



HAL
open science

Generating a synthetic population of individuals in households: Sample-free vs sample-based methods

Maxime Lenormand, Guillaume Deffuant

► **To cite this version:**

Maxime Lenormand, Guillaume Deffuant. Generating a synthetic population of individuals in households: Sample-free vs sample-based methods. 2012. hal-00725531v1

HAL Id: hal-00725531

<https://hal.science/hal-00725531v1>

Preprint submitted on 27 Aug 2012 (v1), last revised 8 May 2015 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Generating a synthetic population of individuals in households: Sample-free vs sample-based methods

Maxime Lenormand¹ and Guillaume Deffuant¹

IRSTEA, LISC, 24 avenue des Landais, 63172 AUBIERE, France
(maxime.lenormand, guillaume.deffuant)@irstea.fr

Abstract. We compare a sample-free method proposed by [Gargiulo et al. \(2010\)](#) and a sample-based method proposed by [Ye et al. \(2009\)](#) for generating a synthetic population, organised in households, from various statistics. We generate a reference population for a French region including 1310 municipalities and measure how both methods approximate it from a set of statistics derived from this reference population. We also perform sensitivity analysis. The sample-free method better fits the reference distributions of both individuals and households. It is also less data demanding but it requires more pre-processing. The quality of the results for the sample-based method is highly dependent on the quality of the initial sample.

1 Introduction

For two decades, the number of micro-simulation models, simulating the evolution of large populations with an explicit representation of each individual, has been constantly increasing with the computing capabilities and the availability of longitudinal data. When implementing such an approach, the first problem is initialising properly a large number of individuals with the adequate attributes. Indeed, in most of the cases, for privacy reasons, exhaustive individual data are excluded from the public domain. Aggregated data at various levels (municipality, county,...), guaranteeing this privacy, are hence only available in general. Sometimes, individual data are available on a sample of the population, these data being chosen also for guaranteeing the privacy (for instance omitting the individual's location of residence). This paper focuses on the problem of generating a virtual population with the best use of these data, especially when the goal is generating both individuals and their organisation in households.

Two main methods, both requiring a sample of the population, aim at tackling this problem:

- The synthetic reconstruction method (SR) ([Wilson and Pownall, 1976](#)). These methods generally use the Iterative Proportional Fitting ([Deming and Stephan, 1940](#)) and a sample of the target population to obtain the joint-distributions of interest ([Beckman et al., 1996](#); [Huang and Williamson,](#)

2002; Guo and Bhat, 2007; Arentze et al., 2007; Ye et al., 2009). Many of the SR methods match the observed and simulated households joint-distribution or individual joint-distribution but not simultaneously. To circumvent these limitations Guo and Bhat (2007); Arentze et al. (2007); Ye et al. (2009) proposed different techniques to match both household and individual attributes. Here, we focus on the Iterative Proportional Updating developed by Ye et al. (2009).

- The combinatorial optimization (CO). These methods create a synthetic population by zone using marginals of the attributes of interest and a sub-set of a sample of the target population for each zone (for a complete description see Voas and Williamson (2000); Huang and Williamson (2002)).

Recently, sample-free SR methods appeared (Gargiulo et al., 2010; Barthelemy and Toint, 2012). These methods can be used in the usual situations where no sample is available and one must only use distributions of attributes (of individuals and households). Hence, they overcome a strong limit of the previous methods. It is therefore important to assess if this larger scope of the sample-free method implies a loss of accuracy compared with the sample-based method.

The aim of this paper is contributing to this assessment. With this aim, we compare the sample-based IPU method proposed by Ye et al. (2009) with the sample-free approach proposed by Gargiulo et al. (2010) on an example.

In order to compare the methods, the ideal case would be to have a population with complete data available about individuals and households. It would allow us to measure precisely the accuracy of each method, in different conditions. Unfortunately, we do not have such data. In order to put ourselves in a similar situation, we generate a virtual population and then use it as a reference to compare the selected methods as in Barthelemy and Toint (2012).

