

# Preventing & Reducing Homelessness: An Integrated Data Project

**Technical Documentation** 





# **Table of Contents**

Executive Summary4	ŀ
Introduction6	5
Background7	7
Project Outline	7
Point-in-Time Counts	,
Cross-Agency Team	3
Data Innovation Program	3
Methods9	)
Overview	)
Phase 1 Datasets	)
Analytical Definition10	)
Results 12	)
Understanding the Data Sources12	2
Homelessness Cohort13	\$
Demographics15	5
Conclusion 20	)

# Purpose

This document summarizes the analytical methods used in Phase one of the <u>Preventing and</u> <u>Reducing Homelessness Integrated Data</u> project. Phase one focused on implementing a crossagency analytic definition of homelessness based on administrative data.

# **Project Overview**

The Ministry of Attorney General and Minister Responsible for Housing, in partnership with BC Housing (BCH), Ministry of Citizens' Services and the Ministry of Social Development & Poverty Reduction (SDPR), have undertaken an integrated data project with a focus on homelessness in British Columbia.

Phase one of this project has estimated the population who experienced homelessness in B.C. in 2019 by building an analytical definition of homelessness defined by SDPR's *Employment and Assistance Program* usage, and integrating BCH's definition of homelessness, defined by *Emergency Shelter Program* usage. Several attribute variables were defined that address the demographic and geographic characteristics of this population. This is a novel approach and represents a new method, using administrative data, of estimating homelessness for British Columbia.

The project is enabled by the <u>Data Innovation Program</u>, a data integration and analytics program for the B.C. government. While every B.C. ministry collects and manages its own data, the Data Innovation Program can securely link and de-identify data from multiple ministries, giving government analysts a broader understanding of complex issues. The Data Innovation Program is based on world-leading best practices for managing safe access to confidential or sensitive data, following the <u>Five Safes model</u> to protect data and reduce the risk of sensitive data being accessed or used inappropriately.

The data science methods and implementation are being conducted by the Data Science Partnership program—a service provided through the Digital Platform and Data Division of the Office of the Chief Information Officer, Ministry of Citizens' Services—which aims to accelerate adoption of data-driven digital government and support evidence-based decision making by providing cross-government data science support.

# Data

Data available through the Data Innovation Program is always de-identified to protect privacy. This means identifiers such as names, driver's licence numbers and personal health numbers are removed and replaced by a project-specific studyid used to link individuals across data sets in the secure analytics environment. Publicly available, open-licensed population estimates from Statistics Canada were used to compare demographic characteristics of the estimated 2019 homeless population to the B.C. provincial population.

Data sets used in Phase one of this project were:

### B.C. Employment and Assistance (SDPR, Data Innovation Program)

- 2019 Involvement data
- · 2019 No Fixed Address data

### **Emergency Shelter Program (BCH, Data Innovation Program)**

- Shelter Stays data
- Clients data

### **Client Registry (Data Innovation Program)**

• B.C. Central Demographics data

### Statistics Canada Table 17-10-0005-01 (open-licensed data)

• 2019 Population Estimates on July 1st, by age and sex

### Statistics Canada Table 17-10-0139-01 (open-licensed data)

Population estimates, July 1, by census division, 2016 boundaries

# **Data Integration & Provisioning**

Data sets are linked and provisioned into a secure analytics environment by <u>Population Data</u> <u>B.C.</u>, an provincially recognized academic organization that has facilitated population-based research for over 20 years.

Population Data B.C. uses a *population directory* as the primary key for data linkage. It includes all the individuals about whom Population Data B.C. has information. The directory has been built with Medical Services Plan (MSP) registration & premium billing data going back to 1985. Since MSP covers most of the B.C. population and issues a unique lifetime Personal Health Number, this directory provides the largest representative coverage of residents in B.C. to support high-quality linkage across data sets.

The data linkage rate, defined here as the percentage of individuals linked to the population directory for a given data set, is dependent on the number and kind of identifiers contained in the data set. For data used to estimate the homeless population, SDPR data contains Personal Health Numbers, while BC Housing's data contains only provided name, date of birth and gender. As a result, the linkage rate of the BC Housing data was lower than the SDPR data.

