

Estimating the population size of persons contending with homelessness using electronic health records

Gracia Y. Dong^{1,2,3} , Kenneth Moselle⁴, Stanley Robertson⁵,
Patrick Brown^{1,3}  and Laura L. E. Cowen² 

¹Department of Statistical Sciences, University of Toronto, 700 University Ave, Toronto, Ontario M5G1X6, Canada

²Department of Mathematics and Statistics, University of Victoria, 3800 Finnerty Road, Victoria, British Columbia V8P5C2, Canada

³Centre for Global Health Research, St Michael's Hospital, 209 Victoria Street, Toronto, Ontario M5B1W8, Canada

⁴Department of Psychology, University of Victoria, 3800 Finnerty Road, Victoria, British Columbia V8P5C2, Canada

⁵Applied Health Data Analytics, Vancouver Island Health Authority, 1952 Bay Street, Victoria, British Columbia V8R1J8, Canada

Address for correspondence: Gracia Y. Dong, Department of Statistical Sciences, University of Toronto, 700 University Ave, Toronto, Ontario M5G1X6, Canada; Department of Mathematics and Statistics, University of Victoria, 3800 Finnerty Road, Victoria, British Columbia V8P5C2, Canada; Centre for Global Health Research, St Michael's Hospital, 209 Victoria Street, Toronto, Ontario M5B1W8, Canada. Email: gracia.dong@utoronto.ca

Abstract

The majority of attempts to enumerate the homeless population rely on point-in-time or shelter counts, which can be costly and inaccurate. As an alternative, we use electronic health records from the Vancouver Island Health Authority, British Columbia, Canada from 2013 to 2022 to identify adults contending with homelessness based on their self-reported housing status. We estimate the annual population size of this population using a flexible open-population capture–recapture model that takes into account (1) the age and gender structure of the population, including aging across detection occasions, (2) annual recruitment into the population, (3) behavioural-response, and (4) apparent survival in the population, including emigration and incorporating known deaths. With this model, we demonstrate how to perform model selection for the inclusion of covariates. We then compare our estimates of annual population size with reported point-in-time counts of homeless populations on Vancouver Island over the same time period, and find that using data extracts from electronic health records gives comparable estimates. We find similarly comparable results using only a subset of interaction data, when using only ER interactions, suggesting that even if cross-continuum data is not available, reasonable estimates of population size can still be found using our method.

Keywords: British Columbia, capture–recapture, hidden population, Jolly-Seber, population size, public health

1 Introduction

Homelessness is a difficult state to define, as individuals in identical situations may or may not identify as homeless based on their unique circumstances. In Canada, homelessness is defined as a lack of stable, safe, permanent, appropriate housing, or the immediate prospect, means and ability of acquiring it (Gaetz et al., 2012). However, individual jurisdictions within Canada have their own categories of homelessness which may not align with these categories. In addition to being

Received: October 18, 2023. Revised: April 2, 2024. Accepted: April 5, 2024

© The Royal Statistical Society 2024.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

difficult to define, homelessness is difficult to measure or enumerate (Williams, 2010; Williams & Cheal, 2002), as persons contending with homelessness are considered ‘rare and elusive’. They constitute only a tiny proportion of the overall population, and not all persons contending with homelessness are ‘visible’ to the system, since they are less likely to appear on any kind of official list due to movement, location, or distrust of the system. For example, a city with more hostels and shelters provides more visibility than one with fewer, whereas poor weather can decrease visibility for point-in-time (PIT) counts (Wright & Devine, 1995). Policies against street homeless will decrease visibility for those living rough, but may change probabilities of detection for different groups who may be more likely to utilize shelter services.

Several methods of enumerating the homeless population are currently used in practice, the most common being PIT counts, including shelter counts and performing a street census. Other methods include snowball sampling (Dávid & Sniijders, 2002) and network scale-up methods (Killworth et al., 1998; McCormick & Zheng, 2013), which involve using a social network approach where the social networks of subjects are used to recruit future subjects. However, these methods can be costly and are prone to inaccuracies, as several issues can arise when attempting to enumerate the population (Williams, 2010; Williams & Cheal, 2002). The first issue is undercounting, since it is very easy to miss persons contending with homelessness because there are not enough volunteers performing counts or not every agency or organization that interacts with persons contending with homelessness is contacted. Individuals may also purposely avoid being enumerated, by obtaining knowledge of a count happening, allowing them to hide or move during enumerations. Street canvases have been estimated to miss 40–70% of hidden homeless, especially hidden from public view at night, which is an issue because one of the conventional definitions of homelessness is lacking a fixed regular adequate night-time residence (Berry, 2007). Snowball sampling, where survey participants recruit additional participants from their social network, is commonly used for homeless populations (Forchuk et al., 2018; Kincaid, 2019). However, this method has the potential to exclude a sizeable portion of the population due to spatially clustered and age-stratified social networks—individuals who do not have overlapping social groups with the initial participants are less likely to be included in the study. The second issue is double counting, which occurs because persons contending with homelessness are often spatially mobile, and can be counted twice in different neighbourhoods or cities. Additionally, counting occurs at multiple locations within the same city, such as food banks, shelters, and motels, which have overlapping populations. Thirdly, since PIT counts are typically done by volunteers with little training, they are prone to bias. For example, there are problems with not correctly identifying persons contending with homelessness, even if they are sleeping rough. There have been attempts at quality control, such as using decoys (Schretzman, 2016) or plants (Hopper et al., 2008). Many studies on homeless enumeration also assume that persons contending with homelessness would come into contact with one or more agencies, including the police, shelters, or places where they receive benefits (food or advice), (Gomez et al., 1999; Williams et al., 1995). Thus, those sleeping with friends or living in squats (Darab et al., 2018) are often invisible to these methods of counting. Finally, another problem arises from the process of how individuals enter and exit the state of homelessness. Homeless populations fluctuate in size, and it is difficult to quantify the speed or extent of the fluctuations.

Therefore, there is a need to develop statistical models to reduce these errors and biases. Improved methodologies for obtaining more accurate estimates are thus crucial for a better understanding of how persons contending with homelessness behave. Capture–recapture is an ecological method used to estimate population parameters of elusive or rare wildlife populations that has also been used for human populations when a census is impractical or impossible (Bird & King, 2018). The method involves sampling the population, marking each sampled individual with a unique identifier, then releasing them back into the population. At subsequent sampling times, new samples are taken, which contain both individuals previously detected and newly detected individuals. Capture–recapture methodologies require two or more independent samples of the same population, for example, two sources that represent approximately the same population or the same source at two-time points. For example, Xu et al. (2014) and van Dam-Bates et al. (2016) used responses from multiple years from the I-track survey in Victoria (Epidemiology and Disease Control and Population Health Surveillance Unit, 2006), developed by the Public Health Agency of Canada to track changes in the prevalence of HIV and hepatitis C to estimate

population size of people who use injection drugs. From these data, it is possible to estimate parameters such as population size, survival rates, and detection probabilities. Capture–recapture has been used in some studies to enumerate homeless populations (Berry, 2007; Cowan et al., 1988; Dávid & Snijders, 2002; D’Onise et al., 2007; Gomez et al., 1999; Williams et al., 1995) and also to estimate the morbidity of homeless populations (Fisher et al., 1994).

Past studies in Plymouth (Williams et al., 1995) and Torbay (Gomez et al., 1999) used a simple two-sample approach for capture–recapture: two samples were taken over a one-week period, repeated three times within a year in different seasons. The authors used a simple model as the benefit of a more complex model would only be apparent with good data, and with the way data was collected, there were likely many violations of assumptions. Individuals were ‘tagged’ using unique identifiers of date of birth, name (or initials) and sex. Information collected included the type of current accommodation and length of stay. Berry (2007) performed another street-based capture–recapture study in Toronto, Ontario, where data was collected on 2 days in 2004 and 2005. Observation criteria in this study was quite informal, observers assessed if individuals ‘looked homeless’ before approaching by looking for criteria such as shopping carts, carrying many bags, sleeping on the sidewalk, panhandling or sifting through trash. They considered different models that allowed detection probabilities to vary by individual and by time. This study notes some issues with using street-based capture–recapture for homeless populations. Persons contending with homelessness cluster spatially and over time, for example, they may cluster around resources or shelters, and the time clustering may be because of weather, or cluster as part of a social group, thus violating the independence of individuals assumption. Street-based capture–recapture methods are the most suited for environments where the homeless have a fairly high probability of being sighted, which doesn’t hold in all urban settings.

The ability of using secondary use data that is already collected, such as electronic health records, which includes demographic information, including living situation, rather than conducting surveys or counts each time, as is the current standard, will substantially decrease costs and can also lead to more granular and frequent estimates. These data extracts also potentially remove one source of bias—where individuals are missed by counters and interviewers because they do not ‘look homeless’. To the best of our knowledge, data extracts from electronic health records have not yet been used with capture–recapture methodology to estimate homeless populations, mostly due to a lack of information about housing and living arrangements for patients. However, increasing amounts of information on patients are recorded and stored electronically, meaning that analyses such as the one within this article can be done more and more frequently in an increasing number of health authorities.

