

Title	Distribution of Synthetic Populations of Japan for Social Scientists and Social Simulation Researchers
Author(s)	Murata, Tadahiko; Harada, Takuya; Ichikawa, Manabu et al.
Citation	Proceedings - International Conference on Machine Learning and Cybernetics. 2020, 2019-July
Version Type	AM
URL	<a href="https://hdl.handle.net/11094/93359">https://hdl.handle.net/11094/93359</a>
rights	© 2020 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.
Note	

***Osaka University Knowledge Archive : OUKA***

<https://ir.library.osaka-u.ac.jp/>

Osaka University

# DISTRIBUTION OF SYNTHETIC POPULATIONS OF JAPAN FOR SOCIAL SCIENTISTS AND SOCIAL SIMULATION RESEARCHERS

TADAHIKO MURATA<sup>1</sup>, TAKUYA HARADA<sup>2</sup>, MANABU ICHIKAWA<sup>3</sup>, YUSUKE GOTO<sup>4</sup>, LEE HAO<sup>5</sup>,  
SUSUMU DATE<sup>6</sup>, MASAHARU MUNETOMO<sup>7</sup>, AKIYOSHI SUGIKI<sup>7</sup>

<sup>1</sup> Department of Informatics, Faculty of Informatics, Kansai University, Takatsuki, Japan

<sup>2</sup> Department of Industrial and Systems Engineering, Aoyama Gakuin University, Kanagawa, Japan

<sup>3</sup> Department of Planning, Architecture and Environmental Systems, Shibaura Institute of Technology, Saitama, Japan

<sup>4</sup> Faculty of Software and Information Sciences, Iwate Prefectural University, Iwate, Japan

<sup>5</sup> College of Informatics, Shizuoka University, Shizuoka, Japan

<sup>6</sup> Cybermedia Center, Osaka University, Osaka, Japan

<sup>7</sup> Information Initiative Center, Hokkaido University, Hokkaido, Japan

E-MAIL: murata@kansai-u.ac.jp, harada@ise.aoyama.ac.jp, m-ichi@shibaura-it.ac.jp, y-goto@iwate-pu.ac.jp,  
lee@inf.shizuoka.ac.jp, date@ais.cmc.osaka-u.ac.jp, munetomo@iic.hokudai.ac.jp, sugiki@iic.hokudai.ac.jp

## Abstract:

In this paper, we describe how synthesized populations are essential in real-scale social simulations (RSSS), and the current situation of the population synthesis for whole populations in Japan. RSSS is simulations using the real number of populations or households in social simulations. This paper describes how we have completed to synthesize multiple sets of populations based on the statistics of each local government in Japanese national census in 2000, 2005, 2010 and 2015. We have started to distribute those multiple sets of the synthesized populations for researchers of RSSSs in Japan. In distributing the synthesized populations, we should set some regulations in order to protect personal or private information in the synthesized populations.

## Keywords:

Synthetic population; Real-scale social simulations; Big data analysis; Agent simulations; Data protection

## 1. Introduction

In this paper, we try to develop a platform for Real-Scale Social Simulation (RSSS) by synthesizing whole households in Japan and providing the data of synthesized households for researchers who try to develop RSSS tools. RSSS is simulations using populations or households in the real scale.

Recently social simulations have attracted from many researchers to tackle with problems in our environments or communities. One of the most influential social simulations is the segregation model proposed by Schelling [1]. In his model, he clearly shows how segregations happen due to the preference of residents to be a neighbor of the same race or group. His model shows that segregations can happen even if there is no hostility among races but only modest preferences

about neighbors. His model is quite interesting and meaningful to give understanding of conflict and cooperation. He was awarded the 2005 Nobel Memorial Prize in Economic Science for his achievement of game theory including social simulations.

Although Schelling's model is quite significant, interpretation is required to apply his model to real situations. If we are able to directly conduct simulations with real-scale environments and real-scale residents, it is easy to draw insight from simulation results. That is why RSSS has much attention from many researchers recently.

In order to conduct RSSS, real-scale populations are required. For example, when Murata & Konishi [2] optimized the number of polling places with considering the voting rate and the number of polling places in Takatsuki City in Japan using a scheme of RSSS, they should synthesize the population of Takatsuki City and measure the distance of polling places from their homes. When Murata & Hamaguchi [3] applied their method to assess the effectiveness of common voting places that were used in Japan since 2016, they employed their method in Hakodate City, Japan. They should also synthesize the population of Hakodate City to calculate the voting rate according to the population distribution and the voting places in the city. When Murata & Du [4] assessed effects of the pension program for each household in Japan, they should create and simulate demographic movement of all prefectures in Japan for 25 or over 100 years according to the statistics of Japanese census conducted in 2010.

