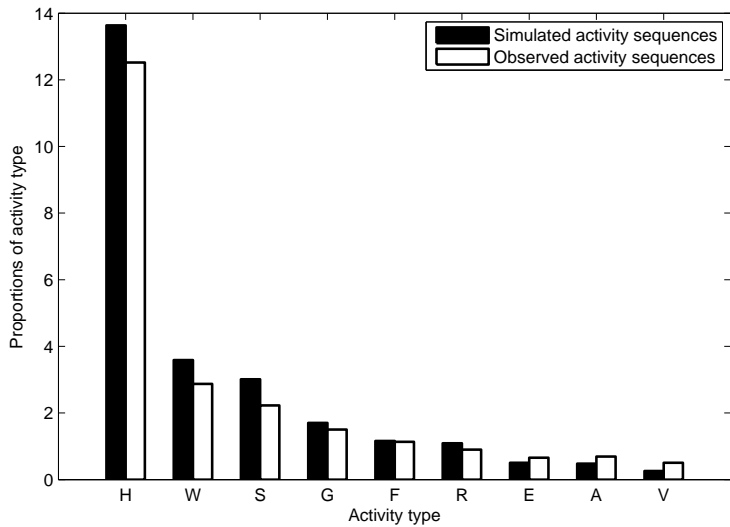


Figure 9: Comparison between the marginal distributions with respect to activity occurrence

405 The synthetic population sets up the connection between multiple micro-datasets. Indeed, the condi-
 406 tional probabilities are built so that they include the information resulting from all the available sources. In
 407 this regard, their correct determination represents a highly important result. Given the high performances
 408 of the simulation-based approach compared to standard methods (e.g., IPF), we opted for this technique to
 409 be introduced within the global modeling chain. Furthermore, the flexibility of this technique is particularly
 410 adapted to address partial micro-datasets. The results presented in Figure 4 clearly indicate that there is a
 411 scope that is able to fit the true population by implicitly merging different micro-datasets while ensuring
 412 high accuracy. Thus, these findings are in complete agreement with the conclusions of Farooq et al. (2013).
 413 A limitation of the simulation-based population synthesis is that the synthesis of households with multiple
 414 individuals is not yet possible. Nonetheless, as discussed in the methodological section, one could take into

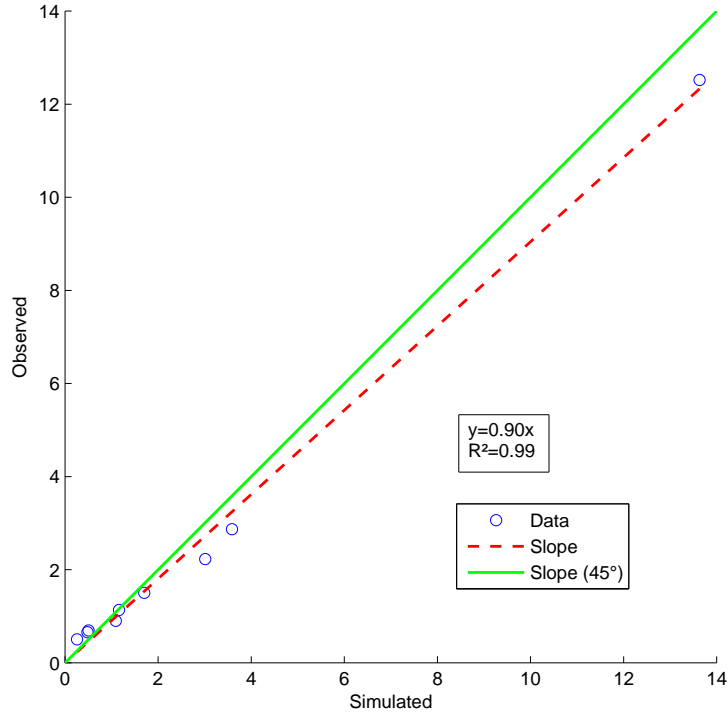


Figure 10: Comparison between the joint distributions with respect to activity occurrence

415 account household effects by using IPU in a similar way or by clustering the target households with their
 416 associated individuals. In this regard, we can conclude that the consecutive steps of our framework, i.e., the
 417 multiple-sequence alignment and the profile-HMM characterization steps, are not significantly affected. In
 418 the future and to preserve the presented framework, an extended population synthesis exclusively based on
 419 the simulation-based approach could be adopted at both levels (households and individuals) by consider-
 420 ing an additional module able to generate associations between households and individuals. An interesting
 421 prelude in this regard can be found in Anderson et al. (2014).

422 The profiling approach enables the characterization of multiple activity sequences (activity regularity
 423 and sequential information) from only one model without neglecting any irregular activities. As a result,
 424 performing a comparison between clusters in terms of activity-travel patterns is much easier or investigating
 425 to what extent the activity-travel patterns of a specific group can be distinguished from the general behavior
 426 of a population.