While the administrative data sets available through the Data Innovation Program are comprehensive, administrative data is still a sample of the actual population and all derived estimates have some associated uncertainty. Data set linkage rates offer quantification of some of this uncertainty across linkages in *all* data sets in the Data Innovation Program. While "unlinked" studyids within a data set represent studyids not linkable to any data in the Data Innovation Program, this is distinct from individuals present in some (but not all) data sets. For example, some individuals were present in both the SDPR and BCH data while others were present in only one or the other. For the definition of homelessness, this enables inclusions of individuals using either service. However, when linking to additional data sets, for example the B.C. Central Demographic file, it means that not all individuals in the SDPR and BCH data are linked or found in the B.C. Central Demographic data. Further complicating linkages, while the *population directory* developed by Population Data B.C. is continuously updated, the data provisioned to a given project may have been linked based on older versions of the directory, resulting in "linked" individuals who are not present in any other data sets. Fortunately, this number remains low and does not presently affect the population-level analysis described here.

# **Data Preparation**

A number of BC Housing clients matched multiple studyids in other administrative data sets. Population Data B.C. would typically collapse these studyids into a single id and restate all domain data sets accordingly before provisioning the data into the secure analytics environment. In this case, to expedite data release, that work was passed along to the project team by providing a PopData Extra ID table that defines the BC Housing client's relationship to other data sets. For the 2019 period, the step-by-step process to subsequently collapse those studyids is as follows:

- Find all unique combinations of studyid and clientid in the PopData Extra ID file
- Find all unique combinations of studyid and clientid in the BC Housing Clients file
- Find all rows where clientid is equal in both PopData Extra ID and BC Housing Clients files

This process, recommended by Population Data B.C., results in a crosswalk table that allows for identifying duplicated ids outside of the BC Housing data and restating data sets by associating duplicated ids with the canonical studyid from the BC Housing data.

### Statistical Disclosure

This project meets or exceeds the statistical disclosure guidelines set out by the Data Innovation Program. To maintain future flexibility, most data outputs have been aggregated well above the cell size thresholds to reduce the risk of future re-identification if more finegrained results are exported.

### **Background Materials**

### **Employment and Assistance Over Time**

The SDPR B.C. Employment and Assistance No Fixed Address (NFA) data was linked with the SDPR B.C. Employment and Assistance Involvement data using clientid as a primary key to determine, of those individuals who collect income assistance, which have NFA status. The data was summarized by counting all individuals making income assistance claims and those only making the claim with NFA status over the year and month of income collection. This was done for the 2010 to 2019 period of Involvement data.

### **Time in Shelter**

This output summarized the number of individuals who have stayed for any amount of time in a shelter in 2019. The time spent in shelter represents a range of length of stays. Because this data is not linked to any other dataset and is a BCH specific data set, the unlinked ids are included here.

# Study Population: Estimated 2019 Homelessness Population

The estimated population experiencing homelessness was defined using service use patterns from both the BCH Shelter Stays data and SDPR B.C. Employment and Assistance data. The 2019 definition period was chosen because it represents a common temporal period in both data sets and because a calendar year facilitates year-over-year comparisons as new data becomes available. The study population was defined at the monthly time scale, primarily because B.C. Employment and Assistance is delivered on a monthly basis. An individual was defined as experiencing homelessness in a given month:

a) if they had three consecutive months of income assistance (including the current month) with no fixed address or;

b) they spent any time in a shelter in the given month

Using this criteria, the population estimate was generated for each month of 2019. For the months of January and February, this required using income assistance data from 2018 to determine if an individual met the definition. For an annual estimate, unique studyids from the monthly estimate were compiled to form the 2019 population. These two temporal resolutions defined whether a person experienced homelessness during a given month in 2019 or whether a person experienced homelessness at any point during 2019.

# **Study Population Attributes**

Several demographic and geographic population attributes were added to the above population on both monthly and annual time scales. For attributes where this makes sense, the annual and monthly attributes are differentiated below.

### **Data Source**

The means of entry into the population was recorded to determine whether an individual entered the study population via service usage definitions from the SDPR data, BCH data or both. At the monthly resolution, this measure directly results from whether a person entered into the population by their usage of income assistance, shelters or both. At the annual resolution, this measure expresses whether a person ever used either or both services sufficiently to meet the homeless definition.

### **Homeless Category**

A chronic homelessness sub-population, distinct from the non-chronic homelessness population was also defined. This category was defined using an individual's past 12 months of service usage for a given month. For some months, this includes reaching into the previous year's shelter and income assistance data. We implemented the following criteria to define homeless categories:

- Calculate the cumulative number of nights spent in a shelter over the past 12 months
- Calculate the cumulative number of shelter visits separated by 30 days (unique visits) over the past 12 months