In the first section we formally present the two methods. In the second section we present the comparison results. Finally, we discuss our results.

2 Details of the chosen methods

2.1 Sample-free method Gargiulo et al. (2010)

We consider a set of n individuals X to dispatch in a set of m households Y in order to obtain a set of filled households P . Each individual x is characterised by a type t_x from a set of q different individual types T (attributes of the individual). Each household y is characterised by a type u_y from a set of p different household types U (attributes of the household). We define $n_T = \{n_{t_k}\}_{1 \leq k \leq q}$ as the number of individuals of each type and $n_U = \{n_{u_l}\}_{1 \leq l \leq p}$ as the number of households of each type. Each household y of a given type u_y has a probability to be filled by a subset of individuals L , then the content of the household equals L , which is denoted $c(y) = L$. We use this probability to iteratively fill the households with the individuals of X .

$$\mathbb{P}(c(y) = L | u_y) \tag{1}$$

The iterative algorithm used to dispatch the individuals into the households according to the Equation 1 is described in Algorithm 1. The algorithm starts with the list of individuals X and of the households Y , defined by their types. Then it iteratively picks at random a household, and from its type and Equation 1, derives a list of individual types. If this list of individual types is available in the current list of individuals X , then this filled household is added to the result, and the current lists of individuals and households are updated. This operation is repeated until one of the lists X or Y is void, or a limit number of iterations is reached.

Algorithm 1 The general iterative algorithm

Input : X and Y

Output : P

Set $P = \emptyset$

while $Y \neq \emptyset$ **do**

 Pick at random y from Y

 Pick at random L with a probability defined in Equation 1

if $L \subset X$ **then**

$P \leftarrow P \cup L$

$Y \leftarrow Y \setminus \{y\}$

$X \leftarrow X \setminus L$

end if

end while

In the case of the generation of a synthetic population, we can replace the selection of the list L by the selection of the individuals one at a time by order of importance in the household. In this case Equation 2 replaces Equation 1.

$$\begin{aligned}
 & \mathbb{P}(x_1 \in y | u_y) \times \\
 & \mathbb{P}(x_2 \in y | u_y, x_1 \in y) \times \\
 & \mathbb{P}(x_3 \in y | u_y, x_1 \in y, x_2 \in y) \times \\
 & \dots
 \end{aligned} \tag{2}$$

The iterative approach algorithm associated with this probability is described in Algorithm 2. The principle is the same as previously, it is simply quicker. Instead of generating the whole list of individuals in the household before checking it, one generates this list one by one, and as soon as one of its member cannot be found in X , the iteration stops, and one tries another household.

In practice this stochastic approach is data driven. Indeed, the types T and U are defined in accordance with the data available and the complexity to extract the distribution of the Equation 2 increases with n_T and n_U . The distributions defined in Equation 2 are called distributions for affecting individual into household. In concrete applications, it occurs that one needs to estimate n_T , n_U and the distributions of probabilities presented in Equation 2. This estimation implies that the Algorithm 2 can not converge in a reasonable time because of the

Algorithm 2 The iterative algorithm

Input : X and Y **Output :** P Set $P = \emptyset$ **while** $Y \neq \emptyset$ **do** Pick at random y from Y Pick at random x_1 with a probability $\mathbb{P}(x_1 \in y|u_y)$ Pick at random x_2 with a probability $\mathbb{P}(x_2 \in y|u_y, x_1 \in y)$ Pick at random x_3 with a probability $\mathbb{P}(x_3 \in y|u_y, x_1 \in y, x_2 \in y)$

...

if $\{x_1, x_2, x_3, \dots\} \subset X$ **then** $P \leftarrow P \cup \{x_1, x_2, x_3, \dots\}$ $Y \leftarrow Y \setminus \{y\}$ $X \leftarrow X \setminus \{x_1, x_2, x_3, \dots\}$ **end if****end while**

stopping criterion ($Y \neq \emptyset$). This stopping criterion is equivalent to an infinite number of "filling" trials by households. In this case, we can replace the stopping criterion by a maximal number of iterations by households and then put the remaining individuals in the remaining households using relieved distributions for affecting individual into household.