427 Furthermore, when the pHMM is calibrated using a training dataset, a non-limited number of activity
 428 sequences can be regenerated from the estimated pHMM according to the size of the studied cluster. This
 429 application is particularly interesting in the context of agent-based micro-simulation models. Indeed, most
 430 of them require a synthetic population describing the attributes and the activity sequences of every indi-
 431 vidual. In this regard, we assume that such a modeling framework can be adapted to handle problems of
 432 multi-agent model generation and, as a result, provide new insights for further research.

433 Regarding the results, we indicated in Section 4.2 that the positions within a profile-HMM give the
 434 general trends of the activity sequencing from a temporal perspective. However, we consider that the
 435 approach presents a limitation at this level. Further developments of the framework should aim at the
 436 inclusion of the activity time dimension in a more explicit way.

437 Furthermore, by isolating the target populations, the model allowed one to characterize the proportion of

438 these sub-populations of the total population as well as the main travel behaviors. With respect to the young
439 population, we have clearly shown that education appears to play an essential role in the need for travel.
440 Note that, in the context of more elaborate analysis, it is possible to estimate quantitatively the disparity
441 between different combinations of explanatory variables in terms of activity characteristics choice.

442 **6. Acknowledgments**

443 The research was funded by the ARC grant for Concerted Research Actions for project n°13/17-01
444 entitled 'Land-use change and future flood risk: influence of micro-scale spatial patterns (FloodLand)'
445 and by the Special Fund for Research for project n°5128 entitled 'Assessment of sampling variability and
446 aggregation error in transport models', both financed by the French Community of Belgium (Wallonia-
447 Brussels Federation).

448 **References**

- 449 Anderson, P., Farooq, B., Efthymiou, D., Bierlaire, M., 2014. Associations generation in synthetic popula-
450 tion for transportation applications: Graph-theoretic solution. *Transportation Research Record: Journal*
451 *of the Transportation Research Board* , 38–50doi:<http://dx.doi.org/10.3141/2429-05>.
- 452 Auld, J., Hope, M., Ley, H., Sokolov, V., Xu, B., Zhang, K., 2015. Polaris: Agent-based modeling frame-
453 work development and implementation for integrated travel demand and network and operations sim-
454 ulations. *Transportation Research Part C: Emerging Technologies* doi:[http://dx.doi.org/10.](http://dx.doi.org/10.1016/j.trc.2015.07.017)
455 [1016/j.trc.2015.07.017](http://dx.doi.org/10.1016/j.trc.2015.07.017).
- 456 Beckman, R.J., Baggerly, K.A., McKay, M.D., 1996. Creating synthetic baseline populations. *Trans-*
457 *portation Research Part A: Policy and Practice* 30, 415–429. doi:[http://dx.doi.org/10.1016/](http://dx.doi.org/10.1016/0965-8564(96)00004-3)
458 [0965-8564\(96\)00004-3](http://dx.doi.org/10.1016/0965-8564(96)00004-3).
- 459 Bhat, C.R., Singh, S.K., 2000. A comprehensive daily activity-travel generation model system for work-
460 ers. *Transportation Research Part A: Policy and Practice* 34, 1–22. doi:[http://dx.doi.org/10.](http://dx.doi.org/10.1016/S0965-8564(98)00037-8)
461 [1016/S0965-8564\(98\)00037-8](http://dx.doi.org/10.1016/S0965-8564(98)00037-8).
- 462 Cornelis, E., Hubert, M., Hunyen, P., Lebrun, K., Patriarche, G., De Witte, A., Creemers, L., Declercq, K.,
463 Janssens, D., Castaigne, M., Hollaert, L., Walle, F., 2012. Belgian Daily Mobility (BELDAM): Enquête
464 sur la mobilité quotidienne des belges. Technical Report. SPF Mobilité & Transports. Brussels.
- 465 Durbin, R., Eddy, S.R., Krogh, A., Mitchison, G., 1998. *Biological sequence analysis: probabilistic models*
466 *of proteins and nucleic acids*. Cambridge university press.
- 467 Farooq, B., Bierlaire, M., Hurtubia, R., Flötteröd, G., 2013. Simulation based population synthesis. *Trans-*
468 *portation Research Part B: Methodological* 58, 243–263. doi:[http://dx.doi.org/10.1016/j.](http://dx.doi.org/10.1016/j.trb.2013.09.012)
469 [trb.2013.09.012](http://dx.doi.org/10.1016/j.trb.2013.09.012).
- 470 Joh, C.H., Arentze, T., Hofman, F., Timmermans, H., 2002. Activity pattern similarity: a multidimensional
471 sequence alignment method. *Transportation Research Part B: Methodological* 36, 385–403. doi:[http://](http://dx.doi.org/10.1016/S0191-2615(01)00009-1)
472 [dx.doi.org/10.1016/S0191-2615\(01\)00009-1](http://dx.doi.org/10.1016/S0191-2615(01)00009-1).
- 473 Joh, C.H., Timmermans, H., 2011. Applying sequence alignment methods to large activity-travel data sets:
474 Heuristic approach. *Transportation Research Record: Journal of the Transportation Research Board* ,
475 10–17doi:<http://dx.doi.org/10.3141/2231-02>.

