

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/355191038>

Creating spatially-detailed heterogeneous synthetic populations for agent-based microsimulation

Article in *Computers Environment and Urban Systems* · October 2021

DOI: 10.1016/j.compenvurbsys.2021.101717

CITATION

1

READS

117

4 authors:



Meng Zhou

Singapore-MIT Alliance for Research and Technology

13 PUBLICATIONS 396 CITATIONS

[SEE PROFILE](#)



Jason Li

Massachusetts Institute of Technology

1 PUBLICATION 1 CITATION

[SEE PROFILE](#)



Rounaq Basu

Massachusetts Institute of Technology

36 PUBLICATIONS 250 CITATIONS

[SEE PROFILE](#)



Joseph Ferreira

Massachusetts Institute of Technology

77 PUBLICATIONS 2,484 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



Energy Transportation and Logistics [View project](#)



Informal and Public Transportation [View project](#)

CREATING SPATIALLY-DETAILED HETEROGENEOUS SYNTHETIC POPULATIONS FOR AGENT-BASED MICROSIMULATION

A PREPRINT

Meng Zhou^{1,2,a}, Jason Li³, Rounaq Basu³, and Joseph Ferreira^{2,3}

¹School of Intelligent Systems Engineering, Sun Yat-sen University, Shenzhen, Guangdong 518107, China

²Future Urban Mobility IRG, Singapore-MIT Alliance for Research and Technology, 138602, Singapore

³Department of Urban Studies and Planning, Massachusetts Institute of Technology, Cambridge, MA 02139, USA

^aCorresponding author; Email: zhoum89@mail.sysu.edu.cn

October 13, 2021

ABSTRACT

1 Agent-based models (ABMs) of urban systems have grown in popularity and complexity due to
2 the widespread availability of high-performance computing resources and large data storage ca-
3 pabilities. Credible synthetic populations are crucial for the application of ABMs to understand
4 urban phenomena. Although several (agent) population synthesis methods have been suggested over
5 the years, the spatial dimension of synthetic populations has not received as much attention. This
6 study addresses this myopic treatment of synthetic populations by creating two distinct components
7 - *agents* and the *built environment* - that are integrated to form a ‘full’ spatially-detailed synthetic
8 population. To generate agents, we used multiple Bayesian Networks (BN) to probabilistically draw
9 pools from the microsample, followed by a Generalized Raking (GR) adjustment to match marginal
10 controls. Using various measures, we demonstrate that our BN + GR framework performs better
11 than more commonly used synthesis methods in both capturing the heterogeneity in the microsample
12 and matching marginal controls. We also highlight the importance of accounting for heterogeneity
13 by using separate type-specific models based on an explicitly defined household typology. For
14 built environment synthesis, we generated various spatial entities such as buildings, housing units,
15 establishments, and jobs at distinct spatial locations by fusing data from various spatial datasets.
16 Their spatial distributions are found to effectively approximate the ‘real’ built environment in our
17 study area. Our proposed framework can be used to generate a ‘full’ synthetic population for use in
18 ABMs with more spatio-demographic heterogeneity than can otherwise be estimated using traditional
19 methods.

20 **Keywords** Synthetic population · Built environment · Agent-based microsimulation · Bayesian Network · Land
21 Use-Transport Interaction (LUTI) model

22 **Citation:** Zhou, M., Li, J., Basu, R. and Ferreira, J. (2021). Creating spatially-detailed heterogeneous
23 synthetic populations for agent-based microsimulation. *Computers, Environment and Urban Systems*. doi:
24 10.1016/j.compenurbssys.2021.101717

25 1 Introduction

26 With the availability of increased computing power, applications of agent-based microsimulations in the fields of
27 transportation and urban studies have burgeoned in recent years (Fagnant and Kockelman, 2014; Waddell, 2002;
28 Salvini and Miller, 2005). In particular, decision support systems such as land use-transport interaction (LUTI) models
29 have increasingly delved deeper to portray the complex interrelationships between urban development and travel
30 behavior with high spatio-temporal resolution through dynamic microsimulations (Acheampong and Silva, 2015; Basu
31 and Ferreira, 2020a; Waddell, 2011). These microsimulation platforms, integrated with econometric models, depict

32 behaviors of various agents (e.g., households or individuals) related to mobility patterns at various spatio-temporal
33 scales in addition to the interactions between agent behaviors and urban systems.

34 To that end, these agent-based models (ABMs) require disaggregate and comprehensive representation of systems
35 they aim to simulate, a major component of which is the *synthetic population*. The purpose of a synthetic population
36 in a microsimulation platform is to characterize the population, including households and individual members, with
37 socio-demographic attributes in rich detail. The extent to which a synthetic population can replicate the ‘real’ (or actual)
38 population has a significant impact on the credibility of the simulation that relies on it. Separate from the transportation
39 domain, a related research stream has focused on constructing *spatial microsimulations* (Tanton et al., 2014; Ballas
40 et al., 2005) that seek to create, analyze, and model individual-level data allocated to geographic zones (Lovelace
41 et al., 2017). Although the two research communities use different terms, spatial microsimulations can be considered
42 analogous to population synthesis as both aim to generate spatially-detailed microdata from samples.

43 While nearly complete information of the real population is collected by national censuses, such data are largely
44 inaccessible even for research purposes due to valid concerns over privacy and security. Instead, microsamples (often
45 referred to as public use microdata samples or PUMS) collected from various types of surveys are often available along
46 with marginal statistics of some key socio-demographic attributes. For example, in the U.S., the PUMS offer detailed
47 data for every individual and household but the spatial resolution is purposely kept low (e.g., a large area with at least
48 100,000 residents) to deter reverse-engineering efforts. Alternatively, marginal distributions of socio-demographic
49 attributes (e.g., number of children in the household) are available at high spatial resolution (e.g., usually up to the block
50 group level). The major challenge in creating a synthetic population for agent-based microsimulations lies in combining
51 agent-based information at coarse spatial resolution with aggregate summary information at high spatial resolution.

52 Despite the growing interest in population synthesis, the spatial dimension of synthetic populations has received limited
53 attention. Most existing approaches assign aggregated zonal information to the synthetic agents and fail to go further
54 in terms of spatial resolution. This may be because the use of ABMs in the LUTI realm has been largely dominated
55 by transportation researchers, who are satisfied with the spatial resolution of aggregated zones (e.g., Traffic Analysis
56 Zones or TAZs) that are adequate for their aim of simulating medium or short-term activity-travel patterns. However,
57 aggregated zones are insufficient for the disaggregate modeling of long-term urban decisions, such as residential and
58 workplace location choices (Zhu et al., 2018). For example, if we are to construct an ABM for exploring housing
59 market dynamics, we would want households to bid on specific housing units in specific buildings at precise locations
60 (not aggregated zones). Thus, we argue that the term ‘*synthetic population*’ has received myopic treatment in the
61 literature and should be extended to include not just agents within the population but also detailed representation of the
62 built environment (e.g., spatial entities such as housing units, buildings, schools, and establishments) that may enable
63 spatially disaggregate allocation of the population. This is in keeping with the rising importance of ‘digital twins’ that
64 seek to include increasingly large and accurate building information models.

65 In this study, we apply state-of-the-art methods to generate a ‘full’ synthetic population accounting for the heterogeneity
66 in household and individual characteristics as well as the marginal controls of key socio-demographics. Additionally, and
67 more importantly, we augment the agent population synthesis by incorporating the construction of the city-wide building
68 population and detailed inventories of housing units and establishments. The integration of these two components,
69 *agents* (i.e., households and individuals) and the *built environment*, results in a spatially disaggregate ‘full’ synthetic
70 population that replicates the urban system at a high spatial resolution. We demonstrate this framework through an
71 application to the city-state of Singapore for the base year of 2016.

72 The remainder of the paper is organized as follows. The next section reviews relevant literature and discusses key gaps
73 and contributions. We briefly introduce the study area and data used for our Singapore application before presenting our
74 framework for the ‘full’ synthetic population generation in Section 3. Section 4 presents the results of our population
75 synthesis framework with comparisons to other popular methods. The paper concludes with discussions and remarks on
76 research extensions in Section 5.

77 2 Literature Review

78 In this section, we first discuss existing population synthesis methods from the transportation literature. Then, we draw
79 on the spatial microsimulation literature to summarize efforts in incorporating spatial detail into synthetic populations.
80 Finally, we reflect on the research gaps within existing literature and propose a few contributions that this study hopes
81 to make.

82 2.1 Population synthesis methods

83 Perhaps the most popular method of population synthesis is the classical Iterative Proportional Fitting (IPF), which
 84 was originally introduced as a technique to adjust contingency tables (Deming and Stephan, 1940) and later widely
 85 applied in urban studies and transportation research for population synthesis (Arentze et al., 2007; Beckman et al., 1996;
 86 Guo and Bhat, 2007; Zhu and Ferreira Jr, 2014). The IPF method fits a multivariate contingency table initialized from
 87 microsample data to the target marginal control distributions in an iterative manner. Despite its conceptual simplicity
 88 and popularity, the IPF algorithm has some notable limitations. Its performance depends heavily on the quality of
 89 microsample data, which are often inadequate due to the recruitment of niche samples or inconsistencies in the sampling
 90 methodology. Additionally, the microsample is likely to reflect only a limited number of attribute combinations, which
 91 limits the heterogeneity of the constructed synthetic population (Sun and Erath, 2015). Trying to obtain unobserved (or
 92 limitedly observed) attribute combinations using IPF results in what is commonly referred to as the ‘zero-cell problem’
 93 (Farooq et al., 2013; Guo and Bhat, 2007). The IPF method also suffers from scalability issues whereby inclusion of a
 94 large number of attributes, especially those with multiple categories, can impose heavy computational burdens (Farooq
 95 et al., 2013; Sun and Erath, 2015). While the IPF usually matches distributions only at one demographic level (i.e.,
 96 either household or individual), a more recent variant known as the Iterative Proportional Updating (IPU) algorithm has
 97 been proposed to allow for matching both household-level and individual-level distributions (Ye et al., 2009). This
 98 algorithm has been implemented in PopGen, an open-source synthetic population generator (Konduri et al., 2016).

99 Another often-used technique - combinatorial optimization (CO) - attempts to reach an optimized solution of population
 100 synthesis by randomly drawing from the microsample while minimizing differences in marginals with algorithms such
 101 as Simulated Annealing (Abraham et al., 2012; Voas and Williamson, 2000). CO-based approaches resemble IPF in
 102 that they also replicate existing agents from the microsample (Sun and Erath, 2015). There are other variants of IPF
 103 or CO such as fitness-based methods (Ma and Srinivasan, 2015) that follow the process of microsample replication.
 104 However, as mentioned earlier, over-dependence on microsample replication can result in several conceptual and
 105 empirical challenges.

106 More recently, researchers have adopted a probabilistic paradigm instead of the deterministic approach of the conven-
 107 tional IPF-based methods (Farooq et al., 2013; Ilahi and Axhausen, 2019; Saadi et al., 2016; Sun and Erath, 2015; Zhang
 108 et al., 2019). These studies break down the synthesis process into two steps: (a) characterization of the joint distribution
 109 of agent attributes, and (b) sampling from the learned joint probability distribution. Thus, the synthesized agents
 110 generated through this approach are *not* replicas of the microsample, and consequently reflect a greater heterogeneity of
 111 agent attributes (Sun et al., 2018).

112 While some studies opted to use Markov Chains for probabilistic population synthesis (Farooq et al., 2013; Saadi
 113 et al., 2016), others adopted data-driven inferential methods such as Bayesian Networks (Sun and Erath, 2015; Zhang
 114 et al., 2019). Markov Chains capture correlations among variables sequentially, which can be challenging to model
 115 when the sequence (or ordering) of variables is unknown and complex interdependencies exist among attributes. In
 116 the two studies using Markov Chains, we observed that the sequence of variables was exogenously pre-determined
 117 instead of being conceptually driven or learnt from the data. Bayesian Networks are comparatively better at inferring
 118 the multivariate probabilistic relationships among attributes, as the joint distributions are determined through graphical
 119 representation.

120 A few studies have tried to adopt the best of both worlds by combining statistical learning techniques and fitting
 121 adjustments to generate synthetic populations that are representative of attribute interrelationships and consistent with
 122 marginal controls. Casati et al. (2015), for example, used MCMC and generalized raking (similar to an augmented IPF)
 123 to synthesize the population. Saadi et al. (2018) combined a Hidden Markov Model (HMM) with IPF and reported
 124 quasi-perfect marginal distributions and relatively accurate multivariate distributions. Ilahi and Axhausen (2019) applied
 125 generalized raking to adjust the BN-based synthesis and reported a good fit to the marginal controls.