In a perfect case where all the data are available and the time infinite, the algorithm would find a perfect solution. When the data are partial and the time constrained, it is interesting to assess how this method manages to make the best use of the available data.

3 The sample-based approach (General Iterative Proportional Updating)

In this approach, proposed by [Ye et al. \(2009\)](#), starts with a sample P_s of P and the purpose is to define a weight w_i associated with each individual and each household of the sample in order to match the total number of each type of individuals in X and households in Y to reconstruct P . The method used to reach this objective is the Iterative Proportional Updating (IPU). The algorithm proposed in [Ye et al. \(2009\)](#) is described in Algorithm 3. In this algorithm, for each type of households or individuals j the purpose is to match the weighted sum sw_j with the estimated constraints e_j with an adjustment of the weights. w_j is the weight of household or individual j in the weighted sample and e_j is an estimation of the total number of households or individuals j in P . This estimation is done separately for each individual and household type using a standard IPF procedure with marginal variables. When the match between the weighted sample and the constraint become stable, the algorithm stops. The procedure then generates a synthetic population by drawing at random the filled households of P_s with probabilities corresponding to the weights. This generation

is repeated several times and one chooses the result with the best fit with the observed data.

Algorithm 3 Iterative Proportional Updating algorithm

Input : P_s, ϵ

Output : P

Set $P = \emptyset$

Generate $D \in M_{|P_s| \times (p+q)}(\mathbb{R})$ described by the light grey table in Table 1

Estimate n_T and n_U using the standard IPF procedure and store the resulting estimate into a vector $E = (e_j)_{1 \leq j \leq p+q}$ as in Table 1

for $i = 1$ to $|P_s|$ **do**

 Set $w_i = 1$

end for

for $j = 1$ to $p + q$ **do**

 Compute $sw_j = \sum_{i=1}^{|P_s|} d_{ij}w_i$

 Compute $\delta_j = \frac{|sw_j - e_j|}{e_j}$

end for

Compute $\delta = \frac{1}{p+q} \sum_{j=1}^{p+q} \delta_j$

Set $\delta_{\min} = \delta$

Set $\Delta = \epsilon + 1$

while $\Delta > \epsilon$ **do**

 Set $\delta_{\text{prev}} = \delta$

for $j = 1$ to $p + q$ **do**

for $i = 1$ to $|P_s|$ **do**

if $d_{ij} \neq 0$ **then**

$w_i = \frac{e_j}{sw_j} w_i$

end if

end for

 Compute $sw_j = \sum_{i=1}^{|P_s|} d_{ij}w_i$

end for

 Compute $\delta = \frac{1}{p+q} \sum_{j=1}^{p+q} \delta_j$

if $\delta < \delta_{\min}$ **then**

 Set $W_{\text{opt}} = (w_i)_{1 \leq i \leq |P_s|}$

$\delta = \delta_{\min}$

end if

$\Delta = |\delta - \delta_{\text{prev}}|$

end while

4 Generating a synthetic population of reference for the comparison

Because we cannot access any population with complete data available about individuals and households, we generate a virtual population and then use it as a reference to compare the selected methods as in [Barthelemy and Toint \(2012\)](#).

Table 1: IPU Table. The light grey table represents the frequency matrix D showing the household (HH) type U and the frequency of different individual (Ind.) types T within each filled households for the sample P_s . The dimension of D is $|P_s| \times (p + q)$, where $|P_s|$ is the cardinal number of the sample P_s , q the number of individual types and p the number of household types. An element d_{ij} of D represents the contribution of filled household i to the frequency of individual/household type j .

