

Population Synthesis for a City in a Developing Country

Devika Babu^{1*}, Sreelakshmi Balan¹, Matha Venkata Lakshmi Ranga Anjaneyulu¹

¹ Department of Civil Engineering, National Institute of Technology Calicut, Calicut Mukkam Road, 673601, Kerala, India

* Corresponding author, e-mail: devikababu.cet@gmail.com

Received: 04 September 2019, Accepted: 19 November 2019, Published online: 30 March 2020

Abstract

Activity-based approach in transportation studies increased the demand for detailed disaggregate data. To overcome the tedious and costly data collection, synthesized data are widely used in activity based modeling practices in developed countries. In the case of developing economies, time, resource and monetary constraints hinder the development of new generation travel demand models. Population synthesis practices are the need of the hour in developing countries. Through this study, the authors proposed a method for population synthesis for a medium sized city, in a developing country. Details from 9901 households, collected through home-interview survey in 2011 formed the database for this work. The proposed synthesis procedure makes use of Monte Carlo process along with logit modeling technique to simulate population attributes using the survey data. The procedure is developed as a Visual Basic Application in spreadsheet platform. By adopting this method, various household level and person level attributes are simulated. Comparison of observed data from survey and simulated data showed consistent results and low differences. It is expected that this study will assist the planning authorities to better understand the city's population characteristics.

Keywords

population synthesis, Monte Carlo simulation, logit model, developing country

1 Introduction

Activity-based approach in travel demand modeling originated from the recognition that individuals travel to perform various activities, which are distributed spatially and temporally (Bhat and Koppelman, 2003). This approach received more attention in the travel behavior research arena as many limitations of conventional trip based models are overcome by the behaviorally oriented activity-based models. These models are primarily based on detailed household socio-demographics and activity-travel characteristics to model individual's travel behavior, rather than the traditional aggregate approach. More realistic demonstration of individual travel choices is possible by means of this approach (Davidson et al., 2007).

The growing interests towards developing disaggregate travel-demand models, considerably increased the need for more detailed disaggregate level data (Timmermans, 2005). This intricate data requirement is one of the major challenges in transportation research, especially for formulating sophisticated micro-level models. However, good quality disaggregate level input data are very important in transportation modeling. It is essentially a tedious, costly and cumbersome task to gather the required data from all

respondents in a population. Hence, transportation studies rely upon data from a sample to draw inferences about the whole population. Sampling a population is advantageous in terms of economy, speed, timeliness, feasibility quality and accuracy (Travel survey manual, 1996).

In order to overcome tedious and costly data collection, activity based modeling practices make use of synthesized/simulated data of the population. For the development of disaggregate models, micro-level data are essential. Population synthesis techniques are a viable substitute for this kind of data collection. The relevance of synthetic population generation for modeling applications is increasing, as the demand for more disaggregate level models is increasing now-a-days.

Müller and Axhausen (2010) presented a review of six population synthesis methods which are in use for different micro simulation models. They are PopSynWin, ILUTE, PopGen, FSUTMS, CEMDAP and ALBATROSS. A comprehensive multimodal activity-based system developed for Florida is Florida Activity Mobility Simulator (FAMOS) (Pendyala et al., 2005). The two main modules of FAMOS are Household Attributes Generation System

and the Prism-Constrained Activity-Travel Simulator. These studies primarily make use of the previous year census data for generating the base year target population using Iterative Proportional Fitting (IPF) method. IPF is an iterative algorithm for estimating cell values of a contingency table (i.e., the complete distribution across all control variables) such that the known marginal totals remain fixed. For using IPF method, accurate population totals and joint distributions (for key variables) from census data or other authorized sources are essential.

Guo and Bhat (2007) discussed the two major issues associated with the conventional Iterative Proportional Fitting (IPF) procedure for synthesizing the base-year population. One of them is zero-cell value problem. The other one is that the IPF procedure is not capable to simultaneously control household level and person level attributes. They proposed a new population synthesis procedure which can generate more accurate true population than the conventional approach can.