126 Over the last couple of years, the popularity of machine learning has motivated the use of deep learning methods
 127 in population synthesis. Methods such as Variational Auto-Encoder (VAE) and Wasserstein Generative Adversarial
 128 Network (WGAN) have been found to work well for high-dimensional cases (Borysov et al., 2019; Garrido et al., 2020).
 129 While deep learning methods provide promising approaches for population synthesis, long-standing issues of machine
 130 learning like the lack of interpretability (i.e., the black-box nature of machine learning models) and the tendency to
 131 overfit the training data remain viable concerns (Basu and Ferreira, 2020c).

132 2.2 Spatial microsimulation methods

133 A related and often overlapping stream of studies beyond the transportation domain is referred to as *spatial microsimu-*
 134 *lation* or, more broadly, small area estimation (Pfeffermann, 2002; Tanton and Edwards, 2012; Tanton et al., 2014).
 135 In essence, spatial microsimulation models seek to simulate the population at spatially disaggregated scales. Such

136 models have been developed for several decades and applied in various domains. Spatial microsimulation extends the
 137 traditional agent-based microsimulation methods to incorporate a spatial dimension (Farrell et al., 2012) and is often
 138 considered to be analogous to population synthesis, or a broader approach of which population synthesis is a crucial part
 139 (Lovelace et al., 2017). Similar to the previously discussed population synthesis approaches, spatial microsimulation
 140 models mainly leverage deterministic or probabilistic sample reweighting techniques - IPF and its variants (Edwards
 141 and Clarke, 2012; Panori et al., 2017), CO (Voas and Williamson, 2000; Farrell et al., 2012), and generalized regression
 142 (Ballas et al., 2007; Tanton et al., 2011; Vidyattama et al., 2013) - to generate synthetic population microdata and assign
 143 them to geographic zones (Lovelace et al., 2017). These models often go beyond the ‘mere’ synthesis of agents and
 144 derive estimates of certain key indicators like income and its inequalities (Vidyattama et al., 2013; Panori et al., 2017)
 145 and obesity (Edwards et al., 2011; Edwards and Clarke, 2012) or model population dynamics over time (Rephann and
 146 Holm, 2004; Kavrouidakis et al., 2012; Birkin et al., 2017). Several spatial simulation models have been operationalized
 147 for policy analysis in various areas such as demography (Ballas et al., 2005; Birkin et al., 2017), healthcare (Edwards
 148 et al., 2011; Edwards and Clarke, 2012), and economics (Campbell and Ballas, 2013; Kavrouidakis et al., 2012).

149 2.3 Research contributions of this study

150 Although the population synthesis literature has continued to evolve in methodological rigor, the methods largely fail
 151 to consider spatial information of the synthetic agents or adequately represent the built environment. Detailed spatial
 152 information is of great value to ABMs seeking to model spatially disaggregate agent behaviors, e.g., housing market
 153 dynamics, evacuation behaviors, pandemic spreads. On the other hand, spatial microsimulation models account for the
 154 spatial dimension but usually assign agents to aggregated geographic zones (Tanton et al., 2014; Lovelace et al., 2017).
 155 Peters et al. (2014) is an exception where housing units are considered but the study has limited explicit representation
 156 of spatial entities. Additionally, most methods utilize reweighting techniques (IPF, CO etc.) that replicate microsamples,
 157 which limits the heterogeneity of the synthetic microdata.

158 The ‘digital twin’ approach that is recently gaining attention aims to provide a digital replication of living as well as
 159 non-living entities that can facilitate the means to monitor, understand, and optimize the functions of all physical entities
 160 and for humans to provide continuous feedback to improve quality of life and well-being (El Saddik, 2018). Translated
 161 to a more ABM-friendly language, this provides an impetus for greater attention to modeling the ‘non-living’ entities
 162 within urban systems (e.g., the built environment comprising housing units, buildings, jobs, schools, and establishments)
 163 by using ‘full’ synthetic populations.

164 This study aims to contribute to the population synthesis and spatial microsimulation literatures on several counts. First,
 165 we propose a combined Bayesian Network and Generalized Raking framework for agent synthesis that can incorporate
 166 microdata heterogeneity and match marginal controls better than more traditional and popular methods such as IPF.
 167 Second, we construct the built environment at a more spatially detailed resolution than in previous studies (e.g., housing
 168 units, buildings, and establishments). Third, we assign synthetic agents to specific housing units and jobs, not just
 169 aggregated zones, that enable us to simulate detailed residential and job location dynamics (although we do not show
 170 these simulation results here). Fourth, by virtue of using a probabilistic sampling design, our agents are truly synthetic
 171 and cannot be traced back to the observations in the microdata, thereby lending an additional layer of privacy to the
 172 original data.

173 3 Research Methods

174 In this section, we first describe the study area of Singapore which we use as a case study to demonstrate the application
 175 of our framework. Second, we provide an overview of the various data sources that are used to construct the ‘full’
 176 synthetic population for Singapore, i.e., both agents and the built environment. Finally, we outline our proposed
 177 frameworks for agent synthesis and built environment synthesis.

178 3.1 Study Area

179 Singapore is a city-state that covers a total area of 719 square-kilometers. Located south of Peninsular Malaysia, it has
 180 a total population of around 5.61 million (as of 2016), among which 3.93 million are local residents (i.e., Citizens and
 181 Permanent Residents) belonging to 1.26 million resident households.¹ The land area of Singapore is divided into six
 182 planning regions and subsequently 55 planning districts (or planning areas), as shown in Figure 1. As of 2016, there are
 183 1,422 TAZs and around 126,000 postcodes in use. Unlike the more conventional definition of postcodes (or ZIP codes)
 184 that most readers may be used to, postcodes in Singapore usually refer to a specific building in most areas (or a block in

¹These statistics are sourced from the Singapore Department of Statistics (commonly referred to as SingStat), available at <https://www.tablebuilder.singstat.gov.sg/publicfacing/mainMenu.action>.

185 less dense and undeveloped areas). Thus, Singaporean postcodes are point features, not polygon features, lending high
 186 spatial resolution to the representation of urban systems (which we will subsequently use for population synthesis).

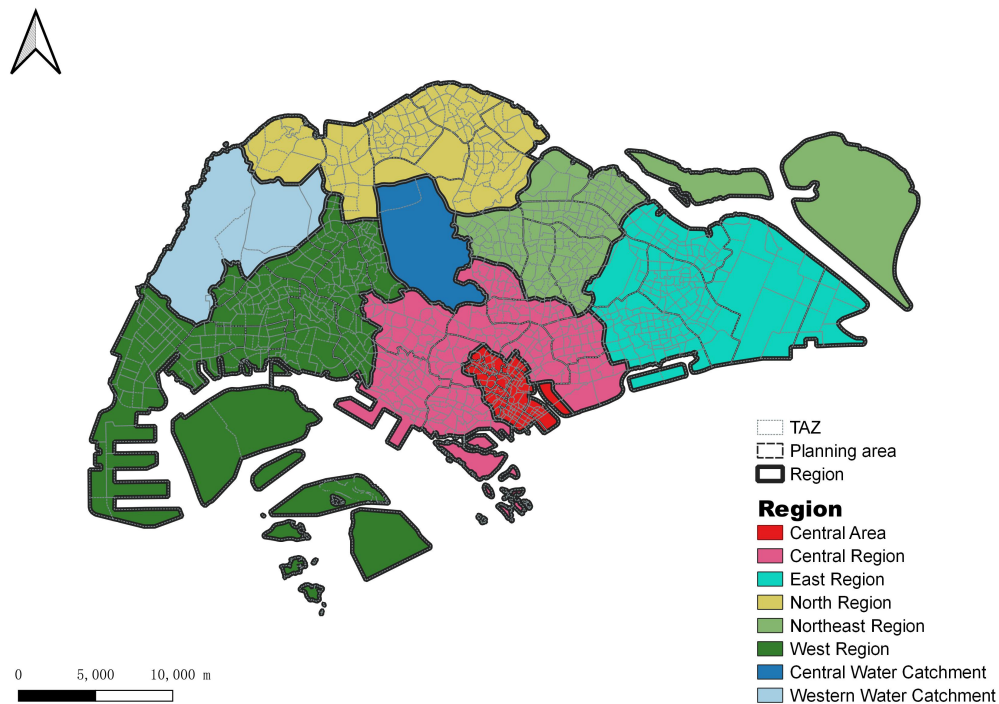


Figure 1: Geographical layout of Singapore ²

187

188 Unlike most of the U.S., land use and housing policies in Singapore prioritize public housing and mixed land use.
 189 Various ‘New Towns’ have been developed that are designed to be self-sufficient in terms of providing everyday
 190 facilities within close proximity. The Housing Development Board (HDB) of Singapore is responsible for public
 191 housing policies and has overseen the implementation of various housing schemes that provide public housing to over
 192 80% of Singaporean households (Singapore Housing & Development Board, 2019). Public housing flats (or HDB flats,
 193 as they are more commonly referred to) can be sold by current owners under certain conditions on the occupancy period.
 194 In addition, there are several types of private housing in Singapore including condominiums, apartments, and landed
 195 properties (i.e., where the property deed includes the land as well as the built structure).

196 Our proposed framework for population synthesis can be generalized to any other metropolitan region with similar data
 197 sources (which we intend to demonstrate through future research). We chose Singapore as the study area in this paper
 198 because of the availability of a rich variety of data sources (some of which are proprietary), and our need for a synthetic
 199 Singapore population to initialize the ABM we have developed for a project to simulate urban futures.

200 3.2 Data

201 In this study, we use the 2016 Household Interview Travel Survey (HITS) as the detailed microsample. The HITS is
 202 a 10% representative sample of Singaporean households that contain at least one resident (i.e., Citizen or Permanent
 203 Resident). This dataset is proprietary and was provided by the Land Transport Authority (LTA) of Singapore that
 204 conducts these travel surveys periodically every four years. That being said, any other representative microsample, such
 205 as a Public Use Microdata Sample (PUMS), a housing survey, or a consumer expenditure survey, would be just as

²The Central Water Catchment is part of the Central Region and the West Water Catchment is part of the West Region. Both are designated as natural reserves that have restrictions on residential, commercial, and industrial uses.

206 viable as long as it contains detailed information on a rich set of variables that are of significance and interest to the
207 phenomena modelers seek to explore through ABMs.

208 IPF-like algorithms also use a set of marginal controls that specify the total number (or proportion) of households or
209 individuals that belong to certain categories across one or more variables. Although most studies in the literature used
210 only socio-demographic marginals, we used marginals that controlled for *both* socio-demographic and spatial variations.
211 These datasets were obtained from the 2015 General Household Survey (GHS), which is available on the open data
212 portal provided by the Singapore government.³ The mid-decade GHS provides comprehensive data on Singapore’s
213 population and households in between the Population Censuses (which are conducted every ten years at the turn of the
214 decade).

215 We used a variety of datasets in this study for built environment synthesis. These data were obtained through
216 collaborative projects from local agency partners or sourced from open sources and include a wide range of information
217 regarding the urban space. Proprietary data such as building addresses (postcodes) and building footprints were provided
218 by the Singapore Land Authority (SLA) for 2016. We also used the open-source land use layer from the 2014 Master
219 Plan created by the Urban Redevelopment Authority (URA). Likewise, public housing building information is openly
220 available on the HDB website. Other open-source third-party data⁴ that provide information on building types, building
221 heights (number of stories), and construction times were also utilized. Additionally, data on zone-to-zone travel times at
222 the TAZ level (i.e., travel skim matrices) were provided by LTA and used for the assignment of jobs to synthetic worker
223 agents.

224 3.3 Full population synthesis

225 Although we present our proposed framework for synthesis of both agents and the built environment in Figure 2, the
226 two components of the framework along with their integration are discussed separately.

227 3.3.1 Agent synthesis

228 In this subsection, we will focus only on the framework for agent synthesis (i.e., the left component of Figure 2).
229 Our synthesis approach for agents, i.e., households and individuals, emulates the two-step process discussed earlier,
230 consisting of sampling from a probabilistic model followed by adjustment to marginal controls using IPF or a similar
231 technique. We chose the Bayesian Network (BN) as our probabilistic model because it has the necessary flexibility to
232 capture heterogeneous multivariate joint distributions, and Generalized Raking (GR) for matching multivariate marginal
233 controls at both the household and individual levels because of its higher efficiency compared to the traditional IPF
234 algorithm.

235 The BN is a graphical model that can learn and represent complex relationships between a large set of variables. It is
236 commonly depicted as a directed acyclic graph⁵ where nodes correspond to variables and edges indicate correlation
237 between variables - this is known as the “structure” of the BN. Each node also possesses a conditional probability
238 distribution that defines the probability of observing certain values, given the values of its parent node(s) - these form
239 the “parameters” of the BN. Using its structure and parameters, a BN can model a complex joint distribution that
240 contains a wide range of variable correlations.