- Calculate the longest period of consecutive monthly income assistance for those individuals that reported no fixed address over the past 12 months
- Apply the following definition for a given month:
  - Non-Chronic Homelessness: 3 to 5 months of consecutive income assistance reporting no fixed address OR 180 or fewer days in a shelter OR 1 or 2 unique visits to a shelter (separated by 30 days)
  - Chronic Homelessness: 6 or more months of reporting no fixed address OR more than 180 days in a shelter OR 3 or more unique visits to a shelter (separated by 30 days)
- In instances where the above criteria resulted in differing non-chronic homelessness and chronic homelessness outcomes between the services, a studyid was associated with chronic homelessness

For the annual estimate, if a person is assigned chronic status during any month (based on their previous 12 months of service usage) in 2019, they were assigned chronic status for all of 2019. A small portion of the BCH shelter data contained individuals with overlapping time intervals at different locations. For example, some people were registered in different shelters at the same time. Those overlapping intervals were merged to create one continuous interval. This was possible when the overlapping intervals were in the same city. For instances where an individual was present in two census subdivisions, the intervals were not merged. This is a very small number of shelter visits and therefore represents an acceptable loss of data accuracy.

### Demography

This project had three sources of age and gender data for the study population: BCH client data, SDPR Income Assistance data and the B.C. Central Demographic data. This section outlines the hierarchy of those data sources and the process to determine which to use.

#### **Date of Birth**

The primary data source for date of birth is the B.C. Central Demographic file. When no date of birth was available in that file, date of birth from the SDPR Income Assistance data was used. These two data sources are the most reliable for date of birth because they require some level of verification for a service to be provided. If an individual still had not been assigned a date of birth, the BCH client data was used, which is a self-reported measure. Age was calculated as the age (as an integer) as of December 31, 2019, which is the end of the definition period. To meet conservative statistical disclosure goals, ages were grouped into three groups: 24 & under, 25 to 55, and 55 & over. To compare age distributions of the homelessness population relative to the general population, population estimates by age from Statistics Canada were used.

#### Gender

We have integrated gender identity data<sup>1</sup> from various sources and specified a hierarchy for data quality. Each of the Central Demographic files , SDPR Income Assistance data and BCH client data have gender indicators of varying quality.

The hierarchy of data quality used for gender identity data for this integrated project is as follows:

- 1. BCH client data is collected using open-ended responses and therefore is the priority gender indicator
- 2. The SDPR income assistance data, which collects binary genders, is the most recent data and is the next priority gender indicator
- 3. B.C. Central Demographic data, which collects binary genders, is the most dated data source and is therefore the last choice for gender identity. Note that while B.C. Central Demographic file is the least-used data source for gender identity in this circumstance, it is still a reliable data source

The number of individuals identifying in the shelter data as non-binary was very small. This low number presented a possible statistical disclosure risk for that population of individuals. Therefore, individuals identifying as non-binary in the shelter data were recoded to their binary gender data present in the SDPR or MSP data.

#### Geography

To estimate the location of individuals experiencing homelessness, the location where a person was accessing services was used as a proxy. The hierarchy for this geographic proxy was:

- 1. Location of a shelter where a person visited
- 2. Location of the Service B.C./SDPR issuing office where a person accessed income assistance

The BCH shelter data provides location information for individual shelter visits while the SDPR income assistance data only provides monthly location information. For both of these data sets, the location was provided at either the municipal or the <u>census subdivision</u> level. Adopting a conservative approach to statistical disclosure, the export aggregated these data to the <u>census division</u> level. Further, the census divisions of 'Mt. Waddington' and 'Central Coast' were merged into one as were 'Northern Rockies' and 'Peace River.' The resulting census divisions and merged census divisions provided geographies to evaluate where services were used and by proxy where individuals were experiencing homelessness. The population was further aggregated to <u>economic region</u> to provide an additional geography. To calculate rates of homelessness relative to the general population, population estimates of census division from Statistics Canada were used.

<sup>1</sup> As datasets are revised across government to recognize non-binary definitions of gender, a more appropriate portrayal of gender diversity can be brought into this project work.

#### **Movers and Non-Movers**

To resolve individuals that moved only within and between census division (or economic region) a geographic flag was added to each studyid indicating either the geography that a person spent all of 2019 in or a 'multiple cd' (or 'multiple er') flag if they moved between or among geographies. This flag was applied at both the annual and monthly time scales. Summaries of the location and rates by location of the homelessness population focused only on the 'Non-Mover' category.

# Appendix

This appendix is provided as specific documentation of the data analysis conducted on provisioned Data Innovation Program data in the secure analytics environment.