In this article, we explore the prospect of using electronic health records to estimate the annual population size of specific hidden or vulnerable populations. We propose a flexible capture–recapture model that can accommodate demographic information from electronic health records. Subsequently, we demonstrate our method through a case study on the adult homeless population of Vancouver Island, using data extracts from electronic health records provided by British Columbia’s Vancouver Island Health Authority (Island Health). We show that our method gives comparable estimates to existing PIT counts reported by the province of British Columbia.

Within Island Health, the homeless population accesses healthcare services more than four times as often as an individual who is housed. From 2013 to 2022, there is an average of approximately 13.3 events in the service system per individual, but when we only account for the homeless cohort, this number increases to 58.8 events per individual. With hidden populations such as the homeless, an estimate of the total population size, in addition to the count of individuals who interacted with the service system, is imperative. This gives the number of individuals who are potentially at risk, for example, of freezing during the winter, or of contracting an emerging infectious disease such as COVID-19. These estimates can inform healthcare providers about healthcare load, aiding in policy decisions and service planning, and can guide shelter development and budgeting for harm-reduction services. The demand on the healthcare system can be estimated by generating cohorts of homeless individuals and looking at the utilization cross-continuum. However, demand will be underestimated if homeless persons are being underestimated.

We begin in Section 2 by describing the data from Island Health. Section 3 provides details our methodology, including model assumptions. Results are presented in Section 4, and current limitations of the model are discussed in Section 5. Finally, Section 6 discusses future work using

data extracts from electronic health records to enumerate and quantify the healthcare needs of hidden populations, concluding the article.

2 Electronic health records from Island Health

Secondary use electronic health records from Island Health contains full longitudinal cross-continuum healthcare data for patients all over Vancouver Island from 1 April 2016 to 31 March 2022. This study and the use of the data was approved by the University of Victoria's and Island Health's research ethics boards via a harmonized review (REB Number H20-02451). Partial data is available from 1 January 2013 to 31 March 2016; not all patient interactions are recorded during this period, and there is incomplete documentation within interactions. Island Health upgraded its data collection protocols on 1 April 2016 and started reporting to the British Columbia Ministry of Health for error catching processes to detect missing data, and thus the data is complete after that date. In the context of capture–recapture, the change in data reporting requirements suggest low detection probabilities as well as less precise estimates from 2013 to 2016, as well as in 2022 as data is only available for the first quarter of the year.

We consider patients as being homeless if they satisfy at least one of the following criteria:

1. Designated as homeless as part of the Decampment Initiative as of April 2020 (627 Patients), or
2. if the patient ever self-reported a housing status of 'Absolutely Homeless' or 'Sheltered Homeless' within the British Columbia Ministry of Health's Minimum Reporting Requirements (MRR) for Mental Health and Substance Use (MHSU) data (1,692 Patients).

Within the MRR for MHSU, there are 17 different living arrangements, which describes where the client currently lives most of the time in the reporting period. Two of these arrangements, 'Absolutely Homeless', referring to individuals living on the streets with no form of shelter, or 'Sheltered Homeless', referring to individuals who sleep at crisis or community shelters, are used to define our cohort of interest. However, there are additional categories, such as, 'Precariously Housed', which refers to individuals with no fixed address but who have temporary shelter through family or friends, such as through couch-surfing. Currently, there are 483 individuals known to be within this housing category within Island Health. These individuals are at risk of becoming homeless, and can be easily included in the cohort of interest if warranted.

Of the 627 patients designated as part of the decampment initiative, 518 had a Forward Sortation Area (FSA, first three digits of postal code) recorded. Of those, 488, or slightly over 94%, were in Greater Victoria. Of the 627 patients designated as part of the decampment initiative, 362 had records in the MRR. Of these 362, 92 currently have a housing status of 'Absolutely Homeless' or 'Sheltered Homeless', 115 have a nonhomeless living arrangement, and 155 have an unknown living arrangement (Table 1). Although the majority of the individuals in the decampment cohort did not indicate a homeless housing status, since these individuals were identified as homeless as part of the Decampment Initiative, this takes precedence over their self-reported living arrangement.

Additionally, individuals within our dataset do not declare their housing status at every capture. As such, a large portion of individuals declare their housing status for the first time after their first interaction with Island Health. This phenomenon is evident in the breakdown of when patients were identified as homeless versus their first capture with Island Health (Table 2). For example, out of the 1,154 patients who had their first interaction with Island Health in 2016, only 413 declared themselves as homeless in 2016, and the rest declared later. There are also some individuals who are identified as homeless prior to their first interaction with Island Health as part of the decampment initiative, or if there was a referral or initial contact (such as a phone call) where patient information was collected without an interaction taking place. Patients in the decampment initiative, if they did not have an earlier housing status record, were considered to be identified as homeless in April 2020. Out of the over 1 million individuals within the Island Health system, only 73,505 have a MHSU MRR record and 24,439 indicated one of 17 possible living arrangements. Of these individuals with a reported housing status, more than half only have one record, and more than three quarters have two or fewer records. Additionally, there are no individuals who reported multiple different housing statuses. Despite how sparse the living arrangement data is

Table 1. Self-reported housing status of the 627 individuals within the decampment cohort

	Homeless	Not homeless	Unknown
In MRR	92	115	155
Not in MRR	0	0	265
Total	92	115	420

Table 2. The number of patients identified as homeless (row) and who had their first interaction (column) in each year

Year declared	Year of first interaction										Total
	2013	2014	2015	2016	2017	2018	2019	2020	2021	2022	
Homeless											
2016	9	40	61	413	1	0	0	0	0	0	524
2017	3	11	41	187	133	0	0	0	0	0	375
2018	2	4	19	105	69	84	0	0	0	0	283
2019	2	6	8	78	20	27	64	0	0	0	205
2020	6	21	29	329	121	51	51	70	7	1	686
2021	1	0	3	36	26	9	6	9	31	0	121
2022	0	0	2	6	3	1	5	1	3	5	26
Total	23	82	163	1,154	373	172	126	80	41	6	2,220

populated, we expect individuals who are unhoused or with precarious housing situations to be indicating their living arrangement (or be prompted to enter their living arrangement) much more often than those with stable living arrangements, and as such we expect the majority of the individuals who are experiencing a homeless living arrangement to update their housing status at some point. Due to the sparsity of housing status and the fact that there are no individuals with a housing status that changes through time as reflected in the data, we applied homeless status retroactively to captures. That is, if an individual ever had a housing status of ‘Absolutely Homeless’ or ‘Sheltered Homeless’, then we consider them to have been captured in any year with at least one interaction, even if that year was before the date where they first indicated their housing status. Since this condition relies on self-reported living arrangements, we do expect bias arising from missing data, which will be discussed in Section 5.

From 2013 to 2022, a total of 1,31,039 events were recorded, with emergency response (ER) being the majority with 61,877 events, which is approximately 47.2% of the recorded events. MHSU and MHSU-Addictions events have a combined total of 46,272, approximately 35.3% of recorded events (Figure 1, top). ER and MHSU-Addictions are the most used services, as they have a low barrier of entry, as opposed to other services such as MHSU or Medical Imaging which requires a referral from a physician to access. This also suggests that individuals who do not have regular access to primary care may be utilizing ER services as an alternative (Kushel et al., 2002; Shortt et al., 2010; Vohra et al., 2022). We also consider the number of patients in the homeless cohort who have had at least one interaction with each service type (Figure 1, bottom). There were 2,136 patients who had at least one interaction with ER, 1,908 with MHSU, and 1,799 with MHSU-addictions.

Self-reported demographic information—each individual’s birthdate, gender, and death date (if applicable)—was used in this analysis. The proportion of individuals who are underhoused, varies greatly by various demographic groups. There is one patient who is under the age of 18 as of 31 March 2022, and six patients, identified by the decampment initiative, who have no recorded interactions with Island Health and no recorded demographic information, including age. We removed these seven individuals to focus on modelling the population size of adults, leaving 2,220 patients observed in this cohort. In this study, we focus on the adult homeless population. This is because other than the patients provided as part of the decampment initiative, living