241 Departing from previous literature, we train different BNs based on explicitly defined household types. As households
242 comprise a diverse set of family structures and living arrangements, we expect different household types (or structures)
243 to exhibit different joint distributions representing unique relationships between variables. To this end, the microsample
244 data are partitioned into six sub-samples (Table 1), each corresponding to a household type: *single-member households*
245 (SM), *multi-generational households* (MG), *married couples without co-residing children* (MC), *single parents with*
246 *children* (SP), *nuclear households* (NH), and *others* (OT). These types mirror the household structures defined by the
247 Ministry of Social and Family Development of Singapore (Singapore Ministry of Social and Family Development,
248 2017). The first 5 types (excluding ‘others’) encompass over 90% of Singaporean households, which makes us
249 reasonably confident that this typology should be able to satisfactorily represent a large amount of the variation within

³<https://data.gov.sg/>

⁴For example, Streetdirectory (<https://www.streetdirectory.com/>) and EMPORIS (<https://www.emporis.com/>)

⁵A directed acyclic graph (DAG) consists of vertices and edges (or arcs), with each edge directed from one vertex to another, such that following those directions will never form a closed loop.

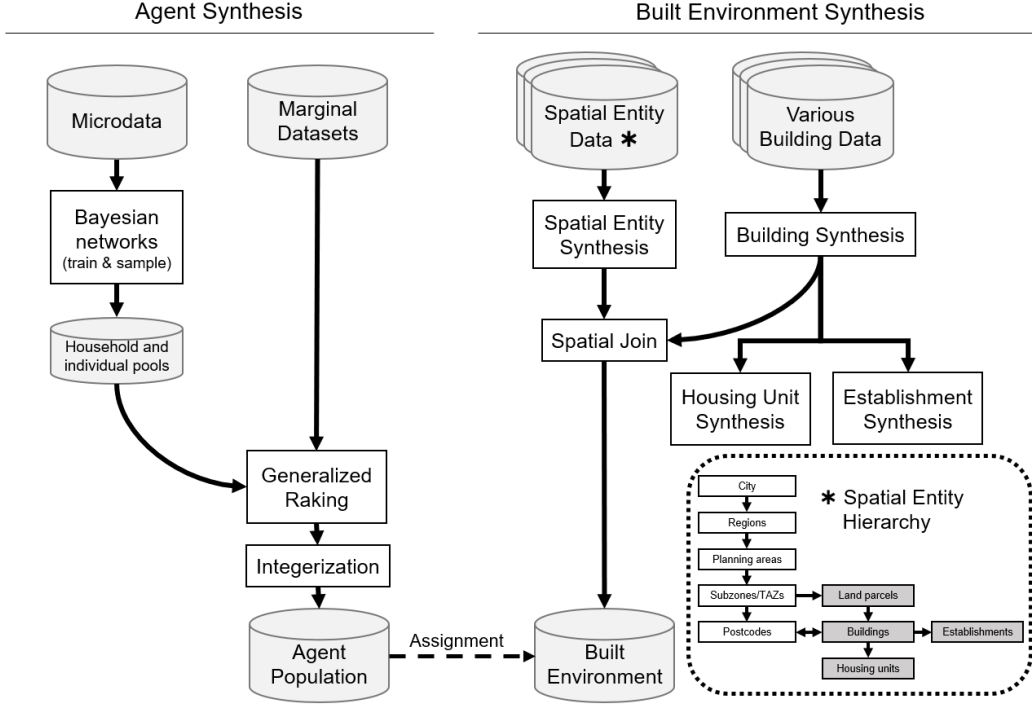


Figure 2: Integrated framework for agent and built environment population synthesis

Table 1: Household typology for agent synthesis

Type	Definition	Count	Share (%)
SM	Single-member (contains exactly 1 individual)	156,950	12.5%
MG	Multi-generational (contains individuals of at least 3 age ranges, each separated by at least 15 years)	118,850	9.5%
MC	Married couple without co-residing children (contains 2 individuals whose age groups are not more than 15 years apart)	186,800	14.8%
SP	Single parent with children (contains 2 or more individuals, exactly 1 of which is at least 15 years older than the others)	87,500	7.0%
NH	Nuclear household (contains 3 or more individuals, exactly 2 of which are at least 15 years older than the others)	598,950	47.7%
OT	Other households (any household that does not fulfill the above criteria, e.g., multiple siblings living together)	108,550	8.6%

250 Singaporean households. We note here that any other classification scheme that is a reasonable representation of the
 251 socio-cultural context of a particular study area will likely be just as appropriate. In addition to modeling heterogeneous
 252 interrelationships, explicitly separating household types also facilitates expert-guided verification and, if necessary,
 253 modification of the BN structure and parameters discovered by the algorithm.

254 BNs can be specified directly based on expert knowledge. In addition, there exist a wide variety of data-driven algorithms
 255 for learning both the structure and parameters of BNs. A popular subset of these are “score-based” algorithms, which
 256 search for a network structure that results in an optimal score measure, such as the Bayesian Information Criterion
 257 (BIC) or Akaike Information Criterion (AIC). In keeping with previous BN-based population synthesis methods, we
 258 chose an algorithm that selects a parsimonious model with good fit based on the AIC measure. A household-level BN
 259 was constructed for each of the six microsample data sub-samples corresponding to a household type. We used the
 260 greedy hill-climbing algorithm provided by the R package *bnlearn* for learning the structure of the BN, followed by

261 maximum likelihood estimation for learning the parameters of the BN (Scutari, 2010). This method is also robust to
 262 missing data, as we can use multiple imputation to obtain several BNs and average them into a single model (not unlike
 263 ensemble-based models in machine learning such as random forests).

264 To create a pool of synthetic households for Singapore, the appropriate number of households of each type (Table
 265 1) were drawn from the corresponding BN using forward sampling, producing a total of 1.26 million households. A
 266 similar, but slightly modified, process was used to create a pool of synthetic individuals. First, the separate household-
 267 and individual-level microsample tables were joined together to create a combined dataset containing all variables
 268 of interest. Modelers may choose to ignore this step for applications where a combined dataset for households and
 269 individuals is directly available. Second, this dataset was also partitioned into the six aforementioned household types,
 270 since differences in variable relationships and distributions at the household level are expected to percolate down
 271 to the individuals within the households as well. Third, individual-level BNs (which contain both household- and
 272 individual-level variables) were learned from the sub-samples. In training the individual-level BNs, we initially used
 273 the entire set of household-level variables used for the household-level BNs and consequently trimmed edges that
 274 were found to be inconsequential or have very weak strength. The resultant individual-level BNs are parsimonious
 275 representations which prohibit edges to household-level variables that do not play a part in the following individual
 276 sampling step.

277 Building on the previous step of sampling households, we drew as many individuals as determined by the sampled
 278 household size variable from the corresponding individual-level BN (that is specific to the household type), by forward
 279 sampling conditional on the sampled household variables. This ensures that the characteristics of household members
 280 correspond to the characteristics of the household they belong to. A pool of roughly 4 million synthetic individuals was
 281 constructed from this individual sampling process.

282 Although the pools of synthetic households and individuals sampled from the BNs can reproduce the microsample’s joint
 283 distributions fairly well, they do not match reported aggregate marginal distributions of certain control variables. This is
 284 because the BNs approximate the joint distributions observed in the microsample, which is an imperfect representation
 285 of the overall population (usually represented through the Census) despite best efforts to obtain representativeness.
 286 Several issues could occur that distort the representativeness of the microsample (e.g., sampling bias, attrition, non-
 287 response), but those are usually beyond the purview of the modeler seeking to construct synthetic populations from a
 288 given microsample and are outside the scope of this paper.

289 For ABMs with a spatial dimension (e.g., location choice simulations), it is imperative to ensure that spatial distributions
 290 of people, housing, and jobs are appropriately represented. Failure to consider the spatial dimension in the population
 291 synthesis approach can affect the veracity of the ABM and may reduce the effectiveness of scenario explorations. In
 292 order to adjust the BN-generated samples spatially, we used the Generalized Raking (GR) procedure from the R package
 293 *MultiLevelIPF* for the household and individual pools to simultaneously fit them to a set of selected multivariate
 294 marginal controls (Mueller, 2018). We chose the planning area as our spatial unit of analysis (recall from Figure 1 that a
 295 planning area in Singapore is equivalent to a neighborhood), as our ABM focuses on location choices that are pertinent
 296 to this level of detail. Thus, the marginal controls we chose to match our synthetic population with are: (a) Planning
 297 Area \times Dwelling Type, Planning Area \times Household Income, and Dwelling Type \times Number of Workers at the household
 298 level, and (b) Planning Area \times Age and Planning Area \times Employment Status at the individual level. Choosing any other
 299 spatial unit (e.g., the commonly used TAZ for traffic simulations) is just as acceptable; we suggest that the modeler
 300 choose the scale of spatial detail for marginals based on the granularity with which they wish to model spatial processes.

301 After convergence, GR produces fractional weights, separately for each household and for each individual. However,
 302 in order to avoid disarranging the grouping of individuals in households, using only the household weights should be
 303 sufficient. To create a synthetic population, an integerization procedure must be performed on the fractional household
 304 weights to convert them to integers. The truncate, replicate, sample (TRS) method is chosen for this purpose (Lovell
 305 and Ballas, 2013). TRS first truncates weights to their integer part, then randomly samples weights to increment by one,
 306 weighted by the fractional part, until the original total weight is restored. We then replicated each household according
 307 to its computed integer weight, replicating the constituent individuals along with it. Since GR and TRS preserve the
 308 total size of the population, we finish the resident synthetic population generation with 1.26 million households and
 309 3.96 million individuals.

310 Finally, since Singapore has a large non-resident population (consisting of 1.67 million individuals in 2016) that factors
 311 significantly into ABMs of mobility choices, we add these households and individuals as a post-processing step. Most
 312 non-residents have residential and job locations that are largely dictated by government and employer policies. For
 313 example, construction workers and single-individual household work permit holders live in assigned dormitories and
 314 work at assigned locations. Therefore, a certain number of individuals holding each work visa type (e.g., Employment
 315 Pass, Student Pass, Construction Work Permit) as determined by statistics from the Ministry of Manpower are inserted
 316 into the synthetic population (Singapore Ministry of Manpower, 2020). Their characteristics are determined by a

317 combination of expert knowledge about the different foreigner demographics and observed variable distributions in the
 318 HITS microsample. Adding non-residents to the synthetic population allows the transportation modeling component of
 319 ABMs to model the full set of daily trips. However, these non-resident households and individuals are not included in
 320 the results and discussion that follow in this paper, since the reference data (i.e., HITS microsample and GHS marginal
 321 controls) that we use for comparison do not include non-residents.

322 3.3.2 Built environment synthesis

323 We represent the urban built environment through a series of spatial entities in a hierarchical structure at different
 324 spatial scales (see the right side of Figure 2). These levels of spatial aggregation (e.g., planning regions, planning
 325 areas, etc.) may vary by the study area but the general strategy of adopting a hierarchical structure of spatial entities
 326 is expected to serve the modeler well for all cases. The shaded boxes in the sub-figure are spatial entities that we
 327 specifically constructed to record the residences and workplaces for households and workers. Their size and location
 328 are estimated using sources independent from the socio-demographic data (on the left side of Figure 2). They are linked
 329 to the demographic data through elements in the spatial hierarchy, usually the postcode.

330 The ontology-based approach is well-suited for the synthesis of the built environment in this study given the variety
 331 of data utilized. Datasets from various sources containing different aspects of the built environment are integrated
 332 based on the ontology that links semantic features that characterize the entities in the built environment. Based on the
 333 created ontology, the integration process then constructs the full list of entities, retrieves attribute values based on the
 334 relationships in the ontology, and imputes missing values with similar entities. Spatial entities are created either in
 335 sequence (in the cases of buildings, housing units, and establishments) or in parallel (in the case of land parcels). We
 336 refer readers to Zhu and Ferreira (2015) for further details of this ontology-based approach.

337 The primary element of built environment synthesis is the creation of buildings, as they form the basis for synthesizing
 338 housing units and establishments. This process includes cleaning building geometry data and inferring various building
 339 attributes. We obtained spatial locations and building geometries directly from building footprint data with relatively
 340 minimal processing (such as merging multiple postcodes that point to the same building). Next, we inferred building
 341 attributes such as building type, height, and space. Building types (e.g., residential, commercial, industrial, etc.) were
 342 inferred first based on the footprint data and third-party datasets matched through postcodes. For residential buildings,
 343 specific types (e.g., public, private, and landed) and subtypes (e.g., condos, apartments, terrace houses, etc.) were also
 344 identified. For most cases, we were able to retrieve the number of stories within each building directly from available
 345 data. For cases with missing data, we used the building heights to estimate the number of stories. The building space
 346 for different types was estimated using the number of stories and the area directly retrieved from the building footprints.
 347 We also estimated the age of the buildings using data on construction and commencement dates, when available.