Filled HH ID	HH Type u_1	...	HH Type u_p	Ind. Type t_1	...	Ind. Type t_q	Weight
1	d_{11}	...	d_{1q}	d_{1q+1}	...	d_{1q+p}	w_1
...
$ P_s $	$d_{ P_s 1}$...	$d_{ P_s q}$	$d_{ P_s q+1}$...	$d_{ P_s q+p}$	$w_{ P_s }$
WS	ws_1	...	ws_p	ws_{p+1}	...	ws_{p+q}	
E	$e_1 = \hat{n}_{u_1}$...	$e_p = \hat{n}_{u_p}$	$e_{p+1} = \hat{n}_{t_1}$...	$e_{p+q} = \hat{n}_{t_q}$	
δ	δ_1	...	δ_p	δ_{p+1}	...	δ_{p+q}	

We start with statistics about the population of Auvergne (French region) in 1990 using the free-sample approach presented above. The Auvergne region is composed of 1310 municipalities, 1,321,719 inhabitants gathered in 515,736 households. In average the municipalities had about 1000 inhabitants with a minimum of 25 and a maximum of 136,180.

4.1 Generation of the individuals

For each municipality of the Auvergne region we generate a set X of individuals with a stochastic procedure. For each individual of the age pyramid (distribution 1 in Table 2), we randomly choose an age in the bin and then we draw randomly an activity status according to the distribution 2 in Table 2.

4.2 Generation of the households

For each municipality of the Auvergne region we generate a set Y of households according to the total number of individual $n = |X|$ with a stochastic procedure. We draw at random households according to the distribution 3 in Table 2 while the sum of the capacities is below n and then we determine the last household to have n equal to the sum of the size of the households.

4.3 Distributions for affecting individual into household

Single

- The age of the individual 1 is determined using the distribution 4 in Table 2.

Monoparental

- The age of the individual 1 is determined using the distribution 4 in Table 2.
- The ages of the children are determined according to the age of individual 1 (An individual can do a child after 15 and before 55) and the distribution 6 in Table 2.

Couple without child

- The age of the individual 1 is determined using the distribution 4 in Table 2.
- The age of the individual 2 is determined using the distribution 5 in Table 2.

Couple with child

- The age of the individual 1 is determined using the distribution 4 in Table 2.
- The age of the individual 2 is determined using the distribution 5 in Table 2.
- The ages of the children are determined according to the age of individual 1 and the distribution 6 in Table 2.

Other

- The age of the individual 1 is determined using the distribution 4 in Table 2.
- The ages of the others individuals are determined according to the age of individual 1.

Table 2: Data description

ID	Description	Level
1	Number of individuals grouped by ages	Municipality (LAU2)
2	Distribution of individual by activity statut according to the age	Municipality (LAU2)
3	Joint-distribution of household by type and size	Municipality (LAU2)
4	Probability to be the head of household according to the age and the type of household	Municipality (LAU2)
5	Probability of having a couple according to the difference of age between the partners (from"-16years" to "21years")	National level
6	Probability to be a child (child=live with parent) of household according to the age and the type of household	Municipality (LAU2)

To obtain a synthetic population P with households Y filled by individuals X we use the Algorithm 2 where we approximate the Equation 2 with the distributions 4, 5 and 6 in Table 2. We put no constraint on the number of individuals in the age pyramid, hence the reference population does not give any advantage to the sample-free method.

5 Comparing sample-free and sample-based approaches

The attributes of both individuals and households are respectively described in Table 3 and Table 4. The joint-distributions of both the attributes for individuals and households give respectively the number of individuals of each individual type $n_T = \{n_{t_k}\}_{1 \leq k \leq q}$ and the number of households of each household type $n_U = \{n_{u_l}\}_{1 \leq l \leq p}$. In this case, $q = 130$ and $p = 17$. It's important to note that p is not equal to $6 \cdot 5 = 30$ because we remove from the list of household types the inconsistent values like for example single households of size 5. We do the same for the individual types (removing for example retired individuals of age comprised between 0 and 5).