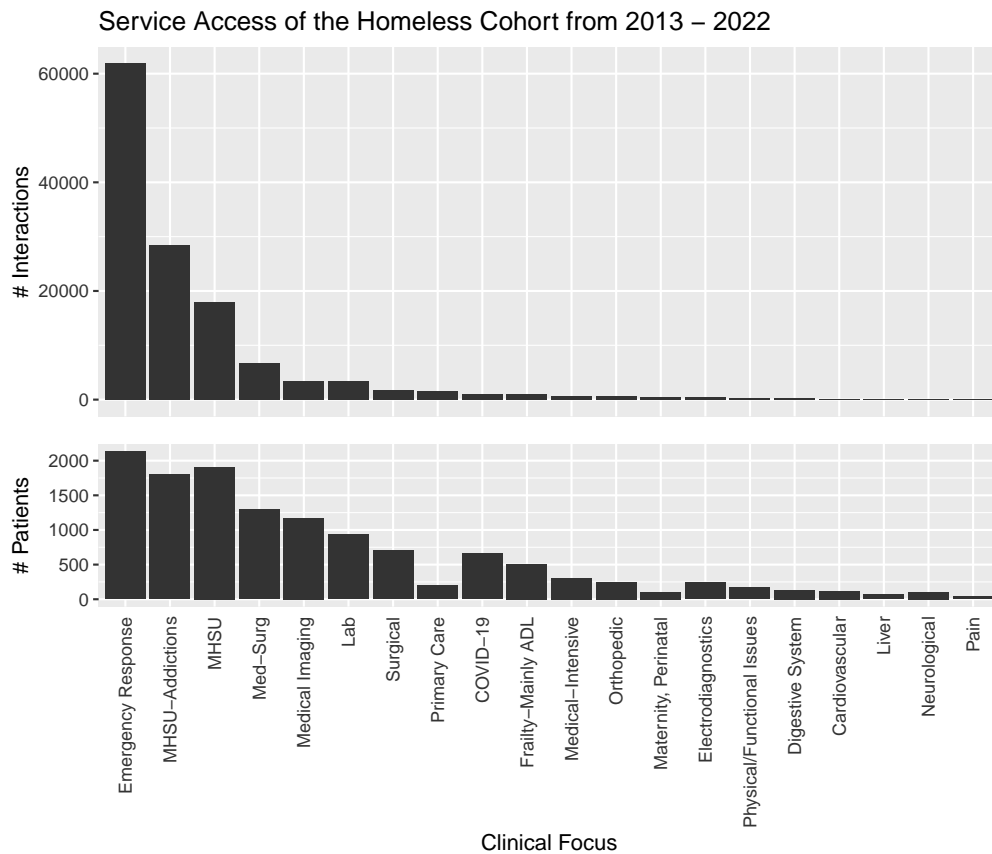


Figure 1. The frequency of service access of the top 20 most frequently used Island Health services from homeless patients (top) and number of patients who have accessed each service (bottom).

arrangement information is only available for patients who have had interactions with MHSU services. Minors who engage with MHSU services who are found to be experiencing precarious housing situations are typically immediately referred to arrangements such as foster care and thus will never have a homeless housing status. Additionally, there are recorded death dates for 112 individuals within this population, that is, 112 individuals are known to have died between 1 January 2013 and 31 March 2022.

In both the 2018 and 2020/2021 reports on homeless counts in British Columbia, published by the Homelessness Services Association of British Columbia ([Homelessness Services Association of BC, 2021](#); [Homelessness Services Association of BC and Urban Matters and BC Non-Profit Housing Association, 2018](#)), the gender breakdown was 68% male and 32% not male, which is similar to the Island Health data which has 69% male and 31% not male within those self-reported as homeless. Specifically, there are 691 females, 1,528 males, and 1 patient of unknown gender. For our purposes, we divide the patients into two groups: 1,528 ‘Male’ and 692 ‘Not male’. [Table 3](#) breaks down the population by gender and their living status as of 31 March 2022, and gives the average age of each gender for those individuals alive on 31 March 2022 (also see [Figure A1](#) in [Appendix A](#)).

The gender structure of the homeless population suggests that individuals of different genders have different risk factors for homelessness, as well as different patterns of interactions with the healthcare service system ([McGrath et al., 2023](#); [Winetrobe et al., 2017](#)). For example, women, particularly young women with children, are more likely to be undercounted by PIT counts, as they are more likely to leverage their social connections to stay out of shelters and off the streets ([Novac et al., 2002](#)). By using secondary use electronic health records, demographic information including gender has already been collected, so the impact of gender can easily be modelled to

Table 3. Summary of gender, deceased status, and mean age on 31 March 2022

Gender	# Not deceased	# Deceased	Mean age on 31 March 2022 (if not deceased)
Female	661	30	40.3
Male	1,446	82	44.3
Other	1	0	49.0

ensure that results are applicable across various healthcare authorities and can be reproduced when new data is obtained, even if the new data does not have the same gender breakdown, ensuring results are representative of all individuals.

3 Methodology

3.1 Model formulation

We use a Jolly-Seber model (see King, 2012) to estimate population size. Individuals are ‘tagged’ with their unique healthcare number and individuals are sampled at the time they access a healthcare service. Although the detection history is on a continuous timescale, rather than using a continuous-time capture–recapture model where each detection is considered an occasion (Altieri et al., 2023; Borchers et al., 2014; Chao & Lee, 1993; Hwang & Chao, 2002; Rushing, 2023; Wilson & Anderson, 1995), we aggregated detections within each year, and considered an individual detected in year t if they had at least one interaction with an Island Health service in that year. We consider each patient to be detected in a specific year if they have at least one interaction with an Island Health service, regardless of what service they used. This violates the instantaneous sampling assumption of the capture–recapture model; we will touch on this in Section 3.3.

We use data augmentation in a Bayesian state-space framework as described in (Kéry & Schaub, 2011, Chapter 10) to estimate the total and the yearly population size; also see (King, 2012; Newman et al., 2023) for more discussion on data augmentation. We have n_{obs} individuals in our dataset, and an augmentation value of M , that is, we add $M - n_{obs}$ rows of patients who have had no detections and unknown demographic information. For this study, we had $n_{obs} = 2,220$ and we augmented an additional 2,780 rows, for a total of $M = 5,000$. Each of the M individuals in the augmented population has a latent inclusion variable z_i , where $z_i = 1$ if individual i is part of the superpopulation and $z_i = 0$ if not. We assume that $z_i \sim \text{Bernoulli}(\psi)$ where ψ is the probability that an individual in the augmented population is a part of the superpopulation. In our case, the superpopulation consists of any adult individual who experienced homelessness in Vancouver Island at any point between 2013 and 2022 and was susceptible (probability > 0) to using the healthcare system at any point between 2013 and 2022. This is in contrast to when we refer to the ‘population’ at a specific time t , which refers to individuals who are susceptible to using an Island Health system at the time t .

We use the superpopulation formation of Schwarz and Arnason (1996), where annual births are modelled using a multinomial distribution from a superpopulation. That is, if an individual is part of the superpopulation, they must be ‘born’ or otherwise be recruited into the population during any one of the $K = 10$ detection occasions, where the entry probability at time t is β_t , and $\sum_{t=1}^K \beta_t = 1$. In our case, since we are only including data from adults, being ‘born’ refers to an ever-homeless individual becoming susceptible to using an Island Health healthcare service for the first time. That is, an individual could have moved to Vancouver Island, or experienced a healthcare event that necessitated accessing care for the first time. For each time t , the conditional probability of entering the population, given that the individual has not already entered, is $\frac{\beta_t}{\sum_{j=1}^K \beta_j}$.

Following the Bayesian state-space formulation described by King (2012), we let the M by K observation matrix, or detection history matrix be ω , where entry $\omega_{i,t}$ is 1 if the individual i is detected at time t and zero otherwise. That is, the $(i, t)^{th}$ element is 1 if the patient i had at least one interaction in year t (i.e. was detected in year t) and 0 otherwise. Further, the state matrix u has elements $u_{i,t} = 1$ if the individual i has entered the population by time t and has not yet left the population by time t , and 0 otherwise.

We assume

$$\omega_{i,t} | z_{i,t}, u_{i,t} \sim \text{Bernoulli}(p_{i,t} z_i u_{i,t}).$$

That is, if an individual i is not in the superpopulation ($z_i = 0$), has not yet entered the population by time t or left the population prior to time t ($u_{i,t} = 0$), the probability of detection is 0.

Further, we assume

$$u_{i,1} \sim \text{Bernoulli}(\beta_1), \text{ if } t = 1,$$

$$u_{i,t} | u_{i,j}, j = 1, \dots, t-1 \sim \text{Bernoulli} \left[\phi_{i,t-1} u_{i,t-1} + \frac{\beta_t}{\sum_{j=t}^K \beta_j} \prod_{l=1}^{t-1} (1 - u_{i,l}) \right], \text{ if } t > 1. \quad (1)$$

Thus, for an individual to be available for detection at time t , they have to either be ‘born’ or recruited at time t if they were never recruited prior to time t , or they were available for detection at time $t-1$ and survived to time t . In our case, ‘dying’ can refer to a physical death, or a permanent emigration from the service area of Island Health.

We modelled the detection probability of an individual i at time t , $p_{i,t}$ by

$$\text{logit}(p_{i,t}) = \alpha_{0,t} + \alpha_1 m_i + \alpha_2 a_{i,t} + \alpha_3 (1 - f_{i,t}), \quad (2)$$

Here, m_i is 1 if individual i is male and 0 otherwise, and $a_{i,t}$ is the age of the individual i at the sampling time t —calculated at 31st December of the year t . The behavioural-response, or first detection indicator, $f_{i,t}$ is 1 if individual i has not been detected prior to time t and 0 otherwise. That is, we model a permanent behavioural-response, rather than an immediate one where individuals can revert to their original naive state if missed, such as in (Pradel & Sanz-Aguilar, 2012). We expect the existence of a behavioural-response to exist as individuals who have utilized services are more likely to re-visit, for example, for follow-up care, or to utilize other services in the Island Health system after being made aware of them through, for example, referrals. The baseline detection rate is time-varying, but the effect of age, gender, and behavioural-response are not. We also modelled the survival probability of an individual i from t to $t+1$, $\phi_{i,t}$ dependent on both age and gender, as

$$\text{logit}(\phi_{i,t}) = a_{4,t} + \alpha_5 m_i + \alpha_6 a_{i,t}. \quad (3)$$

Since demographic information is only known for detected individuals and not augmented individuals, we treated age and gender of individual i at year 1 as latent random variables if unknown, similar to how age is treated in (Hostetter et al., 2021). Specifically, we have

$$m_i \sim \text{Bernoulli}(\pi_m) \text{ and}$$

$$a_{i,0} \sim \text{Truncated Normal}(\mu_a, \sigma_a^2, 0, 100),$$

where $a_{i,0}$ is the age for individual i at 1 January 2013 and each individual aging one year per capture occasion. We used standard model selection techniques, discussed in Section 3.2, to determine which covariates, of age, gender, or behavioural-response, to include in equations (2) and (3).