348 Based on the synthetic buildings and their use types (e.g., residential, commercial, industrial), we then created housing
 349 units of different types and establishments and firms in various industry sectors. Counts of these entities were retrieved
 350 directly from available data or estimated based on building space for the specific use type. We estimated other entity
 351 attributes (e.g., sizes, ages, etc.) based on the relevant characteristics of the buildings. Additionally, we synthesized land
 352 parcels using open-source land use data from URA. We also used spatial data on different amenities (such as public
 353 transit facilities, top schools, shopping malls, and expressway access points) to compute ‘local’ accessibility measures
 354 for each building (and postcode) along the road network.

355 3.3.3 Assigning agents to the built environment

356 Synthetic household agents need to be assigned to housing units, and worker agents (individuals) need to be assigned to
 357 jobs to produce a spatially detailed synthetic population. In this study, we matched housing units to synthetic households
 358 based on their planning areas (neighborhoods) and dwelling types, which are known from the HITS microsample
 359 data. A rule-based matching heuristic was implemented for this assignment. First, a predefined percentage of units
 360 was reserved in each zone-dwelling type bin based on expert knowledge and historical trends, reflecting the vacancy
 361 rate. Then, within each bin, housing units were randomly assigned to households, with larger units more likely to be
 362 assigned to larger and wealthier households. If we ran out of units before all households in that bin were matched, such
 363 households were matched to either a unit of similar dwelling type in the same neighborhood, or, if still unavailable, to a
 364 unit of same dwelling type in a nearby neighborhood.

365 We assigned jobs to worker agents using a destination choice model estimated using the HITS microsample. The
 366 destination choice set of each worker comprises a set of 30 TAZs that contain at least one job pertaining to their
 367 industry sector. The explanatory variables include the number of jobs in that sector, the log-transformed commuting
 368 cost (travel time) between the worker’s home and destination (workplace), and interactions between commuting cost
 369 and individual-level attributes of the worker (e.g., age, gender, income, etc.). The estimated model was used to predict

370 the probabilities of choosing the 30 TAZs for each worker, following which we used a probability-weighted random
 371 assignment to match a worker with a job within the chosen TAZ.

372 4 Results and Discussion

373 We first evaluate the performance of our framework against more commonly used population synthesis methods using a
 374 variety of metrics. Then, we describe the synthetic agents we generated by focusing on the different BNs we learned at
 375 both the household- and individual-levels. Finally, we conclude this section with a discussion of the spatial entities that
 376 are derived from the built environment synthesis.

377 4.1 Model performance

378 Since the built environment synthesis is largely a data integration and fusion process sans the use of any statistical
 379 modeling approaches, we will focus exclusively on the agent synthesis component of our synthetic population generation
 380 framework to evaluate model performance. As a reminder, our synthetic agent population consists of 1.26 million
 381 households and 3.96 million individuals, accurately reflecting the resident population of Singapore in 2016. The
 382 agent population synthesis process, including data imputation and processing, BN learning and sampling, GR, and
 383 integerization, takes about 35 minutes on a x64 laptop PC with 16 GB of RAM and a 1.90GHz Intel Core i7-8650U
 384 CPU.

385 The population synthesis framework used in this study is a combination of Bayesian Networks (BN) and Generalized
 386 Raking (GR), which we call the ‘BN + GR’ method for simplicity. We compared this framework with the BN-sampled
 387 pools generated before the application of GR (i.e., *BN only*), as well as iterative proportional updating (IPU), a multilevel
 388 IPF algorithm that is used as part of the popular open-source software PopGen. We used the IPU implementation from
 389 the R package *ipfr*, with the same marginal controls listed previously to ensure a fair comparison (Ward, 2020). The
 390 synthetic populations generated by these three methods (i.e., BN + GR, BN only, and IPU) are evaluated using two
 391 criteria: (a) similarity to the joint distribution of the weighted microsample, and (b) similarity to the marginal control
 392 distributions.

393 4.1.1 Similarity to the joint distribution of the weighted microsample

394 We evaluate the similarity of the generated synthetic populations to the joint distribution of the weighted microsample
 395 using two methods. First, we use an objective error measure to quantify the extent of the similarity, whereby a lower
 396 error value indicates a greater similarity. Second, we use a graphical method to understand the performance of each
 397 method in greater detail. It is worth noting here that we used the unweighted microsample to generate our synthetic
 398 populations. Sampling weights are usually calculated and provided by the agency conducting the survey in order to
 399 account for stratified sampling or other known deviations from a purely random sample. In our case, there is no need
 400 for sampling weights since the BN sampling process generates a full synthetic population based on the multivariate
 401 correlations observed in the travel survey.

402 As an objective measure of differences between each of our three generated synthetic populations and the one generated
 403 using the HITS microsample with sampling weights, we use the standard root mean square error (SRMSE) as defined
 404 by Sun and Erath (2015):

$$SMRSE = \sqrt{\sum_{m_1=1}^{M_1} \dots \sum_{m_n=1}^{M_n} (f_{m_1, \dots, m_n} - \hat{f}_{m_1, \dots, m_n})^2 \times \prod_{i=1}^n M_i} \quad (1)$$

405 where n is the total number of variables upon which the joint distribution is defined; f_{m_1, \dots, m_n} and $\hat{f}_{m_1, \dots, m_n}$ are
 406 the relative frequencies of a particular variable combination in the weighted microsample and synthetic population,
 407 respectively; and M_i is the total number of categories of the i th variable. AA SRMSE value of zero represents a perfect
 408 match between the two joint distributions under comparison, while higher values represent greater mismatch. However,
 409 we neither expect nor desire a zero SRMSE, since the HITS, even with sampling weights, is an imperfect representation
 410 of the population (as evidenced by HITS’ discrepancies with the marginal controls). SMRSE values for the three
 411 different methods at both the household- and individual-levels are reported in Table 2.

412 We expect the BN-only method to have the lowest SRMSEs, because the BNs directly attempt to model the joint
 413 distribution of the microsample. When GR is applied to match marginal controls, the joint distribution is altered slightly,
 414 which leads to a small increase in SRMSE observed for the BN + GR method. IPU yields the greatest error of all;

Table 2: SRMSE values for different agent population synthesis methods

SynPop method	SRMSE	SRMSE
	(household-level)	(individual-level)
BN + GR (this study)	23.66	7.90
BN only	22.34	6.58
IPU	102.86	9.67

415 this is also unsurprising since it tries to replicate the microsample, which constitutes a significantly smaller sample
 416 size than the BN-sampled pools. IPF-like methods tend to achieve greater accuracy with larger sample sizes that can
 417 capture greater heterogeneity (Wong, 1992), but microsamples are expensive to collect and ‘real-world’ surveys rarely
 418 go beyond a 10% sampling rate (as is the case with HITS in Singapore). Finally, we note that the larger magnitudes of
 419 the household-level SRMSEs are not due to any particular biases of the methods in balancing household-level versus
 420 individual-level fits, but are simply because we define nine variables of interest at the household-level but only six at the
 421 individual-level.

422 The second comparative approach visualizes the fit of the synthetic population to the weighted microsample through a
 423 frequency plot, where the frequencies of every unique variable combination in the two datasets are plotted against each
 424 other. Each point in the frequency plot represents a unique variable combination. A perfect match is represented by a
 425 line of best fit with zero intercept, unit slope, and an R^2 value of one. Additionally, for each plot, we report the Standard
 426 Error around Identity (SEI), which resembles R^2 in that higher values are better but instead measures the extent of
 427 dispersion from the perfect line of best fit (Tanton et al., 2011). Thus, unlike R^2 , the SEI measure can account for
 428 systematic biases in the synthetic population. We present frequency plots at both the household- and individual-levels
 429 in Figure 3.

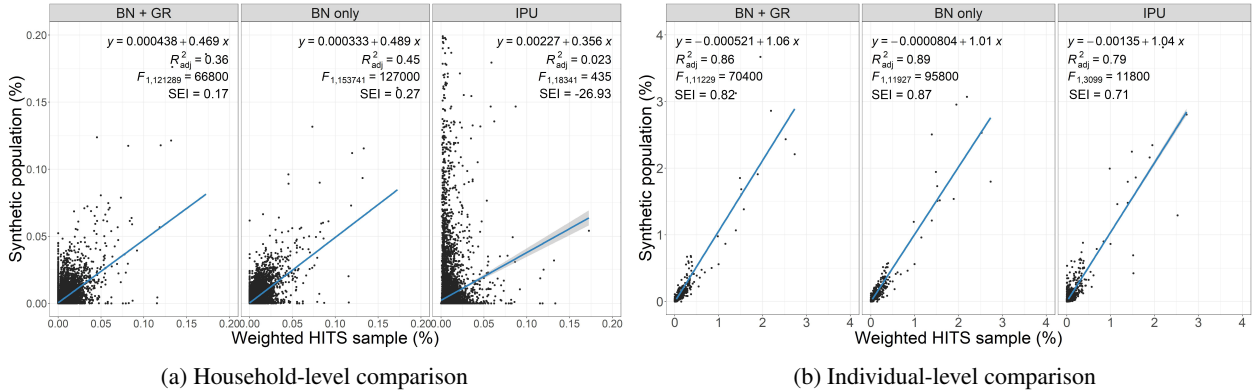


Figure 3: Frequency plots to assess joint distribution match between synthetic populations and weighted microsample

430 Three important observations emerge from this visual analysis. First, the synthetic population generated by the BN-only
 431 method achieves the best fit to the joint distribution, as evidenced by the slope and adjusted R^2 values being closer
 432 to one and the SEI being higher than the rest. Second, the IPU method performs the worst among the three, possibly
 433 because it cannot generate as many variable combinations (there are only 22,494 degrees of freedom for IPU at the
 434 household-level as compared to over 190,000 d.o.f. for BN-based methods). Third, the individual-level distribution
 435 is much easier to fit to than the household-level distribution, as the latter contains more variables (each with more
 436 categories) that leads to increased complexity. In summary, this visual analysis reinforces our observations from the
 437 SMRSE comparisons and shows that BN + GR performs reasonably well at approximating the joint distribution, but not
 438 as well as BN-only. We will explain in the following subsection why BN + GR is a better choice, although it may seem
 439 counterintuitive at this point.

440 **4.1.2 Similarity to the marginal control distributions**

441 We also evaluated the similarity of the synthetic populations generated by the three methods (i.e., BN + GR, BN-only,
 442 and IPU) to the reported marginal control distributions. As a reminder, the marginal controls we chose to match to our

443 synthetic agent populations are: (a) Planning Area \times Dwelling Type, Planning Area \times Household Income, and Dwelling
 444 Type \times Number of Workers at the household level, and (b) Planning Area \times Age and Planning Area \times Employment
 445 Status at the individual level. We explore the similarities to the marginal controls through two methods. First, we look
 446 at bar plots that highlight how the distributions match up for the socio-demographic variables in the aforementioned set
 447 of marginal controls. Second, we look at maps that highlight the spatial variation in the differences (or errors) between
 448 the distributions of marginal controls. To maintain a parsimonious representation of the large set of comparisons that
 449 can be generated (given the number of marginal controls we use), we show bar plots for household income and number
 450 of workers (at the household-level) and age (at the individual-level), and maps for selected dwelling types (at the
 451 household-level) and employment status (at the individual-level).

452 The three bar plots for household income, number of workers, and individual age are shown in Figure 4. In general,
 453 we find that the synthetic populations generated by the BN + GR and IPU methods are able to match the marginal
 454 control distributions almost perfectly. This is unsurprising as these methods include an explicit IPF-like process
 455 that aims to match the reported marginals. On the other hand, the microsample with sampling weights (or weighted
 456 HITS) and the synthetic population generated by the BN-only method have skewed distributions that are often very
 457 different from the marginal controls (see, for example, the case of household income where higher-income households
 458 are under-represented in the microsample). Using sampling weights or only the BN creates an overreliance on the
 459 microsample, which usually falls short of being representative of the population despite best design efforts, and does
 460 not align the microsample with the reported marginal distributions. Thus, it seems clear that both BN + GR and IPU are
 461 equally adept at matching marginal controls and perform much better at this objective than the BN-only method. This,
 462 in combination with our findings regarding the similarity to the joint distribution of the weighted microsample, leads us
 463 to conclude that the BN + GR method ('our' method) achieves a fine balance between the two objectives, which the
 464 other two methods cannot since they perform well on only one of the two objectives.