When they try to conduct their RSSS, they should first synthesize the populations of their target cities or communities. Since it is a tough work to synthesize the population according

the real statistics that are distributed by the government, we started to synthesize whole populations in Japan for researchers of RSSS.

## 2. Synthetic reconstruction methods

Since RSSS researchers should face to synthesize populations in the target area of their social simulation sooner or later, we have synthesized populations using the available statistics in each local government such as city, town and village in Japan according to the national census in 2000, 2005, 2010 and 2015. The number of cities, towns and villages in Japan is 1,741 in the national census in 2015.

Methods synthesizing populations with individual attributes are known as Synthetic Reconstruction method (SR method) [5]. Originally an SR method employs real samples from the real statistics. That method increases the number of individuals from the samples in order to fit the real statistics. Here, we prefer using the term “synthesize” to the term “reconstruct” in this paper. Since a reconstruction method is expected to generate exactly the same attributes of each individual in the population, however, it is impossible to reconstruct exactly the same attributes from a small number of statistics. Therefore, we can only synthesize a population that has the same statistical characteristics using SR methods. Lenormand & Deffuan [6] compared SR methods that employ samples with a synthetic method without samples. They showed the synthetic method without samples is better than the former one.

In order to keep real households in the synthesized populations, SR methods using real samples have advantages since they generate households using real households. However, when we employ SR methods with samples in Japan, the synthesized populations cannot be used by others because Japan Statistics Bureau prohibits the distribution of generated data from the samples and the users of samples should delete the anonymized individual samples and their derivatives after the prespecified period. They allow to use their samples up to three years.

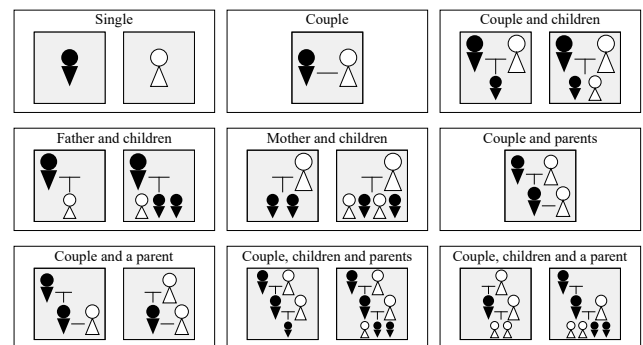
Since our aim to synthesize populations to distribute them for other RSSS researchers, we employ a synthetic method without samples in this paper. The basis of our method is a method proposed by Ikeda et al. [7]. They proposed a method for synthesizing the reduced number of households of nine family types according to the nine real statistics of whole Japan using a simulated annealing method [8]. **Figure 1** shows the nine family types they synthesized. 95% households in Japan come from these family types. Each family member has attributes of sex, age, kinship in its household. Murata & Masui [9,10] modified the objective function and a transition method in their simulated annealing method. Although their

methods [7,9,10] can synthesize a population that has the same statistical characteristics with the real statistics, the number of households in the synthesized population was only 500 or 1,000 households. The synthesized population is too small for social simulations in a real city, town or village.

In order to cope with the problems arisen in the reduced number of populations, we tried to synthesize exactly the same number of individuals in a target area such as states, counties, and prefectures using statistics of prefectures [11,12]. We first increased the number of real statistics for each family type and modified a transition method (Age-Changing method) in their SA method by considering role in a family [11]. We then proposed another transition method (Age-Swap method) in their SA method that keeps the distribution of the initial population that is fit to the real statistics [12]. In the age-swap method, we first initialize the population that is fit to the real demographic pyramid, then exchange age between household members in order to adjust the age difference between couples, or a parent and a child. It has a better performance in reducing the error when the number of transitions in an SA method is relatively small. On the other hand, Age-swap method can reduce better than Age-changing method when the number of transitions in an SA method is relatively large.

When we increase other attributes such as geographical characteristics [13] or occupation and income [14] to the synthesized population, populations by local governments such as city, town or village are required [15]. There are finer statistics that are statistics for each “basic unit block.” The number of “basic unit blocks” in Japan is around 1.9 million. A population synthesis method using statistics of “basic unit block” is proposed by Harada & Murata [16]. We have conducted population synthesis using the above algorithms with high performance computers in Osaka University. **Figure 2** shows an example of a household projection on buildings in a map of Japan. Each figure in a circle shows the number of households residing in the corresponding building.