Mueller and Axhausen (2011) included Iterative Proportional Updating (IPU) in their study to overcome the inability of IPF to simultaneously control for household and person level attributes. They proposed a new procedure which brings in heterogeneity in the synthetic population and applied it for synthesizing the population of Switzerland.

Mohammadian et al. (2010) used a two stage process to develop synthetic population and simulate household level travel survey data for some areas of the State of New York where actual survey is not possible. Ma and Srinivasan (2012) empirically assessed target year populations which are generated using different base year populations and data-fusion methods for Florida. It is observed from their study that the precision of generated target year population is significantly affected by the inaccuracy of base year population.

Most of these studies focused on replicating the population for developed countries, whereas, population synthesis practices in the context of developing countries are very few. Unlike developed economies, there is no dedicated agency to conduct nation-wide household travel survey in India, to the best of authors' knowledge. Davidson et al. (2007) report funding as one of the critical factor for organizations to build up travel demand forecast models. Expensive home-interview surveys, developing models and software are major budget constraints. This is particularly significant in third world nations, where, the lack of professional expertise with in the regional /town planning

agencies, along with time, resource and monetary constraints hinder the development, experimentation and use of new generation travel demand models.

Despite of all these issues, data collection and model development are clearly inevitable in transportation studies, as these form the foundation for devising travel demand management strategies to provide better transportation facilities to the society. This paper is an attempt to simulate the population details of a medium sized city in the context of a developing country. The synthesis of population characteristics is performed primarily with the help of Monte Carlo simulation technique in spreadsheet platform. Logit modeling technique is also adopted to simulate certain person level attributes. Care has been taken to make the synthesis procedure simple, so that the simulation of population is possible with basic expertise.

2 Description of database

The following sections describe the study area, survey administration and analysis on the collected data. For carrying out this research, Calicut city in Kerala State, India is chosen as the study area. It is a medium sized city spread over 118.59 square kilometers. It has a large number of business and commercial establishments and is popular for education facilities. As per Census, the total population of the city is around 0.61 million (adapted from Master plan for Kozhikode urban area-2035 (draft)).

2.1 Activity-travel survey

The data used as the base for population synthesis in this study was collected by home-interview survey using an activity-travel diary. Trained enumerators visited households within the limits of Calicut city, on a non-working day and the household details, personal details and information related to activities and travel were collected by face-to-face interview technique. The survey was organized by National Institute of Technology Calicut between December 2010 and February 2011. Around 150-200 sample households were interviewed from each of the 75 traffic analysis zones, so as to ensure the geographical distribution of samples. Intimation regarding the data collection was done through the leading newspapers prior to the survey, which enhanced the response rate from the people. The contact information from the household was also collected for future clarifications in the responses. The collected data were entered manually in spreadsheets. The data were thoroughly checked for typos and other inconsistencies before analysis. Incorrect and incomplete samples were deleted

and the final dataset consisted of details from 9901 households and activity-travel details of 39637 persons.

2.2 Data analysis

The computerized data was analyzed for getting more insight into the socio-demographic details and travel characteristics of the study area. On an average, 4 members are present in a household. In majority of the households (about 90 %) at least one employed person is present and in 57 percent of households, at least one student is present. 57 percent of households own one or more automobiles. 51 percent of sample population are females and 49 percent are males (as per Census 2011 data, 52 percent are females and 48 percent are males). Working people constitute 34 percent in the sample. The purposes for which travel performed are grouped into work, education, religious, recreation, shopping, medical and other purposes. Work and education trips together constitute 68 percent of the total trips. Majority of the trips, 39 percent, are made by bus and 24 percent by two-wheeler. Share of car trips is observed to be less compared to bus and two-wheeler trips.