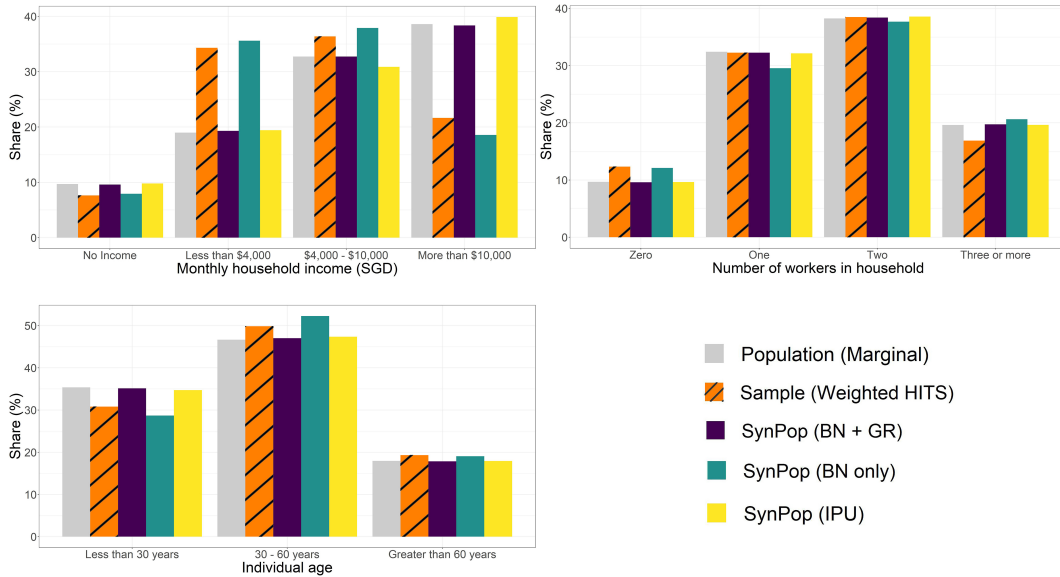


Figure 4: Bar plots to assess distribution match between synthetic populations and socio-demographic marginal controls

465 In addition to exploring the match with socio-demographic marginal controls, we also explored the match with spatial
 466 marginal controls. We present the differences (or errors) between the spatial distributions of the synthetic population
 467 and the marginal controls for selected dwelling types and employment statuses (for brevity) in Figure 5. Instead of
 468 a comparative analysis as earlier, we only show error maps for the synthetic population generated by the BN + GR
 469 method to understand the extent to which systematic spatial biases might be generated by our method, if any. Looking
 470 at the most popular dwelling types for public and private housing (i.e., HDB flats with 3 rooms, and condominiums and
 471 apartments respectively), we do not find any observable non-random patterns of spatial errors. In particular, we find that
 472 almost all planning areas have dwelling type distribution errors between -2.5% to 2.5% , which are quite reasonable.
 473 Errors greater than 5% are infrequent and occur in different planning areas across the maps.

474 We obtain similar observations from the spatial error distribution of the employment status of individuals. In particular,
 475 we find that we are able to predict inactive individual counts (i.e., individuals who are not in the labor force) within
 476 a 2.5% error margin. Our predictions for employed individuals are within the 2.5% error margin for the majority of

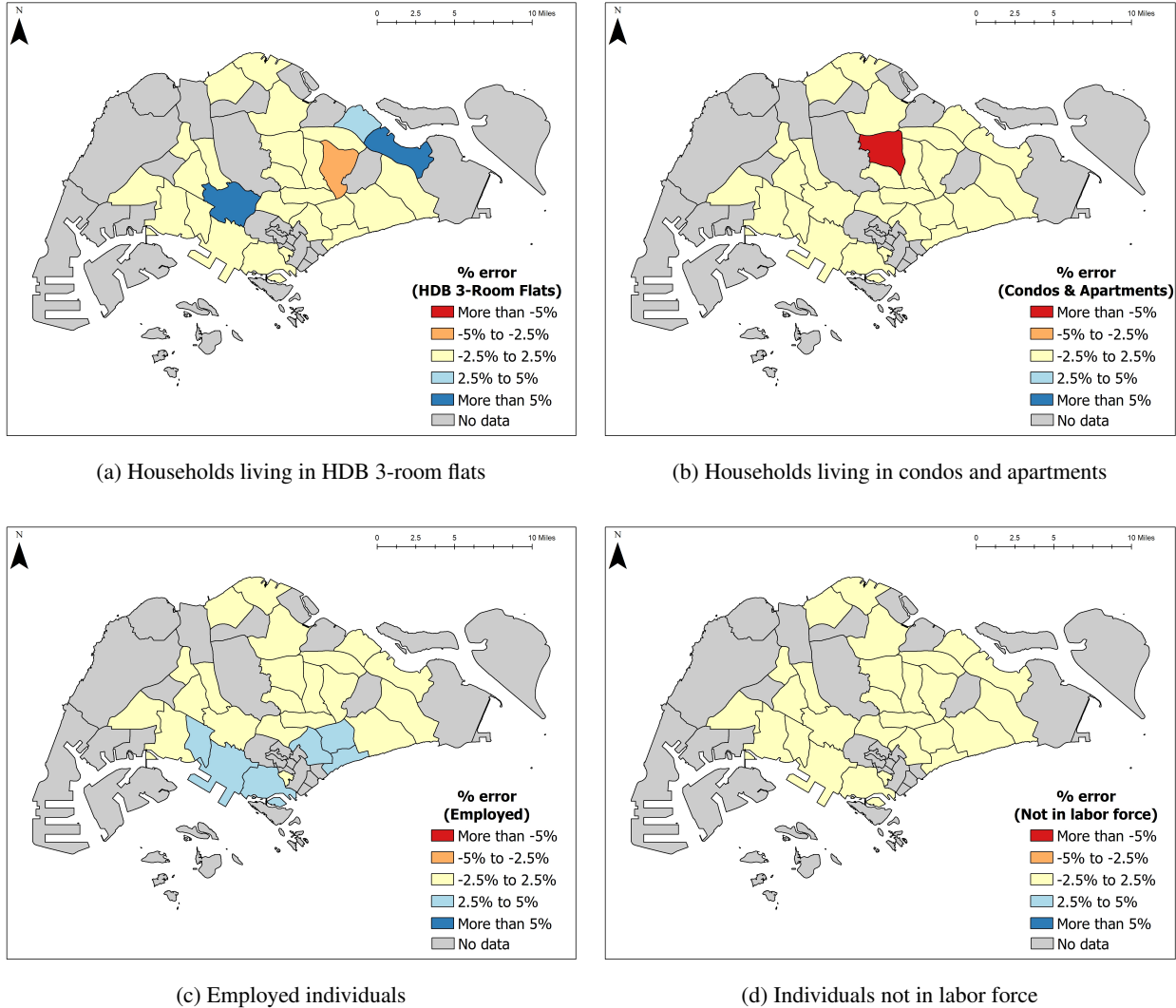
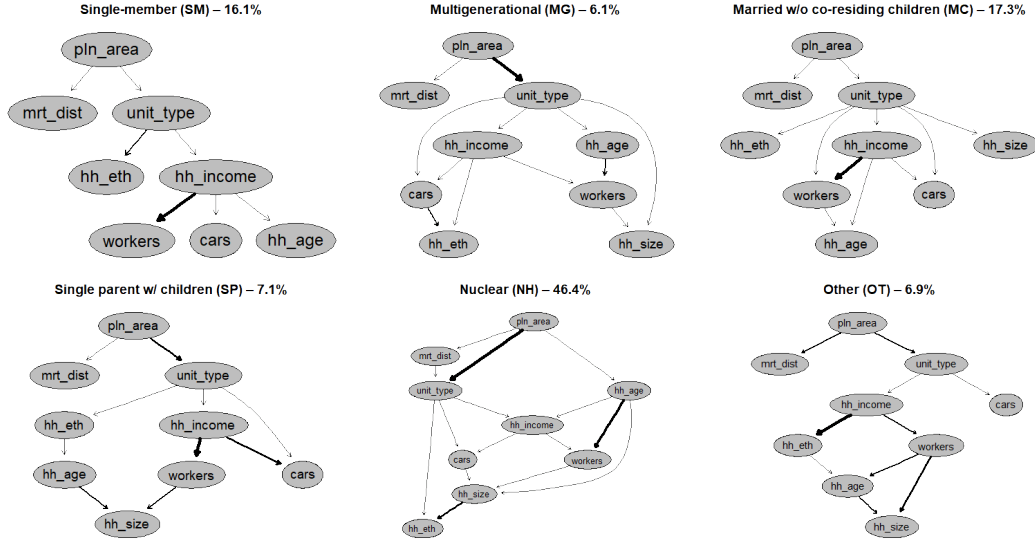


Figure 5: Maps to assess spatial distribution match between BN + GR synthetic population and selected socio-demographic marginal controls

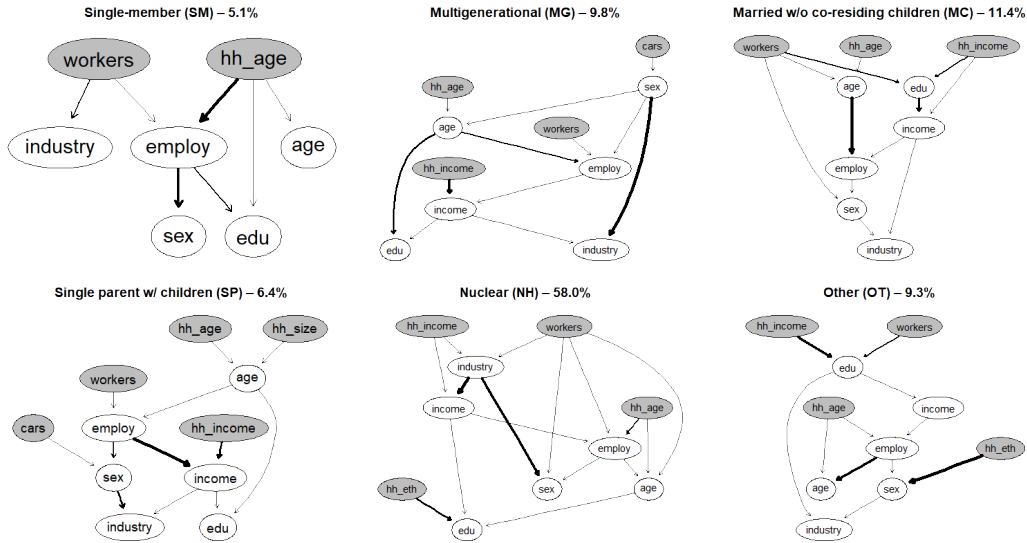
477 planning areas, with the exception of a few cases in the Central Region where the error margin is slightly higher but still
 478 less than 5%. In general, we do not find any reason to suspect that our method may have introduced systematic spatial
 479 biases within the synthetic agent population. Additionally, upon detailed examination of our results, we found that our
 480 predictions after the GR step are quite accurate with less than 1% difference. The errors increase (but not by much, as
 481 demonstrated through the maps) due to the integerization process. An improved integer programming algorithm might
 482 find a better solution, but finding one with adequate performance on our scale of problem is beyond the scope of this
 483 paper.

484 **4.2 Synthetic agents: Households and individuals**

485 In this subsection, we briefly discuss the Bayesian Networks we obtained for households and individuals at the end of
 486 the BN training step in our BN + GR framework. Recall that we had defined a typology of six household categories
 487 based on which we trained six BNs for households and an additional six BNs for individuals. We represent these BNs
 488 in Figure 6, where the nodes are variables and the edges have varying thickness that are directly related to the strength
 489 of the relationship between the nodes they connect. The shaded nodes in the figure are household-level variables, while
 490 the unshaded nodes are individual-level variables. Each variable is treated as a categorical (i.e., ordinal or nominal)
 491 variable with multiple levels (or categories), which are listed in detail in the appendix.



(a) Household-level Bayesian Networks
 (Variable dictionary - *pln_area*: Planning Area, *mrt_dist*: Distance to nearest MRT station, *unit_type*: Dwelling type, *hh_size*: Household size, *hh_income*: Household income, *hh_eth*: Ethnicity of head of household, *hh_age*: Age of head of household, *workers*: Number of workers, *cars*: Number of cars)



(b) Individual-level Bayesian Networks
 (Variable dictionary - *age*: Age, *sex*: Gender, *income*: Income of individual, *industry*: Industry sector of job, *edu*: Highest educational qualification, *employ*: Employment status)

Figure 6: Trained Bayesian Networks for the six household categories at both household and individual levels

492 The household-level BNs shown in Figure 6a share the same ‘root’ relationships, whereby the dwelling type and distance
 493 to the nearest MRT station are both conditional on the planning area where the household resides. The importance
 494 of considering the spatial dimension in population synthesis is underscored by the fact that the planning area is the
 495 root node for all six BNs, implying that socio-demographic correlations are largely dependent on (and vary by) spatial
 496 locations. We observe further commonalities as we proceed ‘deeper’ down the BNs, e.g., the number of workers and
 497 the number of cars owned by the household are conditional on the household income, which in turn is conditional on
 498 the dwelling type. However, the nature of these relationships vary by household type. In single-member households,
 499 single-parent households with children, and married households without co-residing children, the strongest link is

500 unsurprisingly between the number of workers and household income as these households are less likely to earn income
 501 through non-work means (e.g., pension or government benefits). The strongest relationship for multigenerational
 502 households and nuclear households is between dwelling type and planning area, as these choices are likely to be driven
 503 by similar preferences (e.g., proximity to primary schools, community centers, or parks). For nuclear households (who
 504 comprise almost 50% of the Singaporean population), we also notice a strong relationship between the number of
 505 workers and the age of the household head. This is because nuclear households with a younger household head are
 506 likely to have more workers (and vice versa).