**TABLE 1** summarizes the SA-based SR Method without samples. Currently we have conducted a method that employs



**FIGURE 1.** Nine family types

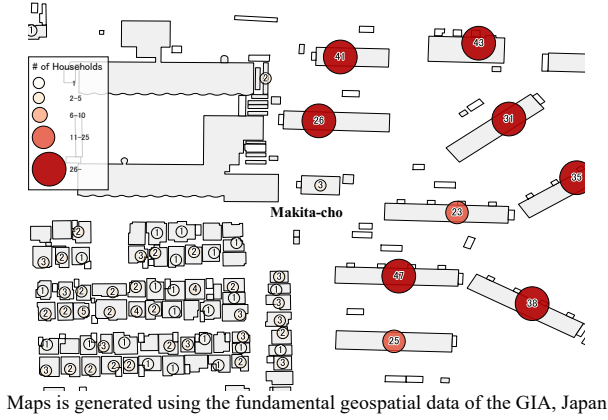


FIGURE 2. Households Projection on buildings

TABLE 1. SA-based SR Method without Samples

Method	Attribute	Population
Ikeda et al. [7]	Age, Sex, Kinship	500-1,000 households
Murata, Masui [9,10]	Age, Sex, Kinship	500-1,000 households
Murata, Harada, Masui [11,12]	Age, Sex, Kinship	A Prefecture
Harada, Murata [13,15,16]	Age, Sex, Kinship, Map	Japan
Murata, Sugiura, Harada [14]	Age, Sex, Kinship, Income, Industry, Company Size	A Prefecture

TABLE 2. Distributed Synthetic Populations

Organization	Synthesized Area	Statistics
RTI International, USA	All states, USA Population: 300 million	2010 US Decennial Census 2007-2011 American Community Survey
CDRC: Consumer Data Research Center, UK	England & Wales, UK Population: (53 + 3) million	2011 UK Census
Kansai University, Japan	Japan Population: 120 million	2000 National Census 2005 National Census 2010 National Census 2015 National Census

[14-16] to synthesize households of whole Japan using the statistics for the smallest district called basic unit blocks and the income estimation procedure.

### 3. Synthesized Population Distribution

Using the above synthetic methods, we have generated

synthesized populations for whole Japan. We are trying to prepare the database of synthesized populations using database. There are only two organizations that distributes synthetic populations in the national level in the world. TABLE 2 shows the distributed synthesized populations. Those organizations distribute nation-wide populations of their country.

Although they are distributing only one set of synthesized populations, we have synthesized 10 to 100 sets of populations of whole Japan now. Since any methods synthesizes populations based on the limited number of statistics, there is no guarantee that the synthesized population is exactly the same as the real population. Therefore, RSSS should be conducted on several sets of populations and find common outcomes from the simulations, or a unique outcome among them. When we conduct simulations on several sets of populations, we may find common results that occur frequently. When we find a common result, it seems to be obtained from any populations with the same statistical characteristics of the real population. On the other hand, when we find some unique results in extreme cases, we should carefully see how the obtained result is caused. In those cases, we may find what we can promote or what we should avoid. In order to conduct such multiple simulations, we should prepare several sets of populations for proper conducts of RSSS.

We are distributing synthetic populations with the notations of the following points.

- 1) The synthetic populations do not contain any data of the real households and individuals.
- 2) The synthetic populations contain only the same statistical characteristics of the real households and individuals.
- 3) The synthetic populations do not contain any statistical characteristics that are not used in the synthetic process.
- 4) The synthetic population will be updated when latest statistics or a modified synthetic method become available.
- 5) Simulations or analysis using the synthetic populations should be conducted on multiple sets of populations.
- 6) Outcomes of simulations and analysis should NOT be released any personal or private information that is relating to real households or individuals.

Although synthesized populations are not real populations, residents may consider that their privacy is offended by releasing their personal information such as their occupations, income or educational back grounds. Therefore, we require researchers to conduct their simulations or analysis using multiple sets of synthesized populations in Item 5). We also require researchers not to release outcomes of their simulations or analysis in any forms that enables others to identify or estimate a private information in a certain household.

#### 4. Conclusions

In this paper, we show the current status of population synthesis of whole populations in Japan. We have developed multiple sets of synthesized populations with the same statistical characteristics of the real populations. In synthesizing the populations, we utilized the statistics conducted in 2000, 2005, 2010, and 2015. After synthesizing the populations with sex, age, kinship in their household, we are increasing attributes of each individual such as geospatial data, occupation, and income. We hope enriching such synthesized population will help researchers who try to develop real-scale social simulations or analyze micro data to see characteristics of our communities or environments.

In order to protect the personal or private information in the synthesized populations, we are planning to employ a cloud service that enables simulations using the synthesized populations. By employing a cloud-style service, we do not have to distribute the synthesized data themselves to researchers but allow them to access the synthesized data in their RSSs. In order to realize such an interface for accessing the synthesized populations, we should develop online programming tools for utilizing the synthetic populations in simulations or analysis.