We used vague noninformative priors (Banner et al., 2020) for our population parameters as follows:

$$\psi, \text{expit}(\alpha_{0,t}), \text{expit}(\alpha_{4,t}) \sim \text{Beta}(1, 1),$$

$$\beta_{1:K} \sim \text{Dirichlet}(1, \dots, 1),$$

$$\alpha_1, \alpha_2, \alpha_3, \alpha_5, \alpha_6 \sim \text{Normal}(0, \sigma = 30),$$

For demographic information, we used the following priors:

$$\begin{aligned}\pi_m &\sim \text{Beta}(1, 1), \\ \mu_a &\sim \text{Normal}(40, 20), \text{ and} \\ \sigma_a &\sim \text{Inverse Gamma}(5, 10).\end{aligned}$$

Here, ‘logit’ and ‘expit’ are the logistic and inverse logistic transformation, respectively.

Finally, our estimates of interest are the annual population size, $N_t = \sum_{i=1}^M u_{i,t} z_i$, and we compare the estimates of the annual population size of N_t with the PIT counts. This parameterization, where ‘alive in year t ’ ($u_{i,t} z_i$) is modelled as the product of two variables, follows (Kéry & Schaub, 2011, Chapter 10.3.3). An alternative equivalent parameterization where the two are modelled together is also possible. In this case, equation (1) would be modified such that z_i is multiplied by the Bernoulli probability, and then $u_{i,t}$ would represent ‘alive in year t ’.

Parameter nonidentifiability, or when parameters are not estimable, occurs in the Jolly-Seber model (Schwarz & Arnason, 1996) and other complex ecological models. In particular, without individual covariates, $\phi_{i,K-1}$, $p_{i,K}$, the last survival and the last detection probability, are not identifiable individually—it is only possible to estimate their product. Likewise, β_1 and $p_{i,1}$, the first detection probability and the first recruitment probability, are not identifiable individually, and it is only possible to estimate their product. Since they are not individually estimable, we remove this source of extra variability by setting $p_{i,1} = \sum_{i=1}^M \omega_{i,1} / \sum_{i=1}^M u_{i,1}$, which is 1 if u is populated using the information from ω for all i . We also set $\phi_{i,K-1} = 1 - (\sum_{i=1}^M (1 - u_{i,K}) - s \sum_{i=1}^M u_{i,K-1}) / M$, which is simply 1—proportion of individuals who were alive at the time $K - 1$ who were known to have died before time K , for all i . Given these constraints to the model, the population size for the first and last years should be discarded—setting $p_{i,1}$ to be close to 1 will inflate the estimates of β_1 and thus deflate the estimates of N_1 . Likewise, setting $\phi_{i,K-1}$ close to 1 inflates the estimates of N_K . As such, estimates for the years 2013 and 2022 are not reported. See (Kéry & Schaub, 2011, Chapter 7.9 and Chapter 10.6) for further discussion on parameter identifiability.

For inference and parameter estimation, we use a Bayesian approach with Markov chain Monte Carlo (MCMC); we used NIMBLE v0.13.0 (de Valpine et al., 2022) software, accessed via R v4.1.3 (R Core Team, 2022) to implement the MCMC simulations and to fit models. Initial values and data must be passed to NIMBLE. In particular, values of $u_{i,t}$ and z_i must be passed to prevent initialization errors and to speed up convergence. If individual i is known to be in the population, i.e. there are records for them in the Island Health database, then z_i is set to 1. Otherwise, z_i is set to NA, as it is unknown. Similarly, if individual i has interactions at time t_1 and t_2 , then we set $u_{i,t}$ to be 1 for all $t_1 \leq t \leq t_2$, and otherwise NA. That is, $u_{i,t}$ is 1 between the first and last time seen (inclusive) and NA otherwise.

3.2 Model selection

A limitation of our model is the computational demands, especially when the number of individuals in the population is large. All model fitting was performed on a virtual machine environment at Island Health with an Intel(R) Xeon(R) Gold 5118 CPU @ 2.30GHz processor with 15.6 GB of memory, as due to data privacy constraints, exporting the data to run on an external machine or server is not possible. Thus, due to these computational limitations, we employed a two-step process for model selection and parameter estimation: preliminary model selection was done with a single chain with a smaller number of iterations and then the selected final model was fit using three longer chains. If computational restrictions were not present, then this two-step process would not be required. For our preliminary model selection step, we used one MCMC chain with 50,000 iterations, with the first 20,000 being discarded for burn-in. The final selected model was then run with three chains, 1,50,000 iterations each chain, and the first 50,000 iterations were discarded for burn-in. The burn-in amount and number of iterations needed for convergence was chosen manually using convergence plots. Convergence was also confirmed with the Gelman–Rubin statistic (Gelman & Rubin, 1992), which was less than 1.1 for almost all the parameters of interest, suggesting convergence (Brooks & Gelman, 1998). We did not thin the observations (Link & Eaton, 2012) as we did not have problems with memory.

Due to systematic changes in how data was collected and stored as well as underlying changes to the risks of homelessness over the years, only time-varying survival and detection probabilities were considered. Thus, all models we considered have time-varying baseline survival and detection probabilities. Notation for referring to models is from (Lebreton et al., 1992). In our notation, t denotes time-varying, g denotes gender, a denotes age, and b denotes behavioural-response. An additional model with time-varying behavioural-response, gender, and age, denoted ϕ_{agt} , p_{agb^*t} was also considered. That is, instead of equation (2), we modelled the capture probabilities with $\text{logit}(p_{i,t}) = \alpha_{0,t} + \alpha_1 m_i + \alpha_2 a_{i,t} + \alpha_3 (1 - f_{i,t})$. In total, we considered the following nine combinations of covariates:

1. Time-varying behavioural-response, gender, age (full model with time-varying behavioural-response) (ϕ_{agt}, p_{agb^*t})
2. Behavioural-response, gender, age (full model) (ϕ_{agt}, p_{agbt}) ,
3. Gender, age (ϕ_{agt}, p_{agt}) ,
4. Behavioural-response, age (ϕ_{at}, p_{abt}) ,
5. Behavioural-response, gender (ϕ_{gt}, p_{bgt}) ,
6. Age (ϕ_{at}, p_{at}) ,
7. Gender (ϕ_{gt}, p_{gt}) ,
8. Behavioural-response (ϕ_t, p_{bt}) ,
9. No covariates (ϕ_t, p_t) .

We used widely applicable information criterion (WAIC) (Gelman et al., 2014; Watanabe & Opper, 2010) for model selection. WAIC in NIMBLE is calculated conditioning on all parent nodes of the data nodes. Since demographic information is treated as a latent random variable and thus part of the model; the uncertainty from populating demographic information contributes to the WAIC for all models except ϕ_t, p_t . For WAIC to be comparable between models, demographic information must be populated even for the models that do not require it, as otherwise, the data within the model are essentially being changed. That is, age and gender are included in the data for all models, regardless of their inclusion in equations (2) and (3). If data augmentation is not used and there are no individuals with missing demographic information, or missing demographic information is imputed prior to model fitting, this step is not necessary, as there would be no uncertainty from populating missing demographic information contributing to the likelihood and thus WAIC. An alternative is to compare WAIC using the conditional likelihood that ignores the contribution to the likelihood that the random unobserved covariate values give, by marginalizing over those nodes using the `controlWAIC()` argument. However, this increases computation time and due to constraints in computational power, we did not opt for this method.

3.3 Assumptions of the Jolly-Seber model

There were several assumptions of the Jolly-Seber model that were violated in our scenario and might have an effect on our estimates.

1. Individuals who are detected have the same patterns of behaviours (detection and survival) as those who are not detected. That is, they are representative of our population of interest. However, this does not hold in reality if we consider our population to be all the persons contending with homelessness within Vancouver Island—there are individuals who will likely never use Island Health services for various reasons, including distrust of the system, or due to systematic barriers to healthcare access. Thus, we must consider our population to be that of adult homeless who are within Island Health jurisdiction and will potentially utilize Island Health resources.
2. Markers do not affect the behaviour of the marked individuals. This assumption is likely violated, as there may be a behavioural-response, where there is a change in the behaviour of subjects' after being 'trapped', or detected for the first time. Individuals can exhibit either trap-happiness or trap-shyness, when they are more or less likely to return after the first detection, respectively. For example, persons contending with homelessness who utilize Island Health services may be incentivized to return. For example, an individual staying overnight

at the sobering assessment centre may be likely to return to have a bed for the night. In addition, follow-up care or long-term treatment for specific health concerns can also cause an individual to be more likely to return. As such, we include behavioural-response in our model, to account for this.