3 Methodology

The IPF method is widely used by modelers for synthesizing population characteristics. One issue regarding the use of this method in the context of a developing country is the difficulty in gathering precise marginal distributions of control variables for the year of study. According to Müller and Axhausen (2010), use of census data as an input to microsimulation models is limited due to its unavailability and removal of several required information. Even in the case of developed countries, for privacy reasons, many national censuses remove spatial details and the linkages between households and persons, while providing data (Pritchard and Miller, 2012). There are few exceptions to this kind of practice. Farooq et al. (2013) reported that the complete census data are available for research in Switzerland.

The lack of joint distribution of control variables for the year of study motivated to find out alternate method for simulating the population characteristics. In this study, most of the household and person level attributes are synthesized using Monte Carlo simulation process. The empirical distributions of population characteristics are identified from the survey data and inverse transformation method is used to generate random variates. The following description of the method is taken from Perros (2009). Assume that, it is required to generate stochastic

variates from a probability density function $f(x)$. Let $F(x)$ be the cumulative density function, which is defined in the region $[0,1]$. First a random number 'r' is generated and is set equal to $F(x)$. That is, $F(x) = r$. Then by inverting F , the value of 'x' is obtained. That is, $x = F^{-1}(r)$, where $F^{-1}(r)$ indicates the inverse transformation of F .

This procedure is developed as macros using Visual Basic for Application (VBA) in spreadsheet platform. The illustration of this method in replicating the survey data (9901 households' details) is provided in the following paragraphs. Finally, the details of entire population of Calicut city are synthesized using this method.

As the first step, the probability values of different household size categories were calculated from the survey data. Table 1 gives the probability and cumulative probability values of household sizes for the sample data.

Next, 1000 uniformly distributed random numbers were generated. Let the random number corresponding to the first household be 0.705548. This value is compared with the computed cumulative probability values. From Table 1, it can be observed that this random number falls between 0.693263 and 0.860317 which is between household size 4 and 5. Hence the first household will be assigned with size equal to 5. In this manner, size of each household is determined.

4 Synthesis of household characteristics

This section describes the synthesis of household level control variables based on the information from the collected data. Guo and Bhat (2007) selected household size, household age, presence of children under 18 years and type of

Table 1 Probability values of household size categories

Household size	Probability	Cumulative probability
1	0.020503	0.020503
2	0.152914	0.173417
3	0.217352	0.390769
4	0.302495	0.693263
5	0.167054	0.860317
6	0.069993	0.930310
7	0.029795	0.960105
8	0.020503	0.980608
9	0.009191	0.989799
10	0.008282	0.998081
11	0.001515	0.999596
12	0.000303	0.999899
13	0.000101	1

family as some of the control variables. The control variables adopted for household synthesis by Mohammadian et al. (2010) are size of the household, income of household, number of workers, presence of children, age of household member, age group, job category, number of vehicles, education level, land use information and ethnicity. For the present study, the key variables identified for household synthesis are household size, number of males, number of females and income.

Household size was simulated for each household in the population following the method described in Section 3. For simulating the number of males and females in the household, conditional probability values of the number of males and females for different household sizes are computed from the survey data. Cumulative probability values of different combinations are also computed (e.g., for household size 5, different combinations can be 0 males and 5 females; 1 male and 4 females and so on). The conditional probability values are given in Table 2.

A different set of random numbers equal to the total number of synthesized households are generated and each random number is compared with the cumulative probability values of different male-female combinations, conditional on the household size. The values of number of males and females are thus assigned to the synthesized household accordingly. Each of the generated household is assigned with a unique Identification Number (ID).

Table 2 Conditional probability values of number of males and females

Household size	No. of males	No. of females	Frequency	Conditional probability	Cumulative probability
1	0	1	103	0.507389163	0.507389163
1	1	0	100	0.492610837	1
2	0	2	109	0.071994716	0.071994716
2	1	1	1364	0.900924703	0.972919419
2	2	0	41	0.027080581	1
3	0	3	49	0.022769517	0.022769517
3	1	2	1034	0.480483271	0.503252788
3	2	1	1045	0.485594796	0.988847584
3	3	0	24	0.011152416	1
4	0	4	21	0.007011686	0.007011686
4	1	3	689	0.230050083	0.23706177
4	2	2	1632	0.544908180	0.78196995
4	3	1	645	0.215358932	0.997328881
4	4	0	8	0.002671119	1
5	0	5	6	0.00362757	0.00362757
		...			
		(and so on)			