507 Individual-level BNs comprise both individual-level variables (unshaded nodes) and household-level variables (shaded
 508 nodes), as shown in Figure 6b. Among all household types, we find common relationships between age and employment
 509 status, employment status and income, and gender and industry sector of job. However, these relationships vary by
 510 strength across the household types. For example, for individuals in nuclear households (forming almost 50% of the
 511 population), there are very strong relationships between the job industry sector and the individual's gender and income.
 512 Additionally, the ethnicity of the household head can influence the highest educational qualification of individuals within
 513 the household (which may reflect ethnic disparities in access to education resources and/or opportunities). For married
 514 households without co-residing children, the age of the individual has a strong impact on their employment status.
 515 For single parents without children, their employment status strongly influences their income. Among single-member
 516 households, the age of the household head (who is also the only individual in the household) relates strongly with
 517 employment status. This perhaps reflects how younger individuals (likely students) are less likely to be employed.

518 Despite all BNs exhibiting a common and 'expected' set of relationships in general, there are important differences
 519 between the BNs for the different types of households. For instance, the household-level BN for nuclear households
 520 is unique in that the number of cars owned is related to household size, which likely reflects the likelihood of car
 521 ownership (and consequently the number of cars owned) being directly proportional to the number of children in
 522 nuclear households. Another example is the individual-level BN for single parents with children, which is unique in that
 523 household size determines individual age. This lines up with our intuition because the age distribution in a single-parent
 524 household depends strongly on the number of children, which is simply one less than household size. Finally, but
 525 importantly, we note that many significant differences between the BNs for different household types are hidden in the
 526 parameters (i.e., node probability distributions) in addition to those observed from the graph structures. For example,
 527 although number of cars owned is determined by household income in both single-member and nuclear households,
 528 single-member households rarely, if ever, own multiple cars, whereas this is not unexpected for nuclear households.
 529 These numerous differences in both structures and parameters between the BNs justify our choice of using a household
 530 typology to learn type-specific BN models.

531 **4.3 Synthetic built environment: Spatial entities and zones**

532 We generated the synthetic built environment comprising various spatial entities in Singapore for the year 2016. First,
 533 our building synthesis process resulted in the creation of 116,415 buildings. We present the spatial distribution of
 534 buildings in Singapore by use type in Figure 7, where we find residential buildings distributed across the island,
 535 commercial buildings mostly located within the central area, and industrial buildings situated mainly in the suburban
 536 areas (particularly along the south-west shore of the island). This spatial distribution lines up with our first-hand
 537 knowledge of Singapore and external data sources (e.g., official land use maps and geospatial services). These buildings
 538 are home to the housing units and establishments of different industries that provide housing and jobs to the household
 539 and individual agents respectively.

540 Using the residential buildings generated through the building synthesis process, we created around 1.66 million
 541 housing units. We purposely create more units than households to allow for a reasonable vacancy rate that can mimic
 542 the 'real' housing market. We also included dormitories reserved for foreign migrant workers and units occupied by
 543 foreigner-headed households (whom we had to explicitly include in our agent population as a post-processing step due
 544 to data unavailability in the microsample). The spatial distributions of different types of housing units are shown in
 545 Figure 8. Public HDB housing units are located within HDB buildings that are located across the island in HDB estates
 546 and New Towns. Private units (i.e., condominiums and apartments), on the other hand, are more clustered within the
 547 Central Region. Landed properties are scattered around the island and other types of units, such as dormitories and
 548 shophouses (i.e., mixed-use landed houses with the ground floor being commercially used while upper levels are used
 549 for residential purposes), are located more in the suburban areas.

550 Similar to the generation of housing units, we generated synthetic establishments using the generated commercial and
 551 industrial buildings to accommodate the employment opportunities for synthetic individuals. 194,044 establishments
 552 are created with around 3.42 million jobs, which is slightly higher than the 3.12 million employed individuals obtained
 553 from the synthetic agent population (to allow for a 'vacancy' rate in the job market, similar to the housing market).
 554 The spatial distribution of establishments proportional to the number of jobs they contain is presented in Figure 9. The

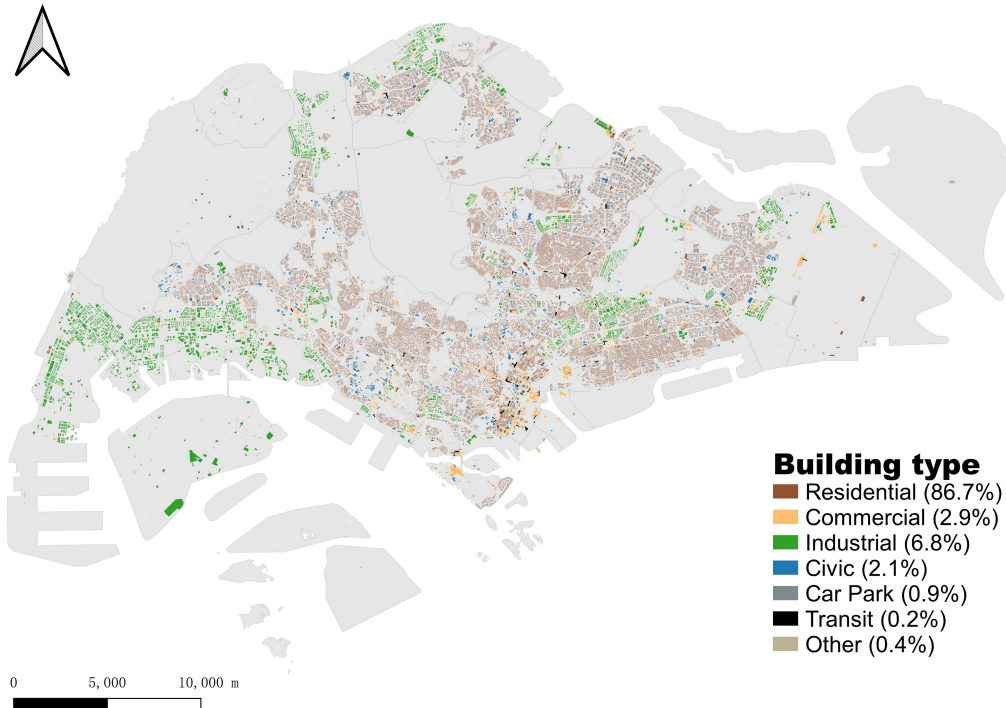


Figure 7: Spatial distribution of synthetic buildings (N = 116,415)

555 distributions of synthetic establishments of different industries are consistent with the distribution of synthetic buildings,
 556 with manufacturing jobs concentrated in suburban areas mostly in industrial buildings and jobs in the finance and real
 557 estate sectors mostly located within the city center in commercial buildings.

558 4.4 Full synthetic population: Linking agents to locations

559 The final step of the population synthesis is to match the synthesized agents with the synthesized spatial entities of the
 560 built environment by assigning housing units to households and jobs to workers. We find that our rule-based heuristic of
 561 matching household agents with housing units works remarkably well. Assuming a vacancy rate of 2%, we were able to
 562 assign housing units to over 99% of the households within their preferred planning area (neighborhood) and dwelling
 563 type. Only 0.8% of households required an adjustment for neighborhood or dwelling type. The spatial distribution of
 564 residential locations of all 1.26 million households is shown in Figure 10a.

565 After assigning housing units to households, we were able to assign a job to each of the 3.12 million workers in their
 566 preferred industry sector. As a recap, we estimated a destination choice model on the HITS microsample whereby
 567 assignment likelihood ratios were (for each industry sector) directly proportional to the number of jobs in the destination
 568 TAZ and inversely proportional to the commute distance. The spatial distribution of job locations of employed
 569 individuals is shown in Figure 10b. Additionally, as a measure of the accuracy of our job assignment, we compared the
 570 job-housing distances for workers in our synthetic population with those in the HITS microsample. We found that the
 571 distributions look quite similar, although we tend to slightly overestimate the commute distances (see Figure A1 in the
 572 appendix). The median commute distance for our synthetic workers is 8.44 kilometers, as compared to 7.93 kilometers
 573 for the HITS sample.

574 After both matching procedures are complete, the spatially assigned agent population is ready for use in large-scale
 575 agent-based microsimulations. These initial assignments can be further adjusted to match predictions of behavioral
 576 models (such as residential and job location choice models) by performing a ‘burn-in’ simulation using the ABM (Basu
 577 and Ferreira, 2020b).

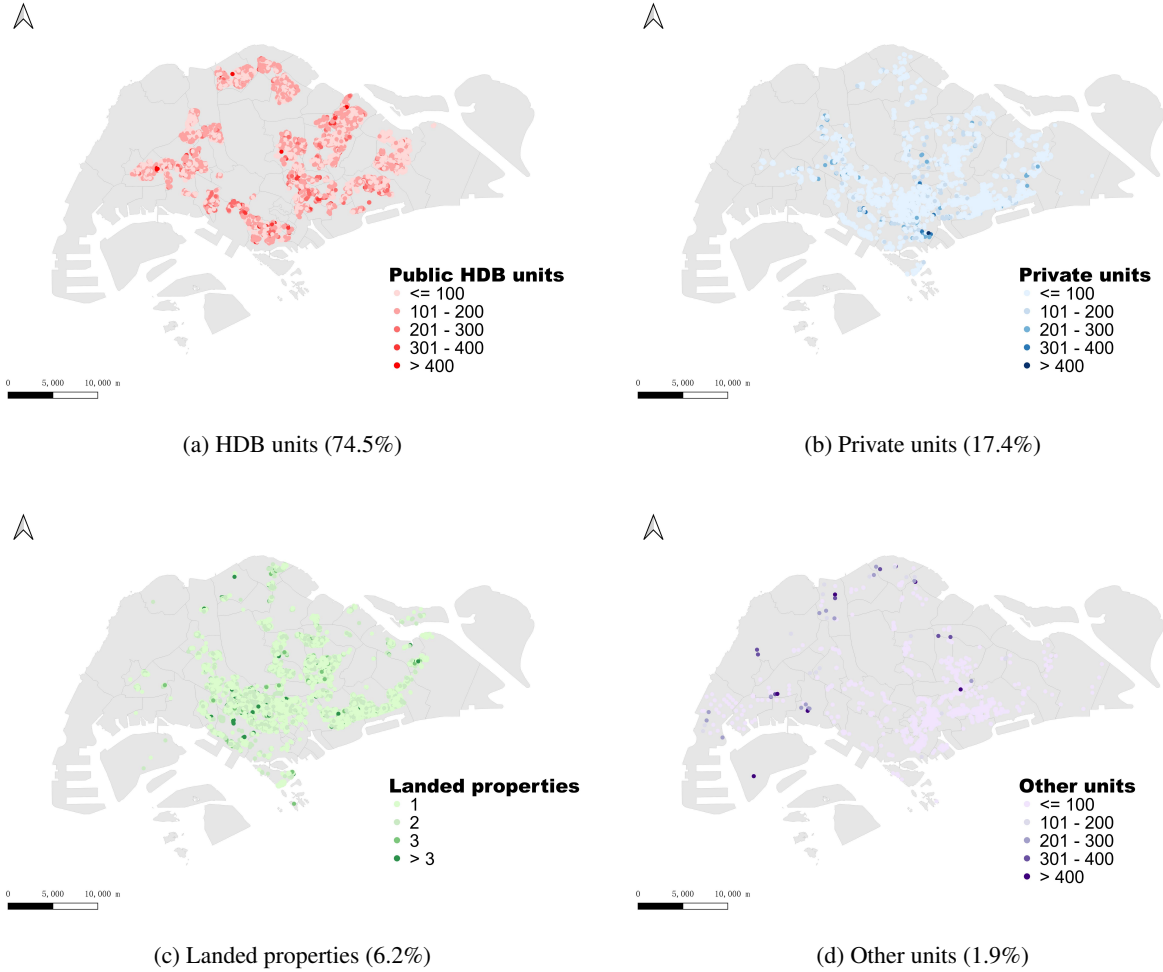


Figure 8: Spatial distributions of synthetic housing units (N = 1,661,284)

588 5 Conclusion

579 Agent-based models (ABMs) of urban systems have been in use for several decades. In recent times, ABMs have grown
 580 in popularity due to the availability of high-performance computing resources and large data storage capabilities. ABMs
 581 also continue to grow in complexity by attempting to model urban systems in increasing spatio-temporal detail. Perhaps
 582 the most crucial component of ABMs is the population they seek to model, thus requiring the creation of a synthetic
 583 population. Data availability challenges affect the resolution at which synthetic populations can be created, whereby
 584 agent-based information at coarse spatial resolution needs to be combined with aggregate summary information at high
 585 spatial resolution. Even though data may be available across agencies, variable definitions and data collection periods
 586 differ, confidentiality issues persist, and considerable time and funding are needed to piece together the elements. We
 587 think it would be worthwhile to invest in periodic construction of detailed synthetic populations that can be used for
 588 many modeling purposes. The construction could be timed to coincide with the periodic travel surveys that many
 589 metropolitan areas conduct every 4-10 years. Several population synthesis methods have been suggested over the years,
 590 starting from iteratively updating weights in a relatively simple manner to complex deep learning models. Despite the
 591 growing research interest in population synthesis, the spatial dimension of synthetic populations has remained largely
 592 neglected. Most existing approaches assign aggregated zonal information to the synthetic agents and fail to go further
 593 in terms of spatial granularity.