3. Marking is accurate and reliable—no lost or misread marks. Since we are using each individual's unique healthcare number as their mark, we can assume that marks are not lost or misread. However, as we will discuss in Section 5, an issue arises where there exists individuals where we have the capture history, but we do not know to include them in our population. This is because there is a lag between an individual's first recorded interaction and the time that they are declared homeless. This primarily affects estimates for later years, individuals who are captured in earlier years have more time to indicate their homeless housing status and have their capture history included. In ecology, a similar, but not identical scenario may arise if sexes of a certain species look identical and are thus indistinguishable until females are seen with young or spawning, such as the Dolly Varden trout (Gallagher et al., 2019), but an interest lies in estimating the population size of females.
4. Every marked individual has the same probability of detection. We address heterogeneity in detection probability by stratifying using demographic information such as age and sex, due to the differing underlying morbidity rates.
5. Individuals are independent. This assumption is likely violated due to the existence of social groups. With data privacy concerns, there is no way to determine family units within the current data; however, we may be able to account for social groups with data that is not anonymized. In addition, even with the existence of social groups, individuals within the groups typically access healthcare services individually, as they are accessed on an as-needed basis. This is different from ecological studies, where family units are often trapped or detected together (Draghici et al., 2021).
6. Sampling is instantaneous. This condition is violated, as we have continuous sampling occasions. This is one of the most commonly violated assumptions, and violation of this assumption results in heterogeneity, as members of the marked population do not have the same probability of survival between sampling occasions. For example, an individual who had an interaction with the healthcare system on 1st January of 1 year will need to survive for 12 months longer than an individual who had an interaction on 31st December of that year to survive to be detected the following year (Smith & Anderson, 1987). However, due to the low death rates in this population, this is likely not a concern.

Heterogeneity and lack of independence arising from the violation of these assumptions is collectively considered to be overdispersion, which, if unaccounted for, can result in underestimated variances for parameters and thus undercoverage in the credible intervals (Lindberg & Rexstad, 2006). Although not all assumptions of the model are met, violation of the assumptions is generally not deleterious (Lindberg & Rexstad, 2006), and the model still gives us a good understanding of the population of interest.

4 Results

Table 4 gives the WAIC values and the Bayesian posterior predictive p -values Gelman et al. (1996) associated with the Freeman–Tukey statistic for evaluating goodness of fit (see (Kéry & Schaub, 2011, Chapter 7.10) and Brooks et al. (2000) for more discussion and an implementation) for each of the above models. There is no evidence of a lack of fit for any of the models, as the p -values are approximately 0.5 for all models. Under the null hypothesis that the model used is in fact the data-generating model, p -values close to 0 or 1 indicate a lack of fit. In our case, we begin with the saturated model, as including demographic information and behavioural-response have been shown to improve precision and bias of estimates, and we perform model selection to determine if any covariates should be removed, and it was found that removing them was not necessary. As such, based on our model selection criteria, the model we proceeded with is ϕ_{agt} , p_{agbt} , with behavioural-response, gender, age incorporated, and time-varying baseline survival and capture probabilities.

Table 4. Model selection results using Island Health electronic health records from 2013 to 2022

Covariates included	Model	WAIC	Δ WAIC	GOF <i>p</i> -value
Time-varying behavioural-response, gender, age	ϕ_{agt}, p_{agbst}	42,683	1,019	0.50
Behavioural-response, gender, age	ϕ_{agt}, p_{agbt}	41,664	0	0.50
Gender, age	ϕ_{agt}, p_{agt}	42,125	461	0.50
Behavioural-response, age	ϕ_{at}, p_{abt}	42,096	432	0.49
Behavioural-response, gender	ϕ_{at}, p_{at}	42,377	713	0.50
Age	ϕ_{gt}, p_{bgt}	42,400	736	0.49
Gender	ϕ_{gt}, p_{gt}	41,997	333	0.50
Behavioural-response	ϕ_t, p_{bt}	42,915	1,251	0.50
No covariates	ϕ_t, p_t	42,232	568	0.50

Note. Δ WAIC denotes the difference in WAIC between that model and the minimum WAIC.

Past 2016, when data is better recorded, the estimates between the nine methods are quite similar, regardless of which model is used (Figure 2, Table C1 in Appendix C). In the earlier years, including demographic information, and notably, including behavioural-response, prevents underestimation of the annual population size. This is important, as we expect the significant difference in the number of observed individuals before and after 2016 to be due to the different data-reporting requirements, resulting in less interactions being recorded, meaning lower detection probabilities, and not due to a significant change in the population size of the homeless population. In all years, including behavioural-response in the model prevents underestimation of the annual population size (triangle markers). This suggests that individuals are indeed more likely to return after their first interaction with the healthcare system. Indeed, not accounting for a trap-happy response can lead to bias; underestimating population size in the Jolly-Seber model (Nichols et al., 1984; Pradel & Sanz-Aguilar, 2012), and we are correct in our decision to model a behavioural-response.

Due to the violation of Assumption (3) in Section 3.3, we observe an artificial drop in population size estimates in the later years. This is discussed in detail in Section 5, but we only display results prior to and including 2019 within the main body of this article. Stability of the homeless population size estimates greatly increases after 2016 (when more stringent data reporting requirements led to complete data being collected). Thus, as we expect more accurate and complete medical records and patient interaction records to be collected going forward, more accurate estimates of the homeless population or other populations will be easier to acquire using these records. Estimates from all years are available in Appendix C.

PIT counts are available for 2018 (1884) and 2020/2021 (2149) for Vancouver Island. Details about which jurisdictions were included for these reported PIT counts can be found in Appendix B. We note that these counts, which enumerate sheltered and unsheltered homeless people, cannot be directly compared with our estimates, for two main reasons: (1) The PIT counts include children, and (2) the PIT counts include those who have not, and will never need to interact with Island Health services. Additionally, the majority of the patients designated as homeless in the Island Health cohort were designated by the MHSU MRR, meaning that they have used a MHSU service in the past. However, there were very few children (2.86% in 2018 and 2.56% in 2020/2021) in the PIT counts in all of British Columbia, and a substantial amount of the individuals surveyed for the PIT counts reported having a mental health issue (44% in 2018 and 51% in 2020/2021) or an addiction (56% in 2018 and 67% in 2020/2021). Thus, although we are not enumerating exactly the same population and neither population is a subset of the other, we expect the populations reported in the PIT counts and our estimates to have significant overlap. The PIT counts are a useful comparison, providing a method to validate our estimates.

From our results from our final model, ϕ_{agt}, p_{agbt} (Figure 2, Table C2 in Appendix C), the capture–recapture model gives estimates that are sensible, implying that enumeration of persons contending with homelessness can be done with electronic health records instead of costly PIT counts in the future. We also note that the point estimates (and 95% CIs) of the coefficients for gender are

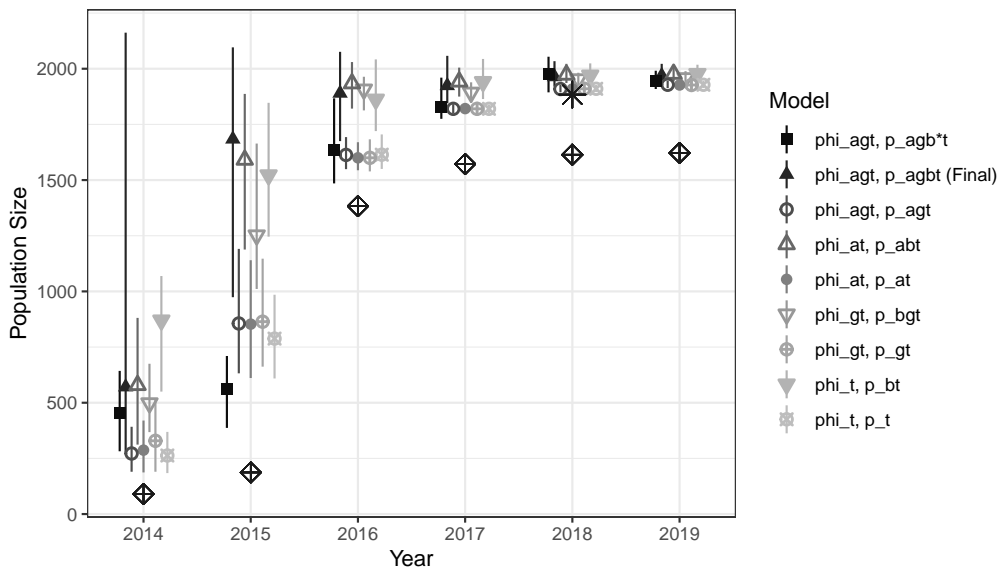


Figure 2. Annual population size estimates for people experiencing homelessness within Island Health with 95% credible intervals from 2014 to 2019 using the nine models specified above. Note that model ϕ_{agt}, p_{agbt} (upright filled triangle) was chosen to be the best model and results displayed for this model were from the final run with three chains, 150,000 iterations each. For the rest of the models, results from the model selection phase, using one chain with 50,000 iterations, are displayed. The 2018 PIT count is shown as the star. Number of observed patients per year is shown as diamonds.