5 Synthesis of person characteristics

Along with the synthesis of household characteristics, person level attributes are also generated. The control variables identified for person synthesis are gender (male/female), age and age group, occupation, occupation category, driving license status and education level. Under occupation, individuals are assigned with an occupation status (worker/ student/ home-maker/ retired person/ infant/ other). Further, a worker is categorized as government employee/ working abroad/ private employee/ daily wage person/ marketing professional/ self-employed person. Similarly, a student is categorized as pre-school/ school/ plus-two/ college student. 'Other' category includes persons seeking employment or not-working. Eleven education levels are considered for synthesis starting from no education, kindergarten and so on up to doctoral level.

For each synthetic household, number of persons is generated based on household size. Each person is assigned with a unique ID. Person details are generated in second spreadsheet. Gender is assigned on the basis of number of males and females' data from the household synthesis spreadsheet to every person. Age is simulated from the conditional probability values of age conditional on gender. Based on the generated value of age, age group is assigned. The age groups considered are less than or equal to 3 years (infants), 4 to 15 years (school going children), 16 to 17 years (plus two students), 18 to 23 years (graduate students), 24 to 40 years (young working group), 41 to 56 years (middle aged working group), 57 to 60 years (retired persons) and greater than 60 years (elderly persons). Based on the conditional probability of occupation with respect to gender and age group, which is calculated from the survey data, occupation is synthesized for each individual. Following this, education level of each person is simulated by comparing the conditional probability value of education level conditional on occupation and age group.

In order to simulate the driving license holding status, a model is developed with dependent variable as license status of a person, that is whether the person holds a driving license or not. As the dependent variable is dichotomous (yes/no), binary logit modeling is adopted. Usually driver's license holding will not be enquired as part of Census data collection in India. Hence, it is decided to develop a model for the same, so that this model can be used for predicting the driving license status. A person's working status and education level are found to be the most influencing variables for holding a driving license. The model summary is provided in Table 3. The goodness of fit measures are also provided in the table.

Table 3 Model for driving license holding status

Variables	Coeff.	Sig.	Exp (B)	95% C.I. for Exp (B)	
				Lower	Upper
Occupation: Worker	2.417	0.000	11.209	10.519	11.945
Education: Degree	1.718	0.000	5.575	4.936	6.296
Education: Diploma	1.830	0.000	6.233	3.602	10.785
Education: High school	0.216	0.016	1.241	1.041	1.478
Education: Higher secondary	0.825	0.000	2.282	1.69	3.082
Education: Post graduation	1.395	0.000	4.035	3.547	4.59
Education: Doctorate	1.432	0.082	4.189	0.834	21.037
Education: Matriculation	0.837	0.000	2.309	2.045	2.606
Constant	-3.437	0.000	0.032		
Goodness of fit measures					
-2 Log likelihood				28021.893	
Cox & Snell R Square				0.247	
Nagelkerke R Square				0.357	
Chi-square				133.576	
Overall Percentage correctly predicted				77.20%	

It is observed that a working person is more likely to possess driving license, compared to a non-worker. The chance of owning driver's license increases with the increase in education level. The probable reason can be workers are likely to own personalized vehicles (and hence driving license) due to their work timings and job responsibilities. Similarly, as the education level of a person increases, due to the wish to improve quality of life, people may prefer to own private vehicles and hence may tend towards having driving license. The overall goodness-of-fit as measured by Nagelkerke R² value is 0.357, which shows reasonably good fit for the model (Hensher et al., 2005) and the model predictability is 77.2 %. Hence, this model is applied to simulate a person's driving license holding status, if the person's age is 18 years or above.

In addition to the simulated household level characteristics, some more household descriptors are derived from the synthesized data of individuals. They are number of employed persons, number of graduates, number of male adults, number of students below 17 years, number of persons above 60 years of age and number of employed females.