Table 3: Individual level attributes

Attribute	Value
Age	[0,5[[5,15[[15,25[[25,35[[35,45[[45,55[[55,65[[65,75[[75,85[85 and more
Activity Statut	Student Active Inactive
Family Statut	Head of a single household Head of a monoparental household Head of a couple without children household Head of a couple with children household Head of a other household Child of a monoparental household Child of a couple with children household Partner Other

Table 4: Household level attributes

Attribute	Value
Size	1 individual
	2 individuals
	3 individuals
	4 individuals
	5 individuals
	6 and more individuals
Type	Single
	Monoparental
	Couple without children
	Couple with children
	Other

5.1 Fitting accuracy measures

We need fitting accuracy measures to evaluate the adequacy between both observed O and estimated E household and individual distributions. The first measure is the Proportion of Good Prediction (PGP) (Equation 3), we choose this first indicator for the facility of interpretation. In the Equation 3 we multiplied by 0.5 because as we have $\sum_{k=1}^p O_k = \sum_{k=1}^p E_k$, each misclassified individual or household is counted twice (Harland et al., 2012).

$$PGP = 1 - \frac{1}{2} \frac{\sum_{k=1}^p |O_k - E_k|}{\sum_{k=1}^p O_k} \quad (3)$$

We use the χ^2 distance to perform a statistic test. Obviously the modalities with a zero value for the observed distribution are not included in the χ^2 computation. If we consider a distribution with p modalities different from zero in the observed distribution, the χ^2 distance follows a χ^2 distribution with $p - 1$ degrees of freedom.

$$\chi^2 = \frac{\sum_{k=1}^p (O_k - E_k)^2}{\sum_{k=1}^p O_k} \quad (4)$$

For more details on the fitting accuracy measures see Voas and Williamson (2001).

5.2 Sample-free approach

To test the sample-free approach, we extract from the reference population, for each municipality, the distributions presented in Table 2. Then we use the procedure used for generating the population of reference but now with the constraints on the number of individuals from the age pyramid derived from the reference (remember that we did not have such constraints when generating the

reference population). Then we have filled the households with the individuals one at a time using the distributions for affecting individual into household. We use limit the number of iterations to 1000 trials by household: If after 1000 trials a household is not filled, we put at random individuals in this household and we change his type for "other". We repeat the process 100 times and we choose, for each municipality, the synthetic population minimizing the χ^2 distance between simulated and reference distributions for affecting individual into household.

In order to assess the robustness of the stochastic sample-free approach, we generate 10 synthetic populations by municipalities, yielding 13,100 synthetic municipality populations in total. For each of them and for each distributions for affecting individual into household we compute the p-value associated to χ^2 distance between the reference and estimated distributions. As we can see in the Figure 1a the algorithm is quite robust.

To validate the algorithm we compute the proportion of good predictions for each 13,100 synthetic populations and for each joint-distribution. We obtain an average of 99.7% of good predictions for the household distribution and 91.5% of good predictions for the individual distribution (Figure 1b). We have also compute the p-value of the χ^2 distance between the estimated and reference distributions for each of the synthetic populations and for each joint-distribution. Among the 13,100 synthetic populations 100% are statistically similar to the observed one at a 0.95% level of confidence for the household joint-distribution and 94% for the individual joint-distribution.

In order to understand the effect of the maximal number of iterations by household, we repeat the previous tests for different values of this parameter (1,10,100,500,1000,1500 and 2000) and we compute the mean proportion of good predictions obtained for both individual and household. We note that after 100 the quality of the results no longer changes.

5.3 IPU

To use the IPU algorithm we need a sample of filled households and marginal variables. In order to obtain these data we pick at random a significant sample of 25% of households from the reference population P and we also extract from P the two one-dimensional marginals (Size and Type distributions) that we need to build the household joint-distributions with IPF and the three two-dimensional marginals (Age x Activity Statut, Age x Family Statut and Family Statut x Activity Statut joint-distributions) that we need to build the individual joint-distributions with IPF. Then we apply the Algorithm 3 using the recommendation of Ye et al. (2009) for the well-know zero-cell and zero-marginal problems to obtain a weighted sample P_s . With this sample we generate 100 times the synthetic population P and choose the one with lowest χ^2 distance between reference and simulated individual joint-distributions.