−0.08 (−0.28, 0.12) in the survival probability (α_4 in equation (3)) and −0.11 (−0.22, 0.00) for the capture probability (α_1 in equation (2)). The estimate of the coefficient for the existence of a prior capture is 1.00 (0.39, 1.52) (α_3 in equation (2)). This suggests that males are less likely to access care, and that once an individual has entered the Island Health system, they are much more likely to return for care in subsequent years.

We now note some additional use cases of our model. If interest is in specific jurisdictions, attention can be focused on them by filtering patients in the areas of interest. This can be especially useful if interest lies in jurisdictions that are more remote, or sparsely populated, where it is not feasible to perform PIT counts. Segmentation by other demographic information is also possible, but not demonstrated within this article. We look at Greater Victoria, as defined using 17 FSAs. Estimates of yearly and total population size are plotted in Figure 3, with values from Table C3 in Appendix C. The 2018 PIT count for Greater Victoria was 931.

Since cross-continuum data is not available for all healthcare authorities, we also explored using ER data only as a suitable alternative to estimate the size of this population. The advantage here is that the ER is a single point of service at each hospital, and thus, data does not need to be aggregated between a variety of service providers. We note that the populations enumerated by only ER interactions and using all interactions are not exactly the same. For example, there are individuals who may access MHSU services through a caseworker or residential care, but may never have an ER interaction. However, since we know the homeless population disproportionately accesses ER services at higher rates and the majority of interactions with the healthcare system is through ER services due to the low barrier to access, we expect the two populations to have significant overlap. In our case, 2,136 out of the 2,220 (96.2%) individuals within the homeless cohort had at least 1 ER interaction. Additionally, using solely ER interactions also aligns more closely with the assumption of independence between individuals: While health service interactions are typically as-needed, nonurgent services such as vaccine clinics can still be influenced by social groupings coordinating appointments. In contrast, ER services are genuinely accessed on an as-needed basis. Another point to note is that within Island Health, housing status is not typically collected at ER interactions—thus, if one wishes to use only ER interactions, then housing status must be collected or provided otherwise.

Using the model ϕ_{agt}, p_{agbt} , we found that population size estimates using only ER interactions give comparable point estimates to those using all interactions, albeit with wider credible intervals

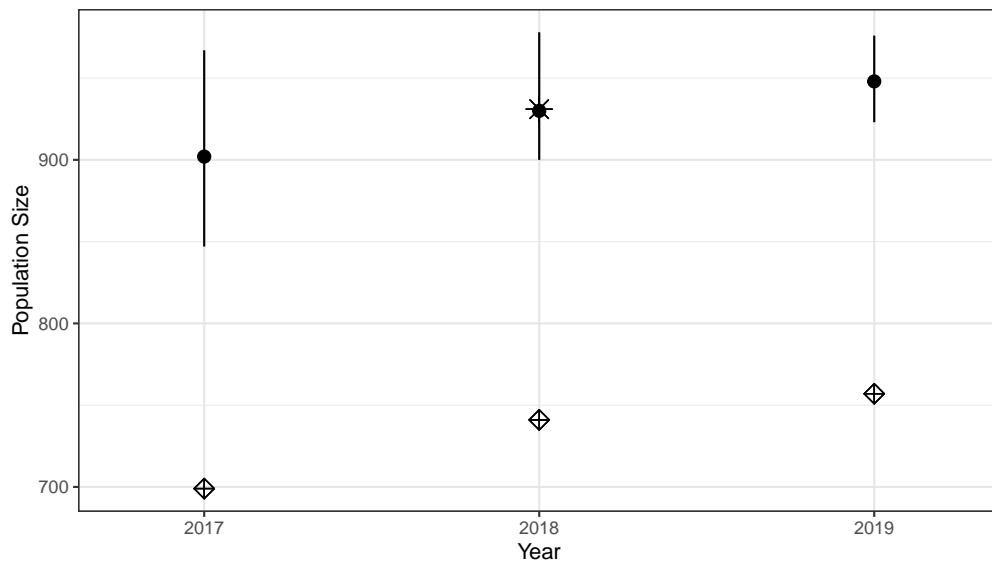


Figure 3. Annual population size estimates for Greater Victoria from 2017 to 2021 using model ϕ_{agt}, p_{agbt} : Medians and 95% credible intervals. The 2018 PIT count is shown as the star. Number of observed patients are shown as diamonds.

(Figure 4, Tables C2 and C4 in Appendix C). Since ER data was only recorded from 1 April 2016 onward, only estimates from 2017 to 2021 are reported.

5 Study limitations

As mentioned, one curious result was that we found the yearly population size estimates to decrease over time, see Table C2 in Appendix C, whereas the 2020/2021 PIT count is higher than the 2018 count. We expect the homeless population size to go up during the COVID-19 impacted years, as COVID-19 has exacerbated housing insecurity in many cities within Canada: 14% of respondents in British Columbia have identified COVID-19 as a reason for their most recent housing loss (Homelessness Services Association of BC, 2021). This is likely the main driver behind the increase between the two PIT counts. We find two possible explanations for this phenomenon.

The first is that this period overlaps with COVID-19, during which Island Health experienced a drastic decrease in service usage after March 2020. The British Columbia government purchased hotels to use as shelters (CBC News, 2020; Chan, 2021). This increased access to shelter may have contributed to a decrease in medical needs for this population. Additionally, COVID-19 reduced access to numerous healthcare services, notably ER services (Canadian Institute for Health Information, 2021) across Canada, due to reduced availability of ER beds and interpretations of public health restrictions or fears of contracting COVID-19 while accessing healthcare services. The Jolly-Seber model is unable to differentiate between nonsurvival and noncapture of individuals in those years, as without post-COVID-19 data, the model cannot differentiate between those individuals who simply were not captured and those who left the population due to being unable or unwilling to access a healthcare service. As post-COVID-19 data is not available, it is unknown if the emigration is permanent (which would result in reduced survival in the population) or temporary (which would result in reduced captures for the COVID-19-impacted years).

Secondly, as mentioned in Section 2, many individuals declare their housing status for the first time after their first interaction with Island Health. This suggests that there are individuals entering the Island Health healthcare system who are homeless, but who have not yet indicated their housing status, and despite having the capture history for these patients, we do not include them in our population as we do not know that they are homeless. This likely in part explains the decrease in estimates over time and why we underestimate the 2020/2021 PIT count—we do not sufficiently account for

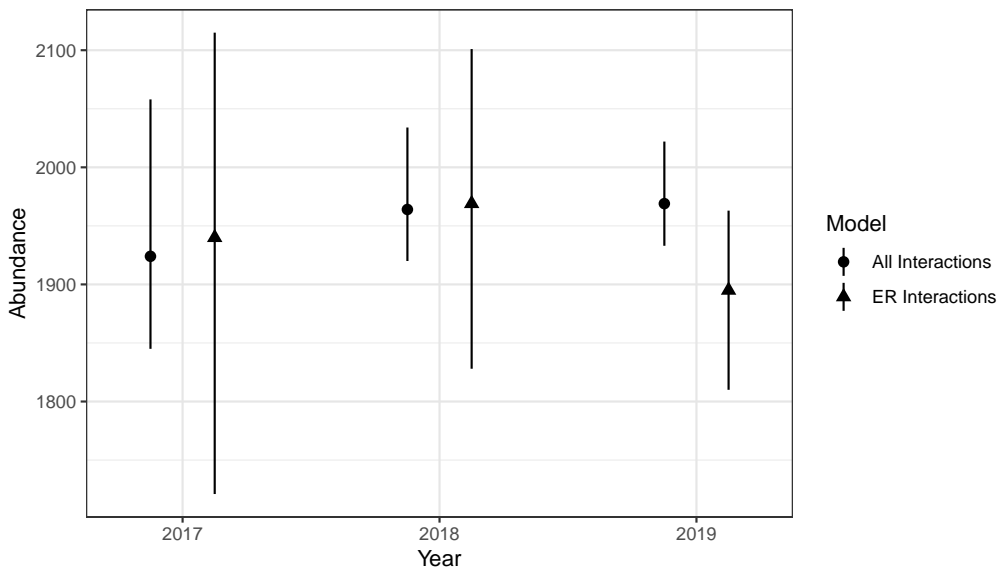


Figure 4. Annual population size estimates for Island Health Authority using model $\phi_{agtr} p_{agbr}$: Medians, along with 95% credible intervals. Estimates and credible intervals generated using all interactions (circle) and ER interactions (triangle) are shown.

individuals entering the population. Estimates in earlier years give more time for patients to update their housing status and retroactively be included in our sample, resulting in less bias in these estimates. A patient who had their first interaction in 2016 has six years to indicate their housing status, whereas a patient with their first interaction in 2022 has less than one.