6 The simulation platform in spreadsheets

The entire simulation procedure is developed in Visual basic for Application (VBA) platform. As an example, the screenshots of the programs simulation of the size of household and number of males and females in the house are given in Fig. 1 and Fig. 2 respectively. In a similar manner, all the synthesis procedures are written as macros and control buttons are assigned for executing the simulation run.

7 Validation of household and person synthesis

Validation of the synthesized household and person level attributes were carried out by comparing the respective attributes with the observed data. The observed and simulated values of household sizes are shown in Fig. 3. Reasonably good match was seen between the observed and simulated percentages of households, for different household sizes. The maximum difference between observed and simulated data is found to be 0.81 (absolute value) percent. Hence the program is observed to accurately simulate the attribute, household size.

Fig. 4 and Fig. 5 present the observed and simulated values of the count of males and females per household respectively. In the attribute 'number of males per

House -hold ID	House -hold size
1	5
2	4
3	4
4	3
5	3
6	5
7	1
8	5
9	5
10	5

Fig. 1 Screenshot of household size simulation function and simulated household sizes

House -hold ID	No. of males	No. of females
1	3	2
2	3	1
3	3	1
4	1	2
5	2	1
6	3	2
7	0	1
8	2	3
9	2	3
10	4	1

Fig. 2 Screenshot of number of males and females simulation function and simulated values