To check the results obtained with the IPU approach, we generate 10 synthetic populations by municipality using different samples of 25% of households randomly selected. For each of these synthetic populations and for each joint-distribution we compute the proportion of good predictions (Figure 2a). We

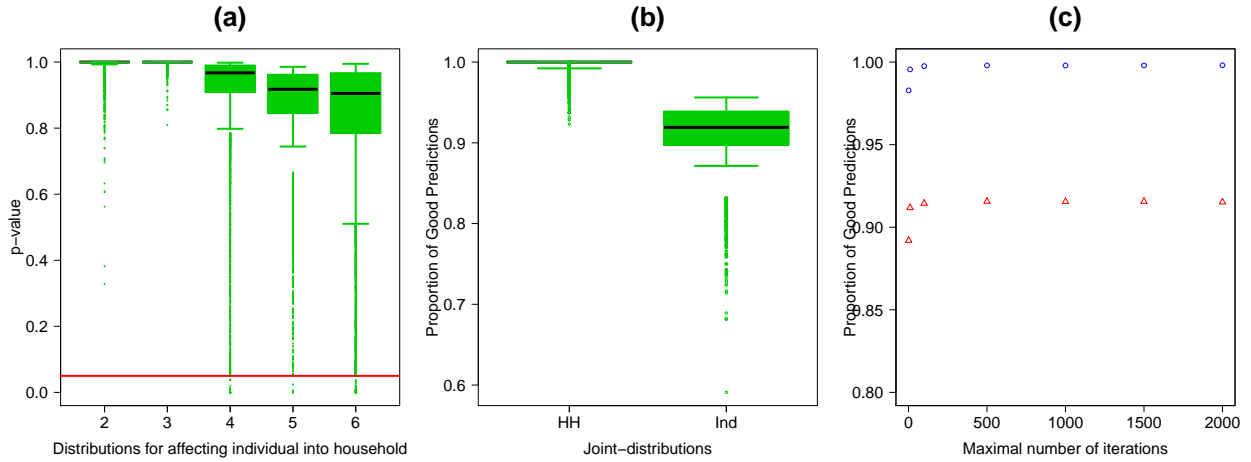


Fig. 1: (a) Boxplots of the p-values obtained with the χ^2 distance between the estimated distributions and the observed distributions for each distributions for affecting individual into household, municipalities and replications. The x-axis represents the controls presented in Table 2. The red line represents the risk 5% for the χ^2 test. (b) Boxplots of the proportion of good predictions for each joint-distribution, municipalities and replications. (b) Average proportion of good predictions in term of the number of maximal iteration by households. Blue circles for the households. Red triangles for the individual.

obtain an average of 98.6% of good predictions for the household distribution and 86.9% of good predictions for the individual distribution. To determine the error of estimation due to the IPF procedure we compute the proportion of good predictions for the estimated and the IPF-reference distributions. As we can see in Figure 2b the results are improved for the household distribution but not for the individual distribution. We also compute the p-value of the χ^2 distance between the estimated and observed distributions for each of the synthetic populations and for each joint-distribution. Among the 13,100 synthetic populations 100% are statistically similar to the observed one at a 0.95% level of confidence for the household joint-distribution and 61% for the individual joint-distribution. We obtained a similarity between the estimated and the IPF-objective distributions of 100% at a 0.95% level of confidence for the household distribution and 64% for the individual distribution.

In order to check the sensitivity of the results to the size of the sample, we plot, on Figure 2c, the average proportion of good predictions of the 13,100 household and individuals joint-distributions for different values of the percentage of the reference households drawn at random in the sample (5, 10, 15, 20, 25, 30, 35, 40, 45 and 50). We note that the results are always good for the household distribution but for the individuals the results are good only from random sample

of at least 25% of the reference household population. Not surprisingly, globally the quality of the results increases with the parameter.