### Software

This analysis is implemented in the R programming language (R Core Team 2021). The code used to generate this analysis was reviewed by three data scientists. Key tools used to complete this work include the Apache Arrow project (Richardson et al. 2021), the tidyverse (Wickham et al. 2019), cansim (von Bergmann and Shkolnik 2021), dipr (Albers and Hazlitt 2020) and the R package targets (Landau 2021) for project organization. All code is stored under the <u>git version control</u> system and shared inside the secure environment in these GitLab repos:

- Parquet and restating: <u>https://projectsc.popdata.bc.ca/shares/data-to-parquet</u>
- Application of definition and creation of output group: <u>https://projectsc.popdata.bc.ca/shares/homelessness.cohort</u>

### Raw Data

All data was converted from compressed fixed width files into parquet files for ease of analysis. Significant testing against Population Data B.C. data provisioning metrics occurred to ensure that conversions were done accurately.

### Restating

Popdata provided a table to resolve duplicate ids from data outside the BC Housing data by using the clientid field as a data bridge. Popdata refers to this process as "restating."

### **ExtraID data**

- filename: "bchousing\_hifis2017-2019.hifis\_extra\_clntid\_popid\_xlk.A.dat.gz"
- columns: studyid, clientid

### **Social Development and Poverty Reduction Data**

2019 Involvement data

- filename: "idosdpr2018-2019.bceainvolvement.A.dat.gz"
- columns used: ym, fileid, studyid, deprltncd, birthdt\_yymm, gender

#### 2019 NFA data

- filename: "idosdpr2018-2019.bceanfa.B.dat.gz"
- columns used: ym, fileid, nfa, csdname
- subsetting: only rows where nfa == 1

### **BC Housing Data**

2019 Shelter Stays Data

- filename: "bchousing\_hifis2017-2019.hifis\_clnts\_shlt\_stays.dat.gz"
- columns used: shelter\_stay\_start\_date, shelter\_stay\_end\_date, shelter\_stay\_ city, studyid, clientid

2019 Clients data

- filename: "bchousing\_hifis2017-2019.hifis\_clients.dat.gz"
- columns used: dob\_yyyymm, gender, studyid, clientid

#### Extra ID data

- filename: "bchousing\_hifis2017-2019.hifis\_extra\_clntid\_popid\_xlk.A.dat.gz"
- columns used: clientid, studyid

### **B.C. Central Demographic data**

- filename: "demographics1986-2019.B.dat.gz"
- columns used: sex, dobyyyy, dobmm, studyid

# Data Processing

### **Restating Method**

The extra id data represent linkages to multiple studyids in other data sets that have not been collapsed to a single id. Popdata refers to this term as "restating" and we have adopted the same terminology here. Restating typically happens with popdata. In this case, however, that work has been passed along to the research team. The step by step process is as follows:

- Find all unique combinations of studyid and clientid in the bchousing\_hifis2017-2019.hifis\_ extra\_clntid\_popid\_xlk.A.dat.gz file
- Find all unique combinations of studyid and clientid in the bchousing\_hifis2017-2019.hifis\_ clients.A.dat.gz file
- Find all rows where clientid is equal in both extra\_id and clients data
- Remove rows where extra\_id studyids and clients ids are already the same

This process results in a crosswalk table that allows us to identify duplicated ids outside of the BC Housing data and associated them with the canonical studyid from the BC Housing data. This process is informed by popdata and is their recommend approach. The restating process occurred with any data outside of the BC Housing data.

### **Geographic Data**

BCH client data was provisioned with the geographic variable shelter\_stay\_city. SDPR data was provisioned with the geographic variable csd\_name. These variables can both be considered census subdivision (CSD) because of the commonality of municipality and CSD. However, because of slight differences in the names of geographic locations between SDPR, BCH and the BCStats population estimates, a crosswalk table was created. This approach a) ensures that population estimates were assigned to the correct CSD and b) resolves slight naming differences between the SDPR and BCH geographic data sets.

Manipulations: A movement flag was derived for a given geography. This binary flag specified when a person (at monthly or annual time scales) moved outside of a given geographic area. For example, if a person moved between or among census divisions during a year, they would receive a flag of "multiple cd."

### SDPR data

Definition

• Homeless: At least 3 months of consecutive income assistance with the nfa == 1 for a given month during a 12 month period.

Attributes (hl\_category)

• Non-chronic homelessness: 3 to 5 months of consecutive income assistance with the nfa == 1 during the past 12 months for any given month in 2019.

• Chronic homelessness: 6 or months of consecutive income assistance with the nfa == 1 during the past 12 months for any given month in 2019.