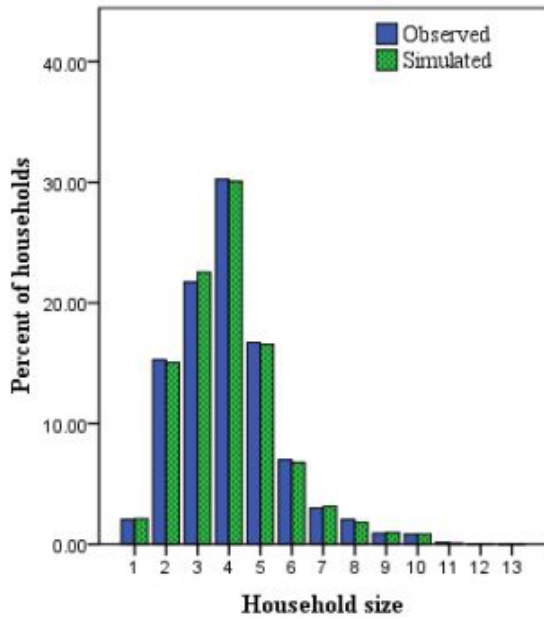


Fig. 3 Observed and simulated values of household size

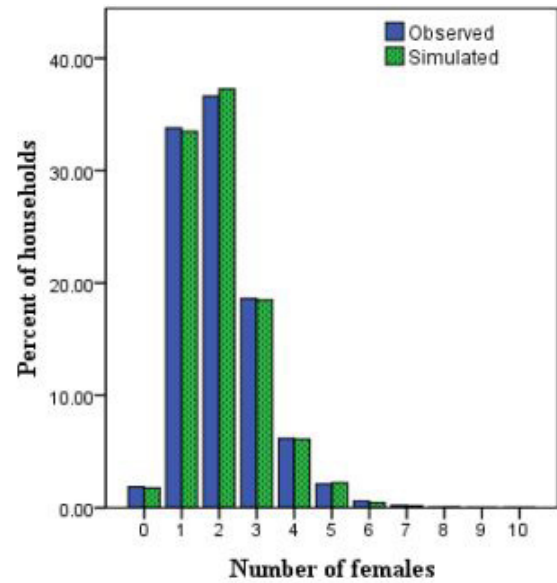


Fig. 5 Observed and simulated values of number of females per household

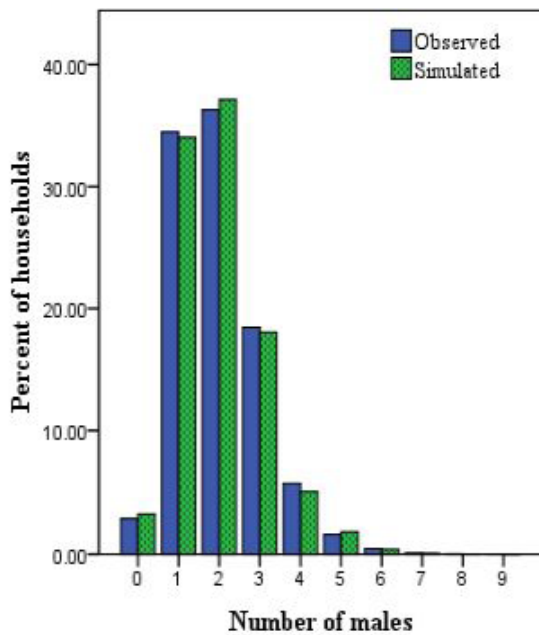


Fig. 4 Observed and simulated values of number of males per household

household’, the maximum difference is observed to be 0.86 (absolute value in percent) for two males per household category. The minimum difference is 0.01 for both 8 and 9 males per household categories.

In the case of number of females per household, the maximum and minimum absolute values of variation are found to be 0.63 and 0.01 respectively. The maximum difference is seen for two females per household category.

7.1 Root Mean Square Error (RMSE) values

Root Mean Square Error (RMSE) values were calculated for checking the goodness of fit between observed and simulated data. The RMSE values computed for various household and person level attributes are given in Table 4. It can be observed from the table that the values are very low.

7.2 Simulation of Calicut city population

Using the above described procedure, the entire population of Calicut city is simulated. Details of 0.15 million households and 0.6 million persons are simulated using

Table 4 Probability values of household size categories

Sl. No.	Attribute	RMSE
1	Household size	0.0026
2	Number of males per household	0.0040
3	Number of females per household	0.0022
4	Number of employed persons per household	0.0628
5	Number of graduates per household	0.0417
6	Number of license holders per household	0.0466
7	Number of two-wheelers per household	0.0212
8	Number of cars per household	0.0449
9	Gender	0.0018
10	Age group	0.0010
11	Education level	0.0031
12	Occupation	0.0018
13	Type of employment	0.0025
14	Driving license holding status	0.0596

the proposed method. The simulated data and Census data totals at the aggregate level are provided in Table 5.

The synthesized data closely matches with the actual population totals. This substantiates the suitability of the proposed method in simulating the population characteristics.

8 Summary and conclusions

This paper aimed to synthesize the population characteristics of a medium sized city, in the context of a developing country. Using Monte Carlo simulation technique and inverse transformation method, household and person level attributes are simulated for the study area. Logit modeling framework is also adopted to model certain attributes. Household, personal and activity-travel details for population synthesis, attributes such as household size, number of males and females, individual's age, gender, education and occupation were synthesized.

This work contributes to the existing body of literature by two fold. First is by synthesizing the population characteristics of a city in the context of a developing country. Another contribution is that the synthesis process is performed in a spreadsheet platform. As the procedure is simple and developed as visual basic application, this can

Table 5 Simulated data of Calicut city

	Simulated data	Census data
Total Population	600629	613255
Total number of males	295286	293003
Total number of females	305343	320252
Male/Female ratio	0.967	0.916

be adopted for simulating the population for regional/local planning agency.

Comparison between observed and simulated data showed very small difference. Efforts are being made to extend the capabilities of this simulator to generate the activity-travel characteristics of individuals. This tool is expected to enable the authorities to understand the city's travel behavior in a better manner and help them to formulate better policy strategies. This will gradually result in improving the quality-of-life of people, particularly in the context of developing countries.

Acknowledgment

Support for this work from Centre for Transportation Research (a Centre of Excellence setup under FAST Scheme of MHRD, Government of India) is greatly acknowledged.