Another way to protect personal information is to employ secure computation [17]. The secure computation enables users to utilize sensitive data without allowing them to see exact values of them. Using those tools, we carefully protect private information with promoting RSSs.

#### Acknowledgements

Part of this research is funded by JSPS KAKENHI 17K03669 in 2017, Foundation for the Fusion of Science and Technology in 2018, and Tateishi Science and Technology Foundation in 2018. Synthesized populations are generated using the large-scale computing systems (VCC) and OCTOPUS of Cybermedia Center, Osaka University, and distributed using Inter-Cloud System, Hokkaido University under the program of “Joint Usage/Research Center for Interdisciplinary Large-scale Information Infrastructures” and “High Performance Computing Infrastructure” in Japan (Project ID: jh190056-MDH).

#### References

- [1] T. Schelling, “Dynamic models of segregation”, *Journal of Mathematical Sociology*, Vol. 1, pp. 143-186, 1971.
- [2] T. Murata, K. Konishi, “Making a Practical Policy Proposal for Polling Place Assignment Using Voting Simulation Tool”, *SICE Journal of Control, Measurement, and System Integration*, Vol. 6, No. 2, pp. 124-130, 2013.
- [3] T. Murata, Y. Hamaguchi, “Examination of the Effectiveness of Common Voting Places Using Voting Model: The Case of House of Councilors Regular Election in 2016 in Hakodate City, Japan”, *Proceedings of Annual Conference of Japanese Association of Electoral Studies*, 6 pages, 2017 (in Japanese).
- [4] T. Murata, N. Du, “Comparing income replacement rate by prefecture in Japanese pension system”, *Advances in Social Simulation*, pp. 95-108, 2015.
- [5] A. G. Wilson, C. E. Pownall, “A new representation of the urban system for modeling and for the study of micro-level interdependence”, *Area*, vol.8, no. 4, pp. 246-254, 1976.
- [6] M. Lenormand, G. Deffuant, “Generating a synthetic population of individuals in households: Sample free vs sample-based methods”, *Journal of Artificial Societies and Social Simulation*, vol. 16, no. 4, pp. 1-9, 2013.
- [7] K. Ikeda, H. Kita, M. Susukita, “Estimation method of individual data or regional demographic simulations”, *Proc. of SICE 43rd Technical Com. on System Engineer*, pp. 11-14, 2010 (in Japanese).
- [8] L. Davies, *Genetic algorithms and simulated annealing; Research Notes in Artificial Intelligence*, Los Altos, CA: Morgan Kaufmann, 1987.
- [9] T. Murata, D. Masui, “Estimating agents’ attributes using simulated annealing from statistics to realize social awareness”, *Proc. of 2014 IEEE Int’l Conf. on System, Man & Cybernetics*, pp. 717-722, 2014.
- [10] T. Murata, D. Masui, “A two-fold simulated annealing to reconstruct household composition from statistics”, *Proc. of 2015 IEEE Int’l Conf. on System, Man & Cybernetics*, pp. 1133-1138, 2015.
- [11] T. Murata, T. Harada, D. Masui, “Modified SA-based household reconstruction from statistics for agent-based social simulations”, *Proc. of 2016 IEEE Int’l Conf. on Systems, Man, & Cybernetics*, pp. 3600-3605, 2016.
- [12] T. Murata, T. Harada, D. Masui, “Comparing transition procedures in modified simulated-annealing-based synthetic reconstruction method without samples”, *SICE Journal of Control, Measurement, and System Integration*, vol. 10, no. 6, pp. 513-519, 2017.
- [13] T. Harada, T. Murata, “Projecting household of synthetic population on buildings using fundamental geospatial data”, *SICE Journal of Control, Measurement, and System Integration*, Vol. 10, No. 6, pp. 505-512, 2017.
- [14] T. Murata, S. Sugiura, T. Harada, Income allocation to each worker in synthetic populations using basic survey on wage structure, *Proc. of 2017 IEEE Symposium Series on Computational Intelligence*, pp. 471-476, 2017.
- [15] T. Murata, T. Harada, Synthetic method for population of a prefecture using statistics of local governments, *Proc. of 2018 IEEE Int’l Conf. on Systems, Man, & Cybernetics*, pp. 1171-1176, 2018.
- [16] T. Harada, T. Murata, Geospatial data additional method using basic unit blocks, *Proc. of SICE Symposium on Systems and Information 2018*, 6 pages, 2018 (in Japanese).
- [17] K. Chida, G. Morohashi, H. Fuji, F. Magata, A. Fujimura, K. Hamada, D. Ikarashi, R. Yamamoto, “Implementation and evaluation of an efficient secure computation system using ‘R’ for healthcare statistics”, *Journal of the American Medical Informatics Association*, Vol. 21, Is. e2, pp. 326-331, 2014.