594 In this study, we addressed this myopic treatment of the synthetic population by creating two distinct components -
 595 *agents* and the *built environment* - that could be integrated to form what we call the 'full' synthetic population. In terms
 596 of creating the built environment, we generated synthetic spatial entities such as buildings, housing units, establishments,
 597 and jobs at various spatial scales (e.g., postcodes, land use parcels, planning areas, planning regions, etc.). We employed

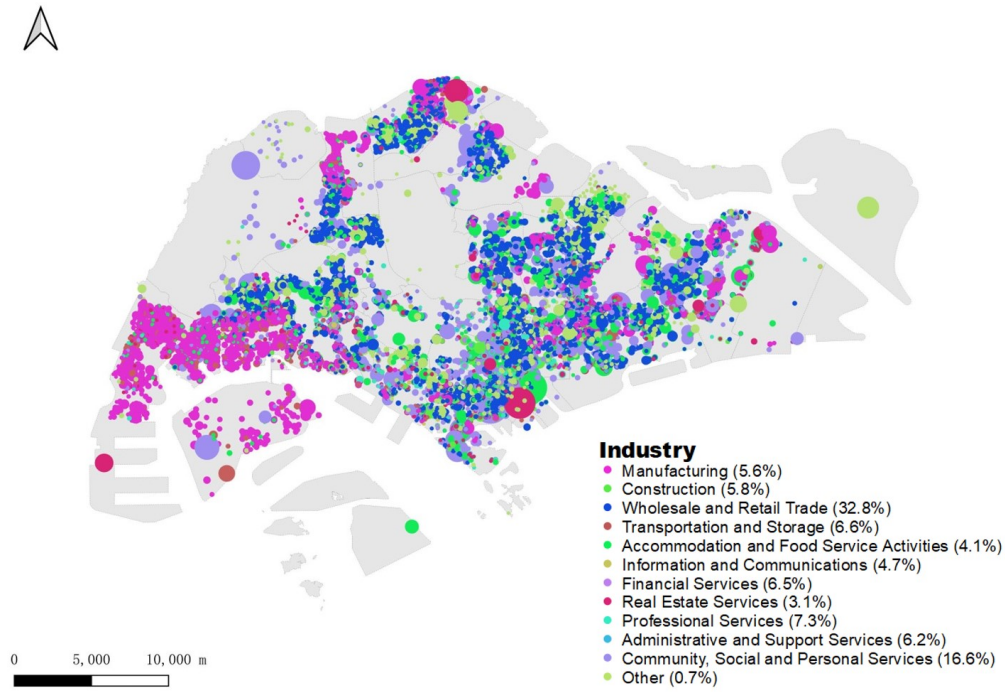


Figure 9: Spatial distribution of synthetic establishments (N = 194,044)

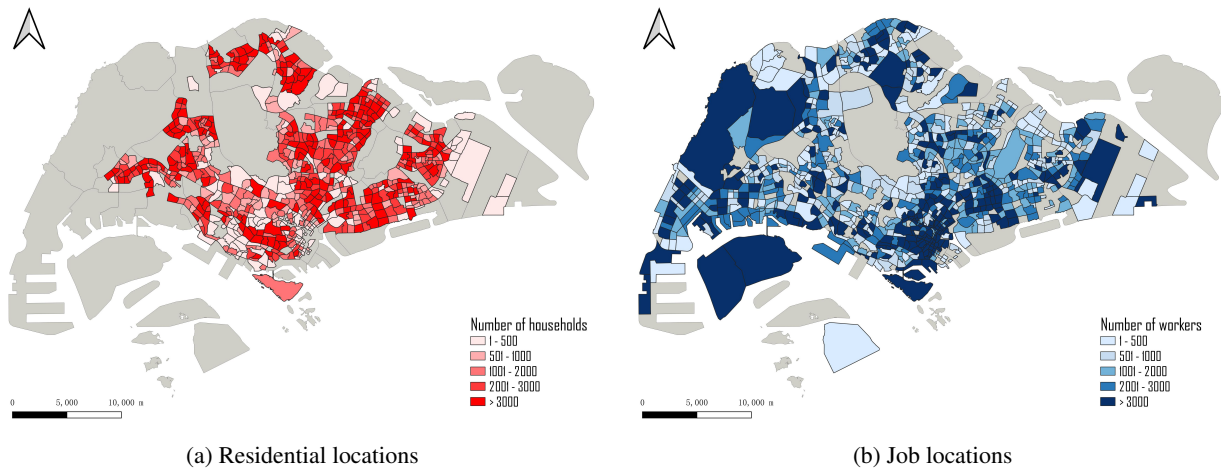


Figure 10: Spatial distributions of residential locations of households and job locations of workers

598 a two-stage framework to probabilistically sample households and individuals from the microsample and subsequently
 599 adjust these pools to match distributions of marginal control variables. Using various measures, we demonstrated
 600 that our BN + GR framework (combining Bayesian Networks and Generalized Raking) performed better than more
 601 commonly used methods (such as IPU and BN only) in both capturing the heterogeneity in the microsample and
 602 matching marginal controls. We also highlighted the importance of accounting for heterogeneity by using separate
 603 type-specific models based on an explicitly defined household typology. Using data fusion techniques on multiple
 604 spatial datasets, we generated various disaggregate spatial entities and found their spatial distributions to match the
 605 ‘real’ built environment in our study area. Thus, we highlighted how our proposed framework can be used to generate a
 606 ‘full’ synthetic population for use in ABMs of any study area of choice.

607 The research presented in this paper can be extended in various ways. One area of future research is the development of
 608 better and faster algorithms for population synthesis. We found that probabilistic models such as the BN can replicate
 609 the microsample well, but an additional proportional update step (using the GR or another IPF-variant) is necessary to
 610 match the marginal controls. The second step of matching reduces the goodness-of-fit of the first step of sampling,
 611 as we highlighted through various error measures. A combined and simultaneous framework (e.g., a multi-objective
 612 optimization routine) could address this issue that arises during sequential adjustment. Additionally, we found that the
 613 integerization process introduced about 2.5-5% errors in the spatial distribution of our synthetic population. Perhaps a
 614 better (and faster) integerization algorithm might be able to reduce these admittedly small errors even further.

615 The exploration of better conceptual frameworks for synthetic populations (and their subsequent use in ABMs) is
 616 another promising research area. While we went further than most in generating disaggregate spatial entities such as
 617 buildings, it is possible to go even further in the pursuit of creating digital twins and synthesizing even the interiors
 618 of buildings. Such detailed synthesis can enable the use of ABMs to model building evacuation techniques in case of
 619 emergencies, building energy use, and airflow within populated buildings (to name but a few applications). The reader
 620 might also wonder if a separate household typology could have resulted in a more ‘accurate’ synthetic population.
 621 We simply chose our typology because it best represented the population in our study area, which is what should
 622 guide modelers. However, it is certainly feasible to explore alternative typologies with different categories or different
 623 numbers of categories, or even sidestep these explicit definitions by deriving the typology from the data (through, e.g.,
 624 latent class analysis). Finally, we note that every population synthesis paper that we reviewed has demonstrated their
 625 proposed framework in only one study area (which we are equally guilty of). It would behoove the ABM community to
 626 begin thinking about extending their population synthesis frameworks to other study areas or to demonstrate the use of
 627 a generalizable framework in multiple study areas.

628 In closing, we hope that we are able to convince readers and the ABM community at large to pay more attention to
 629 standardizing easily repeatable methods for creating synthetic populations in greater spatial detail and with adequate
 630 representation of heterogeneity. We anticipate that ‘full’ synthetic populations (comprising both agents and the built
 631 environment) can enable the exploration of hitherto unanswered research questions about urban processes with high
 632 spatio-temporal granularity.

633 Acknowledgements

634 This research was funded in part by the Singapore National Research Foundation through the Future Urban Mobility
 635 group at the Singapore-MIT Alliance for Research and Technology Center. We appreciate the support of our agency
 636 partners in Singapore for sharing relevant data and information.

637 References

- 638 J. E. Abraham, K. J. Stefan, and J. Hunt. Population synthesis using combinatorial optimization at multiple levels.
 639 Technical report, 2012.
- 640 R. A. Acheampong and E. A. Silva. Land use–transport interaction modeling: A review of the literature and future
 641 research directions. *Journal of Transport and Land use*, 8(3):11–38, 2015.
- 642 T. Arentze, H. Timmermans, and F. Hofman. Creating synthetic household populations: Problems and approach.
 643 *Transportation Research Record*, 2014(1):85–91, 2007.
- 644 D. Ballas, G. P. Clarke, and E. Wiemers. Building a dynamic spatial microsimulation model for ireland. *Population,*
 645 *Space and Place*, 11(3):157–172, 2005.
- 646 D. Ballas, G. Clarke, D. Dorling, and D. Rossiter. Using simbritain to model the geographical impact of national
 647 government policies. *Geographical Analysis*, 39(1):44–77, 2007.
- 648 R. Basu and J. Ferreira. A LUTI microsimulation framework to evaluate long-term impacts of automated mobility on
 649 the choice of housing-mobility bundles. *Environment and Planning B: Urban Analytics and City Science*, 47(8):
 650 1397–1417, 2020a. doi: 10.1177/2399808320925278.
- 651 R. Basu and J. Ferreira. Planning car-lite neighborhoods: Examining long-term impacts of accessibility boosts
 652 on vehicle ownership. *Transportation research part D: transport and environment*, 86:102394, 2020b. doi:
 653 10.1016/j.trd.2020.102394.
- 654 R. Basu and J. Ferreira. Understanding household vehicle ownership in singapore through a comparison of econometric
 655 and machine learning models. *Transportation Research Procedia*, 48:1674–1693, 2020c. doi: 10.1016/j.trpro.2020.0
 656 8.207.