References

- Bhat, C. R., Koppelman, F. S. (2003) "Activity-Based Modeling of Travel Demand", In: Hall, R.W. (ed.) Handbook of Transportation Science. International Series in Operations Research & Management Science, Vol. 56, Springer, Boston, MA, USA, pp. 39–65.
https://doi.org/10.1007/0-306-48058-1_3
- Davidson, W., Donnelly, R., Vovsha, P., Freedman, J., Ruegg, S., Hicks, J., Castiglione, J., Picado, R. (2007) "Synthesis of first practices and operational research approaches in activity-based travel demand modeling", Transportation Research Part A: Policy and Practice, 41(5), pp. 464–488.
<https://doi.org/10.1016/j.tra.2006.09.003>
- Farooq, B., Bierlaire, M., Hurtubia, R., Flötteröd, G. (2013) "Simulation based population synthesis", Transportation Research Part B: Methodological, 58, pp. 243–263.
<https://doi.org/10.1016/j.trb.2013.09.012>
- Guo, J. Y., Bhat, C. R. (2007) "Population Synthesis for Microsimulating Travel Behavior", Transportation Research Record: Journal of the Transportation Research Board, 2014(1), pp. 92–101.
<https://doi.org/10.3141/2014-12>
- Hensher, D. A., Rose, J. M., Greene, W. H. (2005) "10 - Getting started modeling: the basic MNL model", In: Applied Choice Analysis: A Primer, Cambridge University Press, Cambridge, UK, pp. 308–373.
<https://doi.org/10.1017/CBO9780511610356.011>
- Ma, L., Srinivasan, S. (2012) "Synthesizing target year populations for input to travel demand models", In: 2012 Transportation Research Board Annual Meeting, Washington, DC, USA, Article number: 12-4355.
- Mohammadian, A. K., Javanmardi, M., Zhang, Y. (2010) "Synthetic household travel survey data simulation", Transportation Research Part C: Emerging Technologies, 18(6), pp. 869–878.
<https://doi.org/10.1016/j.trc.2010.02.007>
- Mueller, K., Axhausen, K. W. (2011) "Occam's Razor and some randomness: generating a synthetic population for Switzerland", In: 2011 European Transport Conference, Glasgow, Scotland, Article number: 01470930.
- Müller, K., Axhausen, K.W. (2010) "Population synthesis for microsimulation: State of the art", ETH Zürich, Institut für Verkehrsplanung, Transporttechnik, Strassen- und Eisenbahnbau (IVT), 638, pp. 1–12.
<https://doi.org/10.3929/ethz-a-006127782>
- Kozhikode Municipal Corporation "Master Plan For Kozhikode Urban Area-2035 (Draft)", [online] Available at: https://kozhikodecorporation.lsgkerala.gov.in/en/master_plan_news [Accessed: 07 March 2017]
- Pendyala, R. M., Kitamura, R., Kikuchi, A., Yamamoto, T., Fujii, S. (2005) "Florida Activity Mobility Simulator: Overview and Preliminary Validation Results", Transportation Research Record: Journal of the Transportation Research Board, 1921(1), pp. 123–130.
<https://doi.org/10.1177/0361198105192100114>

- Perros, H. G. (2009) "The generation of stochastic variates", In: Perros, H. Computer Simulation Techniques: The definitive introduction!, NC State University, Raleigh, NC, USA, pp. 47–49.
- Pritchard, D. R., Miller, E. J. (2012) "Advances in population synthesis: fitting many attributes per agent and fitting to household and person margins simultaneously", *Transportation*, 39(3), pp. 685–704. <https://doi.org/10.1007/s11116-011-9367-4>
- Systematics, C. (1996) "Travel Survey Manual", Prepared for US DOT and US EPA. TMIP Program, Washington, DC, USA, Rep. FHWA-PL-96-029. <https://doi.org/10.21949/1404543>
- Timmermans, H. J. P. (2005) "Activity-based approaches: models, data and applications", *Activity-Based Approaches: Progresses in Activity-Based Analysis*, pp. 19–26. <https://doi.org/10.1016/B978-008044581-6/50003-X>