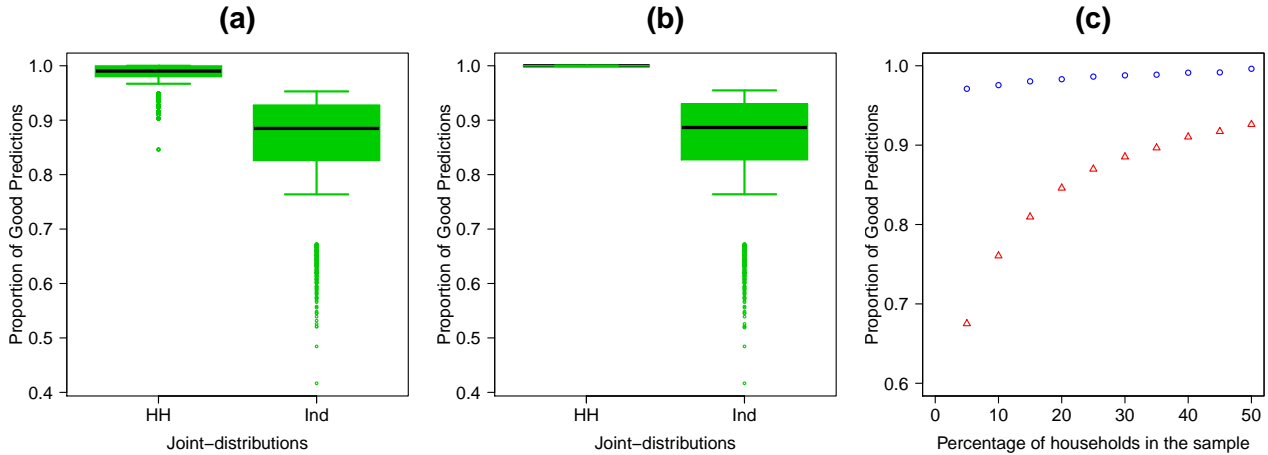


Fig. 2: (a) Boxplots of the proportion of good predictions for a comparison between the estimated distribution and the observed distribution for each municipality and replication. (b) Boxplots of the the proportion of good predictions for a comparison between the estimated distribution and the IPF-objective distribution for each municipality and replication. (b) Average p-values obtained with the χ^2 distance between the estimated distribution and the observed distribution in term of the sample percentage. Blue circles for the households. Red triangles for the individual.

6 Discussion

The sample-free method is less data demanding but the data requires much pre-processing. Indeed, this approach requires to extract the controls from data. The sample-free method gives better fit between observed and simulated distribution for both household and individual distribution than the IPU approach. We can observe in Figure 3 that, for both methods, the goodness-of-fit is negatively correlated with the number of inhabitants. This observation is especially true for the IPU method because it depends of the number of individuals in the sample. Indeed, the lower is the number of individuals, the higher is the number of sparse cells in the individual distribution. The results obtained with the IPU approach depend of the quality of the initial sample. The execution time on a desktop machine (PC Intel 2.83 GHz) is almost the same for 100 maximal iterations by household for the sample-free method and 25% reference households drawn at random in the sample reference households for the sample-based approach.

To conclude, the sample-free method gives globally better results in this application on small French municipalities. These results confirm those of [Barthelemy and Toint \(2012\)](#) who compared their sample-free method for working with data from different sources with a sample-based method ([Guo and Bhat, 2007](#)), and obtained to similar conclusions. Of course, these conclusions cannot be generalized to all sample-free and sample-based methods without further investigation. However, these results confirm the possibility to initialise accurately micro-simulation (or agent-based) models, using widely available data (and without any sample of households).