Manipulations: The nfa data was linked with the involvement data using clientid as a primary key to add studyid to the nfa data. The merged nfa-involvement data was subset for only rows where nfa == 1. Each row represents a monthly income assistance payment. We calculated the number of months of consecutive income assistance with an nfa == 1 for 2019 for a given month. This allowed for monthly assignment of the hl\_category attribute. To assign an annual attribute of hl\_category, the occurrence of any chronic homelessness status during any month resulted in an annual status of chronic homelessness.

### **BC Housing data**

#### Definition

• All individuals who spent any time in a shelter during 2019

#### Attributes (hl\_category)

- Non-chronic homelessness: all individuals who spent 180 or fewer days in a shelter or who had 1 or 2 shelter visits separated by 30 days during the past 12 months for any given month in 2019.
- Chronic homelessness: all individuals who spent more than 180 days in a shelter or who had 3 or more shelter visits separated by 30 days during the past 12 months for any given month in 2019.

Manipulations: The shelter attribute was calculated on a monthly basis. For a given month, shelter visits were subset into that month and the preceding 11 months to assign the hl\_ category attribute. Visits separated by 30 days and the cumulative number of days spent in a shelter were calculated and the BC Housing definition was applied to classify individuals as experiencing non-chronic homelessness or chronic homelessness. In instances where cumulative number of days and number of shelter visits resulted in different classifications, chronic homelessness was taken as the classification. To assign an annual attribute of hl\_ category, the occurrence of any chronic homelessness status during any month resulted in an annual status of chronic homelessness. A small portion of the BCH shelter data contained individuals with overlapping time intervals at different locations. That is, some people are registered in different shelters at the same time. Those overlapping intervals were merged to create one continuous interval. This was possible when the overlapping intervals were in the same city. For instances where an individual was present in two cities, the intervals were not merged. This is a very small number of shelter visits and therefore represents an acceptable loss of data accuracy.

### Demography data

#### Age data

Age data was derived from three ranked sources. Age for an individual was chosen in the following order:

- B.C. Central Demographic file: most authoritative source for accurate birth data.
- SDPR data: largely also derived from registry data but also providing more values due to slightly more complete data
- BCH data: least accurate age data as it is entirely self reported

#### **Gender data**

We established a hierarchy of our data source:

- BCH data collects gender using open ended response and is therefore the priority gender indicator.
- SDPR income assistance data, which collects binary gender, and is the most recent data (2019) is therefore the next priority.
- B.C. Central Demographic file, which also collects binary gender, is the least current data source and therefore, while still very reliable, is the lowest priority.

Manipulations: Studyids with two birthdays were removed from the demographic data. To address concerns of statistical disclosure, the category of "Non-Binary" was ultimately removed from the dataset. All instances of this category of gender were in the BCH data and therefore replaced by either SDPR and demographics data. Age was calculated as the age as of December 31, 2019 which is the end of the accrual window and then rounded down to the nearest integer.

### Merging BC Housing and SDPR data

- Data was merged in the primary studyid key keeping track of whether an individual was present in the BC Housing data, SDPR data or both.
- In instances where the BC Housing and SDPR resulted in different homeless classifications, chronically homeless was chosen.
- To define the annual data\_source variable, any service usage (from either SDPR or BCH data) resulted in a data\_source flag for that individual.

### GitLab repo:

- Parquet and restating: <u>https://projectsc.popdata.bc.ca/shares/hl-data-to-parquet</u>
- Application of definition and creation of output group: https://projectsc.popdata.bc.ca/shares/homelessness.cohort

# References

Albers, Sam, and Stephanie Hazlitt. 2020. Dipr: Provide Functions to Efficiently Import SRE Data.

Landau, William Michael. 2021. "The Targets r Package: A Dynamic Make-Like Function-Oriented Pipeline Toolkit for Reproducibility and High-Performance Computing." *Journal of Open Source Software* 6 (57): 2959. <u>https://doi.org/10.21105/joss.02959</u>.

R Core Team. 2021. R: *A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <u>https://www.R-project.org/</u>.

Richardson, Neal, Ian Cook, Jonathan Keane, Romain François, Jeroen Ooms, and Apache Arrow. 2021. Arrow: Integration to 'Apache' 'Arrow'. <u>https://CRAN.R-project.org/package=arrow</u>.

von Bergmann, Jens, and Dmitry Shkolnik. 2021. *Cansim: Accessing Statistics Canada Data Table and Vectors*. <u>https://CRAN.R-project.org/package=cansim</u>.

Wickham, Hadley, Mara Averick, Jennifer Bryan, Winston Chang, Lucy D'Agostino McGowan, Romain François, Garrett Grolemund, et al. 2019. "Welcome to the tidyverse." *Journal of Open Source Software 4* (43): 1686. <u>https://doi.org/10.21105/joss.01686</u>