The bias caused by this lag between the first interaction and declaration of housing status causes the assumption of ‘Marking is accurate and reliable—no lost or misread marks’. to be violated, but not in a way that is seen in ecological studies. Studies quantifying the bias arising from tag loss (Cai et al., 2021; Cowen & Schwarz, 2006; Cowen et al., 2009; Malcolm-White et al., 2020), tag misidentification, for example, in photographs (Morrison et al., 2011) or in genetic tags (Lukacs & Burnham, 2005; Yoshizaki et al., 2011), where multiple tags are assigned to the same individual or a tag is read as an already existing mark, or unreported recovered tags from deceased individuals in mark-recovery studies (Brook et al., 2022) exist in literature, and it has been shown that even small tag loss rates of 5% can cause severe overestimation of population size and inflated standard errors (Lee, 2002). However, in this study, we are correctly recording each individual’s capture history, but there are cases where we do not know that an individual should be included as part of the population until after their initial capture.

Since post-COVID-19 data is not currently available, it is not possible to separate the effects of the lag in reporting with the effects of COVID-19. Thus, for any modelling of the reporting lag or the effect of the reporting lag to be possible, we must have ample post-COVID-19 interaction history.

If using the model as described in this article, a solution that works in theory is to include interactions as capture histories for all patients, not just those who are recorded as homeless. Then, similarly to data augmentation, a parameter to indicate if that individual is in the population (i.e. homeless or not) can be modelled, possibly using their capture histories or complete pattern of service utilization and demographic information. For example, we may be able to identify proxies for a homeless living arrangement, such as extended stays in sobering centres. This would also allow for modelling a changing homeless status through time, as each patient’s patterns of service utilization will change over time. As an alternative to the model described within this article, we would now have a mixture of declared and undeclared homeless individuals, a mixture model on detection probabilities, such as those described within (Pledger & Phillpot, 2008) could also be used. However, in practice, this method of including capture histories for all patients may not be feasible, due to computational limitations and data privacy. The current computational

environment where Island Health data is allowed to be stored, accessed, and processed, cannot process much more than 5,000 rows, including both patients and augmented individuals, with the model described within this article. Thus, including all patients' capture histories when fitting the model and performing MCMC cannot be easily applied in practice, especially using data extracts from healthcare records, as many jurisdictions have similar data privacy safeguards in place. Thus, different or more efficient statistical methods to quantify the underestimation resulting from lag in reporting homelessness must be developed, especially since having up-to-date estimates of the population size is imperative for service planning.

6 Conclusions and future work

In this article, we used data extracts from electronic health records to estimate the annual homeless population size within the Island Health region in British Columbia. We note again that the estimates are not an estimate of the total population size of the Vancouver Island homeless population, but rather the population of adults who may interact with healthcare services within Island Health. We proposed a flexible capture–recapture model, which gives estimates that are sensible compared to existing PIT counts, meaning that enumeration of persons contending with homelessness can be done with reduced cost and increased frequency compared with PIT counts, contingent on the quality of the data. Using this model, we demonstrated how to perform model selection for the inclusion of covariates using WAIC. Our model can be used at different time scales and for specific segments of the population, allowing for better planning of health services. Finally, we explored using only ER interactions instead of full cross-continuum data, and found that estimates using only ER interactions give comparable point estimates to those using all interactions, albeit with wider credible intervals, suggesting that this model can be extended to other health jurisdictions which may not have records for all cross-continuum healthcare interactions, but have ER interaction records.

The issues discussed in Section 5 motivates us to find ways to use the PIT counts to correct for this underestimation. This would require the use of complex integrated data techniques (Abadi et al., 2010; Besbeas et al., 2005, 2002, 2003; Borysiewicz et al., 2009; Chandler & Clark, 2014; Schaub & Abadi, 2011), which would allow us to use multiple dependent data sources in the model. The PIT counts referenced within this article are an undercount and represent only those individuals who were available at the site of the count during the 24-hr count period. Thus, they represent a lower bound on the number of individuals in our population at that time, and can be used to correct for any individuals who have not yet indicated their housing status to Island Health. Using integrated data modelling techniques will also allow us to better understand the current data source, which is technically comprised of two sources—self-reported housing status from the MRR, and designated homeless as part of the decampment initiative—with an overlap. In this case, we cannot assume independence between these two groups, nor can we assume that these samples are from a closed system.

Additional work on addressing the healthcare needs of the homeless population by using electronic health records can be done by expanding our methodology to achieve more accurate estimates of the population size, which would require more data than we currently have. This study illustrates a potential use of electronic health records, and provides the organizations who are collecting the data—the regional health authorities and the provincial Ministries of Health—with an incentive to upgrade the accuracy and frequency of data collection. We stress the importance of healthcare authorities collecting more granular data that tracks each individual's housing status through time, and for this information to be collected on intake at a greater variety of services, such as at ER interactions, so data can be collected more often for more patients.

Firstly, homelessness is a continuum and can be seen as a latent state of an individual, for example, a homeless state can be viewed as a Markov process. Many individuals experience homelessness as a fluid experience (Czechowski et al., 2022), where their housing circumstances change dramatically and frequently (Gaetz et al., 2012), and it is estimated that up to 11% of homeless individuals experience episodic or cyclic homelessness (Aubry et al., 2013), where they experience multiple homeless episodes for relatively short periods of time. The homeless population constitutes an open system, as people move in and out of the state of homelessness or between the different modes of homelessness, and are also spatially mobile as they move to and from other

neighbourhoods/enumeration areas. Typical open population models in ecological studies account for an underlying birth and death process and migration process, an additional process where individuals can enter and exit homelessness multiple times must be considered. These considerations were not in the scope of this article, as currently, the housing status for each individual is not updated if it changes, making this impossible to model presently. If the data required becomes available, future work in this area involves modelling the way individuals transition between different housing states through time using a multi-state model, similarly to accounting for temporary emigration; see for example, (Kéry & Schaub, 2011, Chapter 9). Alternatively, if we can identify patterns of service utilization (service class or a combination of services classes) that act as a proxy or predictors (i.e. for a classification algorithm) for the homelessness state through time, we can develop a multi-event capture–recapture model which models uncertainty in state assignment (Pradel, 2005). In particular, we are interested in transitions from being housed to being homeless, as the effects of being homeless can be permanent—the supports an individual needs often do not go away after becoming housed. Healthcare-related risks of being homeless, as well as changes in patterns of service use before and after a homeless individual becomes housed, should be explored as well.

Secondly, our future direction also involves modelling not only the population size, but the load on the healthcare system by the homeless. There is large variance in how much the homeless population uses various healthcare services, and as such, the load on the system cannot be calculated using a simple count of the number of services used. Approximately one in four homeless have been reported to be unable to receive necessary medical care when required (Kushel et al., 2002, 2001; Vohra et al., 2022), and this research may identify and quantify the gaps in care. Infectious disease spread and the prevalence of emerging infectious diseases, including COVID-19 and pandemic response for the homeless population, can also be explored.

Finally, these methods could be used to explore homeless populations interacting with healthcare systems in other cities, and results can be rolled out with data from other jurisdictions for similar studies. This includes jurisdictions not in Canada, where rates of access to care may be lower, for example, in San Francisco, where only 40% of homeless respondents had one or more ER interactions in the previous year (Kushel et al., 2002).

Funding

G.D. gratefully acknowledges the financial support of the Canadian Statistical Sciences Institute (CANSSI) via the CANSSI Distinguished Postdoctoral Fellowship program.

Conflict of interests: None declared.

Data availability

Data and computational resources used in this project were provided by the Vancouver Island Health Authority (Island Health) in British Columbia, Canada. The data used within this project cannot be made publicly available due to Island Health privacy restrictions. R code to implement the model ϕ_{agt}, p_{agt} is available in the [online supplementary material](#).

Author contributions statement

G.D. conceived and developed the models for analysis and wrote the article. S.R. and K.M. provided the data. P.B. and L.C. supervised the project.

Supplementary material

[Supplementary material](#) is available online at *Journal of the Royal Statistical Society: Series A*.

Appendix A. Age and gender distribution

[Figure A1](#) shows the distribution of ages and gender of the homeless population as of 31 March 2022 for those individuals that are not known to be deceased at that date. Each bar corresponds with 1 year.