- 476 Joh, C.H., Timmermans, H., Arentze, T., 2006. Measuring and predicting adaptation behavior in multi-
477 dimensional activity-travel patterns. *Transportmetrica* 2, 153–173. doi:http://dx.doi.org/10.
478 1080/18128600608685659.
- 479 Liu, F., Janssens, D., Cui, J., Wets, G., Cools, M., 2015. Characterizing activity sequences using profile
480 hidden markov models. *Expert Systems with Applications* 42, 5705–5722. doi:http://dx.doi.
481 org/10.1016/j.eswa.2015.02.057.
- 482 Mohammadian, A.K., Javanmardi, M., Zhang, Y., 2010. Synthetic household travel survey data simulation.
483 *Transportation Research Part C: Emerging Technologies* 18, 869–878. doi:http://dx.doi.org/
484 10.1016/j.trc.2010.02.007.
- 485 Pendyala, R., Goulias, K., 2002. Time use and activity perspectives in travel behavior research. *Transporta-*
486 *tion* 29, 1–4. doi:http://dx.doi.org/10.1023/A:1012909228433.
- 487 Rasouli, S., Cools, M., Kochan, B., Arentze, T., Bellemans, T., Janssens, D., Timmermans, H., 2012.
488 Uncertainty in forecasts of complex rule-based systems of travel demand: Comparative analysis of the
489 albatross/feathers model system, in: 13th International Conference on Travel Behaviour Research, Inter-
490 national Association for Travel Behaviour Research (IATBR).
- 491 Rasouli, S., Timmermans, H., 2014. Activity-based models of travel demand: promises, progress and
492 prospects. *International Journal of Urban Sciences* 18, 31–60. doi:http://dx.doi.org/10.
493 1080/12265934.2013.835118.
- 494 Saadi, I., Mustafa, A., Teller, J., Cools, M., 2016a. An integrated framework for forecasting travel behav-
495 ior using markov chain monte carlo simulation and profile hidden markov models, in: Proceedings of
496 the 95th Annual Meeting of the Transportation Research Board, Transportation Research Board of the
497 National Academies, Washington, D.C.
- 498 Saadi, I., Mustafa, A., Teller, J., Farooq, B., Cools, M., 2016b. Hidden markov model-based population syn-
499 thesis. *Transportation Research Part B: Methodological* 90. doi:http://dx.doi.org/10.1016/
500 j.trb.2016.04.007.
- 501 Spissu, E., Pinjari, A.R., Bhat, C.R., Pendyala, R.M., Axhausen, K.W., 2009. An analysis of weekly
502 out-of-home discretionary activity participation and time-use behavior. *Transportation* 36, 483–510.
503 doi:http://dx.doi.org/10.1007/s11116-009-9200-5.
- 504 Sun, L., Erath, A., 2015. A bayesian network approach for population synthesis. *Transportation Research*
505 *Part C: Emerging Technologies* 61, 49–62. doi:http://dx.doi.org/10.1016/j.trc.2015.
506 10.010.
- 507 Voas, D., Williamson, P., 2000. An evaluation of the combinatorial optimisation approach to the creation
508 of synthetic microdata. *International Journal of Population Geography* 6, 349–366. doi:http://dx.
509 doi.org/10.1002/1099-1220(200009/10)6:5<349::AID-IJPG196>3.0.CO;2-5.
- 510 Vovsha, P., Hicks, J.E., Paul, B.M., Livshits, V., Maneva, P., Jeon, K., 2015. New features of population syn-
511 thesis, in: Proceedings of the 94th Annual Meeting of the Transportation Research Board, Transportation
512 Research Board of the National Academies, Washington, D.C.
- 513 Williamson, P., Birkin, M., Rees, P.H., et al., 1998. The estimation of population microdata by using data
514 from small area statistics and samples of anonymised records. *Environment and Planning A* 30, 785–816.
515 doi:http://dx.doi.org/10.1068/a300785.

- 516 Wilson, W.C., 1998. Activity pattern analysis by means of sequence-alignment methods. *Environment and*
517 *Planning A* 30, 1017–1038. doi:<http://dx.doi.org/10.1068/a301017>.
- 518 Ye, X., Konduri, K.C., Pendyala, R.M., Sana, B., Waddell, P., 2009. Methodology to match distributions of
519 both household and person attributes in generation of synthetic populations, in: *Proceedings of the 88th*
520 *Annual Meeting of the Transportation Research Board, Transportation Research Board of the National*
521 *Academies, Washington, D.C.*