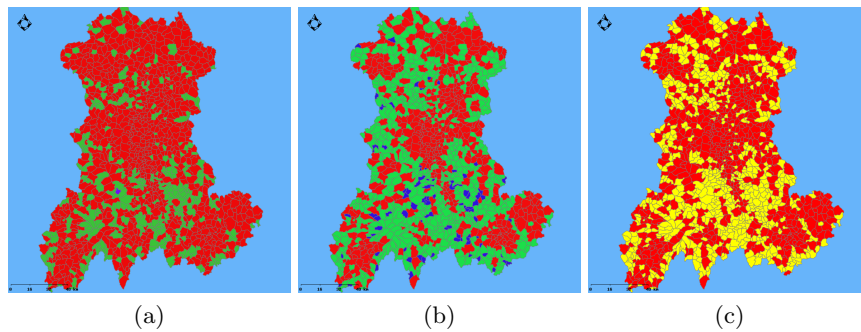


Fig. 3: Maps of the average proportion of good predictions ((a) free-sample and (b) IPU) and the number of inhabitants ((c)) by municipality for the Auvergne case study. For (a)-(b), in blue $0.5 \leq PGP < 0.75$; In green $0.75 \leq PGP < 0.9$; In red $0.9 \leq PGP$. For (c), in yellow, the number of inhabitants is lower than 350. In red, the number of inhabitants is upper than 350. Base maps source: Cemagref - DTM - Développement Informatique Système d'Information et Base de Données : F.Bray & A.Torre IGN (Géofla , 2007)

Table 5: Average execution time for the two approaches for different parameter values.

IPU		Iterative	
Sample size	Time	Iterations	Time
5	13min	1	40min
10	24min	10	41min
15	29min	100	45min
20	38min	500	58min
25	45min	1000	66min
30	53min	1500	78min
40	74min	2000	88min

Bibliography

- Arentze, T., Timmermans, H., and Hofman, F. (2007). Creating synthetic household populations: Problems and approach. *Transportation Research Record: Journal of the Transportation Research Board*, 2014:85–91.
- Barthelemy, J. and Toint, P. L. (2012). Synthetic population generation without a sample. *Transportation Science*.
- Beckman, R. J., Baggerly, K. A., and McKay, M. D. (1996). Creating synthetic baseline populations. *Transportation Research Part A: Policy and Practice*, 30(6 PART A):415–429.
- Deming, W. E. and Stephan, F. F. (1940). On a least squares adjustment of a sample frequency table when the expected marginal totals are known. *Annals of Mathematical Statistics*, 11:427–444.
- Gargiulo, F., Ternes, S., Huet, S., and Deffuant, G. (2010). An iterative approach for generating statistically realistic populations of households. *PLoS ONE*, 5(1).
- Guo, J. Y. and Bhat, C. R. (2007). *Population synthesis for microsimulating travel behavior*. Number 2014 in Transportation Research Record.
- Harland, K., Heppenstall, A., Smith, D., and Birkin, M. (2012). Creating realistic synthetic populations at varying spatial scales: A comparative critique of population synthesis techniques. *Journal of Artificial Societies and Social Simulation*, 15(1):1.
- Huang, Z. and Williamson, P. (2002). A comparison of synthetic reconstruction and combinatorial optimization approaches to the creation of small-area microdata. Working paper, Department of Geography, University of Liverpool.
- Voas, D. and Williamson, P. (2000). An evaluation of the combinatorial optimisation approach to the creation of synthetic microdata. *International Journal of Population Geography*, 6(5):349–366.
- Voas, D. and Williamson, P. (2001). Evaluating goodness-of-fit measures for synthetic microdata. *Geographical and Environmental Modelling*, 5(2):177–200.
- Wilson, A. G. and Pownall, C. E. (1976). A new representation of the urban system for modelling and for the study of micro-level interdependence. *Area*, 8(4):246–254.
- Ye, X., Konduri, K., Pendyala, R., Sana, B., and Waddell, P. (2009). A methodology to match distributions of both household and person attributes in the generation of synthetic populations. In *88th Annual Meeting of the Transportation Research Board*.