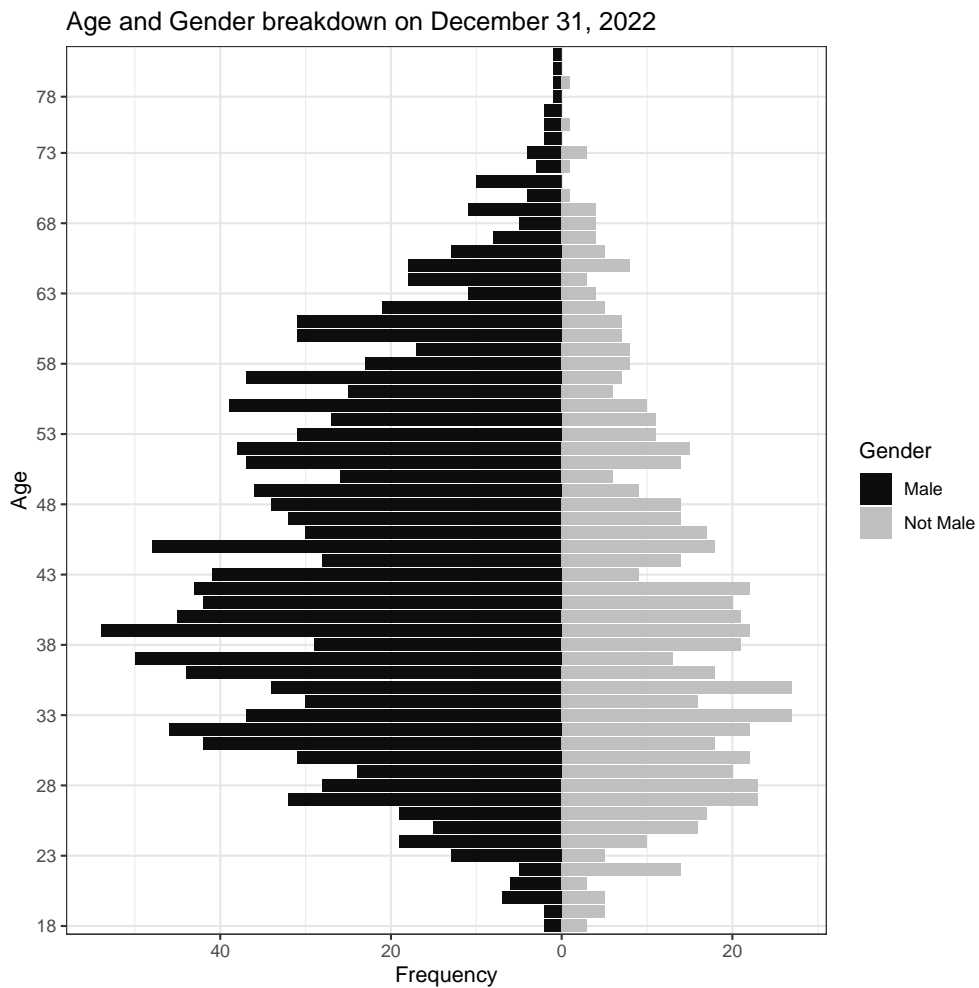


Figure A1. Distribution of ages and gender on 31 March 2022

Appendix B. PIT counts for vancouver Island

PIT counts for homeless populations are available for 2018 ([Homelessness Services Association of BC and Urban Matters and BC Non-Profit Housing Association, 2018](#)) and 2020/2021 ([Homelessness Services Association of BC, 2021](#)) for the province of British Columbia. These reports aggregate results from provincially funded counts, federally funded counts, and independent counts. We can compare our estimates with PIT counts for the areas within the Island Health jurisdiction. This spans several PIT counts—in particular, we consider the counts in eight communities: Campbell River, Comox Valley, Parksville/Qualicum, Port Alberni, Nanaimo, Duncan, Greater Victoria, and Salt Spring Island. This is not an exhaustive list of all areas on Vancouver Island; however, other jurisdictions are less populated and would not contribute significantly to the total count on the island. The PIT counts are summarized in [Table B1](#): 1,884 in 2018 and 2,149 in 2020/2021.

Table B1. PIT counts for several locations within Island Health for years 2018 and 2020/2021 (Homelessness Services Association of BC, 2021; Homelessness Services Association of BC and Urban Matters and BC Non-Profit Housing Association, 2018)

Area	2018 count	2020/2021 count
Campbell river	81	116
Comox valley	117	132
Parksville/Qualicum	42	87
Port Alberni	147	125
Nanimo	301	406
Duncan	150	129
Greater Victoria	931	1,008
Salt Spring Island	115	146
Total	1,884	2,149

Appendix C. Annual estimates

Annual population size estimates from the various models are presented in [Tables C1, C2, C3, and C4](#).

Table C1. Annual population size estimates and 95% CIs from 2014 to 2021 using the models from the model selection phase

Year	ϕ_{agt}, P_{agb*st}	ϕ_{agt}, P_{agbt}	ϕ_{agt}, P_{agt}	ϕ_{at}, P_{abt}	ϕ_{at}, P_{at}
2014	456 (282, 643)	2,465 (2,009, 2,738)	272 (190, 392)	579 (312, 881)	287 (187, 420)
2015	563 (387, 710)	2,045 (1,886, 2,215)	856 (632, 1,191)	1,591 (1,188, 1,887)	853 (611, 1,140)
2016	1,636 (1,485, 1,867)	2,066 (1,974, 2,139)	1,613 (1,549, 1,693)	1,934 (1,821, 2,030)	1,600 (1,544, 1,670)
2017	1,828 (1,775, 1,960)	2,038 (1,948, 2,115)	1,820 (1,801, 1,841)	1,942 (1,875, 2,005)	1,821 (1,802, 1,843)
2018	1,975 (1,894, 2,054)	2,027 (1,955, 2,085)	1,909 (1,894, 1,927)	1,973 (1,936, 2,022)	1,909 (1,894, 1,927)
2019	1,943 (1,909, 1,991)	2,007 (1,972, 2,049)	1,928 (1,913, 1,946)	1,977 (1,943, 2,014)	1,927 (1,912, 1,944)
2020	1,850 (1,825, 1,884)	1,878 (1,854, 1,902)	1,839 (1,825, 1,856)	1,862 (1,841, 1,887)	1,840 (1,825, 1,856)
2021	1,644 (1,625, 1,681)	1,651 (1,629, 1,677)	1,635 (1,621, 1,652)	1,644 (1,626, 1,663)	1,636 (1,619, 1,655)

Year	ϕ_{gt}, P_{bgt}	ϕ_{gt}, P_{gt}	ϕ_{t}, P_{bt}	ϕ_{t}, P_{t}
2014	498 (368, 675)	329 (190, 517)	870 (550, 1,069)	263 (184, 369)
2015	1,252 (1,011, 1,664)	864 (662, 1,147)	1,522 (1,246, 1,847)	788 (609, 985)
2016	1,907 (1,813, 1,964)	1,600 (1,539, 1,683)	1,862 (1,720, 2,042)	1,614 (1,550, 1,705)
2017	1,893 (1,847, 1,939)	1,820 (1,802, 1,840)	1,941 (1,864, 2,044)	1,820 (1,801, 1,839)
2018	1,953 (1,923, 1,982)	1,910 (1,895, 1,928)	1,970 (1,930, 2,024)	1,910 (1,895, 1,927)
2019	1,958 (1,934, 1,989)	1,927 (1,912, 1,944)	1,977 (1,941, 2,018)	1,928 (1,913, 1,946)
2020	1,855 (1,836, 1,877)	1,838 (1,824, 1,855)	1,862 (1,839, 1,887)	1,838 (1,824, 1,855)
2021	1,639 (1,624, 1,661)	1,635 (1,619, 1,655)	1,641 (1,622, 1,663)	1,634 (1,619, 1,651)

Table C2. Homeless population size estimates for Island Health from 2014 to 2021: Medians and 95% credible intervals for annual population size

n_{obs}	Year	Estimates (CI)
90	2014	571 (276, 2,162)
187	2015	1,684 (974, 2,096)
1,383	2016	1,890 (1,675, 2,076)
1,572	2017	1,924 (1,845, 2,058)
1,614	2018	1,964 (1,920, 2,034)
1,622	2019	1,969 (1,933, 2,022)
1,587	2020	1,861 (1,837, 1,889)
1,524	2021	1,642 (1,624, 1,665)

Note. Here, n_{obs} is the number of observed patients and the model used was ϕ_{agt}, p_{agbt} .

Table C3. Homeless population size estimates for Greater Victoria from 2014 to 2021: Medians and 95% credible intervals for annual population size

n_{obs}	Year	Population size estimate
47	2014	225 (103, 467)
75	2015	686 (323, 956)
630	2016	853 (733, 959)
699	2017	902 (847, 967)
741	2018	930 (900, 978)
757	2019	948 (923, 976)
760	2020	895 (880, 911)
718	2021	787 (773, 806)

Note. Here, n_{obs} is the number of observed patients and the model used was ϕ_{agt}, p_{agbt} .

Table C4. Homeless population size estimates for Island Health, using only ER interactions, from 2017 to 2021: Medians and 95% credible intervals for annual population size

n_{obs}	Year	Estimates (CI)
1,281	2017	1,940 (1,721, 2,115)
1,343	2018	1,969 (1,828, 2,101)
1,328	2019	1,895 (1,810, 1,963)
1,277	2020	1,726 (1,677, 1,776)
1,193	2021	1,426 (1,388, 1,473)

Note. Here, n_{obs} is the number of observed patients and the model used was ϕ_{agt}, p_{agbt} .

References

- Abadi F., Gimenez O., Arlettaz R., & Schaub M. (2010). An assessment of integrated population models: Bias, accuracy, and violation of the assumption of independence. *Ecology*, *91*(1), 7–14. <https://doi.org/10.1890/08-2235.1>
- Altieri L., Farcomeni A., & Fegatelli D. A. (2023). Continuous time-interaction processes for population size estimation, with an application to drug dealing in Italy. *Biometrics*, *79*(2), 1254–1267. <https://doi.org/10.1111/biom.13662>
- Aubry T., Farrell S., Hwang S. W., & Calhoun M. (2013). Identifying the patterns of emergency shelter stays of single individuals in Canadian cities of different sizes. *Housing Studies*, *28*(6), 910–927. <https://doi.org/10.1080/02673037.2013.773585>