- 657 R. J. Beckman, K. A. Baggerly, and M. D. McKay. Creating synthetic baseline populations. *Transportation Research*
658 *Part A: Policy and Practice*, 30(6):415–429, 1996.
- 659 M. Birkin, B. Wu, and P. Rees. Moses: dynamic spatial microsimulation with demographic interactions. In *New*
660 *frontiers in microsimulation modelling*, pages 53–77. Routledge, 2017.
- 661 S. S. Borysov, J. Rich, and F. C. Pereira. How to generate micro-agents? a deep generative modeling approach to
662 population synthesis. *Transportation Research Part C: Emerging Technologies*, 106:73–97, 2019.
- 663 M. Campbell and D. Ballas. A spatial microsimulation approach to economic policy analysis in Scotland. *Regional*
664 *Science Policy & Practice*, 5(3):263–288, 2013.
- 665 D. Casati, K. Müller, P. J. Fourie, A. Erath, and K. W. Axhausen. Synthetic population generation by combining a
666 hierarchical, simulation-based approach with reweighting by generalized raking. *Transportation Research Record*,
667 2493(1):107–116, 2015.
- 668 W. E. Deming and F. F. Stephan. On a least squares adjustment of a sampled frequency table when the expected
669 marginal totals are known. *The Annals of Mathematical Statistics*, 11(4):427–444, 1940.
- 670 K. L. Edwards and G. Clarke. Simobesity: combinatorial optimisation (deterministic) model. In *Spatial Microsimulation:*
671 *A reference guide for users*, pages 69–85. Springer, 2012.
- 672 K. L. Edwards, G. P. Clarke, J. Thomas, and D. Forman. Internal and external validation of spatial microsimulation
673 models: small area estimates of adult obesity. *Applied Spatial Analysis and Policy*, 4(4):281–300, 2011.
- 674 A. El Saddik. Digital twins: The convergence of multimedia technologies. *IEEE multimedia*, 25(2):87–92, 2018.
- 675 D. J. Fagnant and K. M. Kockelman. The travel and environmental implications of shared autonomous vehicles, using
676 agent-based model scenarios. *Transportation Research Part C: Emerging Technologies*, 40:1–13, 2014.
- 677 B. Farooq, M. Bierlaire, R. Hurtubia, and G. Flötteröd. Simulation based population synthesis. *Transportation Research*
678 *Part B: Methodological*, 58:243–263, 2013.
- 679 N. Farrell, K. Morrissey, and C. O’Donoghue. Creating a spatial microsimulation model of the Irish local economy. In
680 *Spatial microsimulation: A reference guide for users*, pages 105–125. Springer, 2012.
- 681 S. Garrido, S. S. Borysov, F. C. Pereira, and J. Rich. Prediction of rare feature combinations in population synthesis:
682 Application of deep generative modelling. *Transportation Research Part C: Emerging Technologies*, 120:102787,
683 2020.
- 684 J. Y. Guo and C. R. Bhat. Population synthesis for microsimulating travel behavior. *Transportation Research Record*,
685 2014(1):92–101, 2007.
- 686 A. Ilahi and K. W. Axhausen. Integrating bayesian network and generalized raking for population synthesis in greater
687 Jakarta. *Regional Studies, Regional Science*, 6(1):623–636, 2019.
- 688 D. Kavroudakis, D. Ballas, and M. Birkin. Simeducation: A dynamic spatial microsimulation model for understanding
689 educational inequalities. In *Spatial microsimulation: A reference guide for users*, pages 209–222. Springer, 2012.
- 690 K. C. Konduri, D. You, V. M. Garikapati, and R. M. Pendyala. Application of an enhanced population synthesis model
691 that accommodates controls at multiple geographic resolutions. In *Proceedings of the 95th Annual Meeting of the*
692 *Transportation Research Board, Washington, DC, USA*, pages 10–14, 2016.
- 693 R. Lovelace and D. Ballas. ‘truncate, replicate, sample’: A method for creating integer weights for spatial microsimula-
694 tion. *Computers, Environment and Urban Systems*, 41:1–11, 2013.
- 695 R. Lovelace, M. Dumont, R. Ellison, and M. Založnik. *Spatial microsimulation with R*. Chapman and Hall/CRC, 2017.
- 696 L. Ma and S. Srinivasan. Synthetic population generation with multilevel controls: A fitness-based synthesis approach
697 and validations. *Computer-Aided Civil and Infrastructure Engineering*, 30(2):135–150, 2015.
- 698 K. Mueller. Multilevelipf: Implementation of algorithms that extend ipf to nested structures. Available online from
699 <https://github.com/krlmlr/MultilevelIPF>, 2018.
- 700 A. Panori, D. Ballas, and Y. Psycharis. Simathens: A spatial microsimulation approach to the estimation and analysis of
701 small area income distributions and poverty rates in the city of Athens, Greece. *Computers, Environment and Urban*
702 *Systems*, 63:15–25, 2017.
- 703 I. Peters et al. Constructing an urban microsimulation model to assess the influence of demographics on heat
704 consumption. *International Journal of Microsimulation*, 7(1):127–157, 2014.
- 705 D. Pfeiffermann. Small area estimation-new developments and directions. *International Statistical Review*, 70(1):
706 125–143, 2002.

- 707 T. J. Rephann and E. Holm. Economic-demographic effects of immigration: results from a dynamic spatial microsimu-
708 lation model. *International Regional Science Review*, 27(4):379–410, 2004.
- 709 I. Saadi, A. Mustafa, J. Teller, B. Farooq, and M. Cools. Hidden markov model-based population synthesis. *Transporta-
710 tion Research Part B: Methodological*, 90:1–21, 2016.
- 711 I. Saadi, B. Farooq, A. Mustafa, J. Teller, and M. Cools. An efficient hierarchical model for multi-source information
712 fusion. *Expert Systems with Applications*, 110:352–362, 2018.
- 713 P. Salvini and E. J. Miller. Ilute: An operational prototype of a comprehensive microsimulation model of urban systems.
714 *Networks and spatial economics*, 5(2):217–234, 2005.
- 715 M. Scutari. Learning bayesian networks with the bnlearn r package. *Journal of Statistical Software, Articles*, 35(3):
716 1–22, 2010. doi: 10.18637/jss.v035.i03.
- 717 Singapore Housing & Development Board. Key Statistics, 2018/2019 Annual Report. Technical report, 2019. URL
718 <https://services2.hdb.gov.sg/ebook/AR2019-keystats/html5/index.html?&locale=CHS&pn=9>.
- 719 Singapore Ministry of Manpower. Foreign workforce numbers. Technical report, 2020. URL <https://www.mom.gov.sg/documents-and-publications/foreign-workforce-numbers>.
- 721 Singapore Ministry of Social and Family Development. Families and Households in Singapore, 2000 - 2017. Technical
722 report, 2017. URL [https://www.msf.gov.sg/research-and-data/Research-and-Data-Series/Docu-
723 ments/Families%20and%20Households%20in%20Singapore%20-%20Statistics%20Series%202019%20%282000%20-%202017%29.pdf](https://www.msf.gov.sg/research-and-data/Research-and-Data-Series/Documents/Families%20and%20Households%20in%20Singapore%20-%20Statistics%20Series%202019%20%282000%20-%202017%29.pdf).
- 725 L. Sun and A. Erath. A bayesian network approach for population synthesis. *Transportation Research Part C: Emerging
726 Technologies*, 61:49–62, 2015.
- 727 L. Sun, A. Erath, and M. Cai. A hierarchical mixture modeling framework for population synthesis. *Transportation
728 Research Part B: Methodological*, 114:199–212, 2018.
- 729 R. Tanton and K. Edwards. *Spatial microsimulation: a reference guide for users*, volume 6. Springer Science &
730 Business Media, 2012.
- 731 R. Tanton, Y. Vidyattama, B. Nepal, and J. McNamara. Small area estimation using a reweighting algorithm. *Journal of
732 the Royal Statistical Society: Series A (Statistics in Society)*, 174(4):931–951, 2011.
- 733 R. Tanton et al. A review of spatial microsimulation methods. *International Journal of Microsimulation*, 7(1):4–25,
734 2014.
- 735 Y. Vidyattama, R. Cassells, A. Harding, and J. Mcnamara. Rich or poor in retirement? a small area analysis of australian
736 private superannuation savings in 2006 using spatial microsimulation. *Regional Studies*, 47(5):722–739, 2013.
- 737 D. Voas and P. Williamson. An evaluation of the combinatorial optimisation approach to the creation of synthetic
738 microdata. *International Journal of Population Geography*, 6(5):349–366, 2000.
- 739 P. Waddell. Urbansim: Modeling urban development for land use, transportation, and environmental planning. *Journal
740 of the American planning association*, 68(3):297–314, 2002.
- 741 P. Waddell. Integrated land use and transportation planning and modelling: addressing challenges in research and
742 practice. *Transport reviews*, 31(2):209–229, 2011.
- 743 K. Ward. ipfr: List Balancing for Reweighting and Population Synthesis. Technical report, 2020. URL [https:
744 //CRAN.R-project.org/package=ipfr](https://CRAN.R-project.org/package=ipfr).
- 745 D. W. Wong. The reliability of using the iterative proportional fitting procedure. *The Professional Geographer*, 44(3):
746 340–348, 1992.
- 747 X. Ye, K. Konduri, R. M. Pendyala, B. Sana, and P. Waddell. A methodology to match distributions of both household
748 and person attributes in the generation of synthetic populations. In *88th Annual Meeting of the Transportation
749 Research Board, Washington, DC*, 2009.
- 750 D. Zhang, J. Cao, S. Feygin, D. Tang, Z.-J. M. Shen, and A. Pozdnoukhov. Connected population synthesis for
751 transportation simulation. *Transportation research part C: emerging technologies*, 103:1–16, 2019.
- 752 Y. Zhu and J. Ferreira. Data integration to create large-scale spatially detailed synthetic populations. In *Planning
753 support systems and smart cities*, pages 121–141. Springer, 2015.
- 754 Y. Zhu and J. Ferreira Jr. Synthetic population generation at disaggregated spatial scales for land use and transportation
755 microsimulation. *Transportation Research Record*, 2429(1):168–177, 2014.
- 756 Y. Zhu, M. Diao, J. Ferreira, and P. C. Zegras. An integrated microsimulation approach to land-use and mobility
757 modeling. *Journal of Transport and Land Use*, 11(1):633–659, 2018.

Table A1: Household-level variables used in the BNs

Variable	Categories	Sample share (%)
Dwelling type (<i>unit_type</i>)	HDB 1- and 2-Room Flats	5.0%
	HDB 3-Room Flats	20.5%
	HDB 4-Room Flats	35.9%
	HDB 5-Room and Executive Flats	25.7%
	Condominiums and Apartments	6.6%
	Landed properties	6.2%
	Others	0.01%
Household size (<i>hh_size</i>)	One	3.8%
	Two	13.8%
	Three	19.6%
	Four	27.9%
	Five	19.2%
	Six or more	15.8%
Monthly household income (<i>hh_income</i>)	No Income	4.5%
	Less than \$1,000	2.8%
	\$1,000 to \$2,000	7.6%
	\$2,000 to \$4,000	23.6%
	\$4,000 to \$6,000	21.0%
	\$6,000 to \$10,000	20.0%
	\$10,000 to \$15,000	12.6%
\$15,000 to \$20,000	3.1%	
Number of workers (<i>workers</i>)	More than \$20,000	5.0%
	Zero	7.2%
	One	27.4%
	Two	41.7%
Number of cars (<i>cars</i>)	Three or more	23.7%
	Zero	61.6%
	One	33.6%
	Two	4.0%
Age of head of household (<i>hh_age</i>)	Three or more	0.7%
	15 to 30 years	1.7%
	30 to 60 years	67.2%
Ethnicity of head of household (<i>hh_eth</i>)	More than 60 years	31.1%
	Chinese	72.7%
	Indian	11.4%
	Malay	13.2%
Distance to nearest MRT station (<i>mrt_dist</i>)	Others	2.7%
	Less than 400 meters	16.6%
	400 to 800 meters	37.0%
	More than 800 meters	46.5%

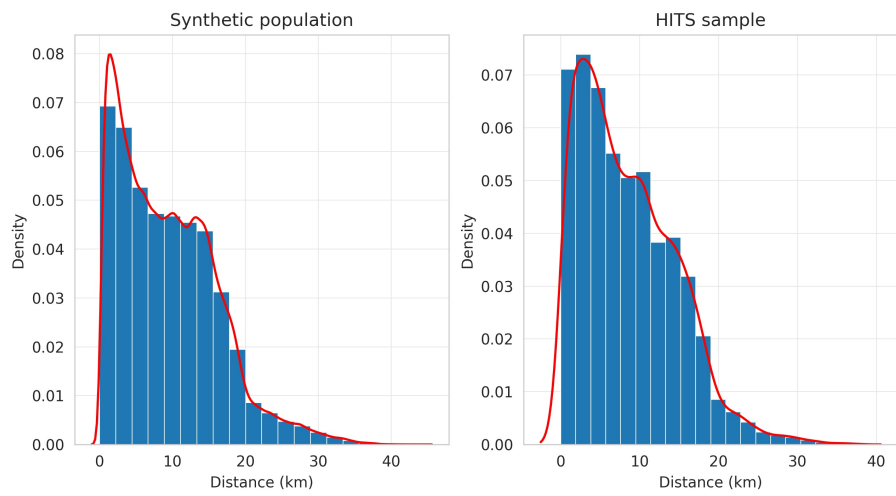


Figure A1: Distributions of job-housing distances for workers

Table A2: Individual-level variables used in the BNs

Variable	Categories	Sample share (%)
Gender (<i>sex</i>)	Male	47.0%
	Female	53.0%
Age (<i>age</i>)	Less than 15 years	6.6%
	15 to 30 years	51.6%
	30 to 60 years	20.4%
	More than 60 years	21.4%
Monthly individual income (<i>income</i>)	No Income	44.1%
	Less than \$1,000	3.8%
	\$1,000 to \$2,000	9.9%
	\$2,000 to \$4,000	22.7%
	\$4,000 to \$6,000	11.2%
	\$6,000 to \$10,000	5.5%
Industry (<i>industry</i>)	More than \$10,000	2.9%
	Accommodation and Food Services	3.6%
	Administrative and Support Services	6.0%
	Community, Social and Personal Services	10.3%
	Construction	3.4%
	Financial Services	4.1%
	Information and Communications	4.1%
	Manufacturing	6.0%
	Professional Driver	1.2%
	Professional Services	7.8%
	Real Estate Services	0.9%
	Transport and Storage	4.5%
	Wholesale and Retail Trade	4.4%
	Others	0.1%
None (for those without a job)	44.4%	
Employment status (<i>employ</i>)	Employed Full Time	46.5%
	Employed Part Time	5.4%
	Self-Employed	3.6%
	Full Time Student	18.0%
	Retired	9.2%
	Homemaker	13.4%
	Unemployed	2.7%
Highest educational qualification (<i>edu</i>)	Others	1.1%
	Primary	6.1%
	Secondary	18.4%
	Post-Secondary	5.6%
	Polytechnic	11.5%
	Bachelor's	17.1%
	Master's/Doctorate	5.4%
	Postgraduate certification	3.6%
	Professional degree	3.7%
Others	19.2%	
None	9.4%	