

...
*
...

Julia Lee

December 3, 2024

...

1 Introduction

2 Data

To simulate, test, download, and clean the Neighbourhoods Profile data, the statistical programming language R was used (R Core Team 2023). Specific libraries that assisted the analysis include `tidyverse` (Wickham et al. 2019), `opendatatoronto` (Gelfand 2022), `tinytex` (Xie 2019), `ggplot2` (Wickham 2016), `knitr` (Xie 2015), `testthat` (Wickham 2011), `here` (Müller and Bryan 2020), `arrow` (Richardson et al. 2024), `modelsummary` (Arel-Bundock 2022), and `sf` (Pebesma 2018).

2.1 Neighbourhood Profiles Data

Neighbourhood Profiles data for the city of Toronto is provided by the Social Development, Finance & Administration from Toronto's Open Data portal (Social Development, Finance & Administration 2024a). The dataset obtained from the city of Toronto's Open Data portal (City of Toronto 2024), is a collection of data from the 2021 Canadian Census from Statistics Canada (Statistics Canada 2024). Records of socio-economic data for 158 geographic regions (i.e. social planning neighbourhoods) in Toronto can be found within this dataset. From age to ethnocultural diversity, the Census data highlights both the socio-economic and demographic characteristics of Toronto residents across its individual neighbourhoods.

The Neighbourhood Profiles data is an extensive dataset with each row representing a demographic or socio-economic characteristic and each column reflecting a Toronto neighbourhood.

*Code and data are available at: https://github.com/jjlee-lee/Toronto_Neighbourhood_Income.git

For each neighbourhood, there is information on its name, identification number, and status. Status refers to a neighbourhood’s Toronto Strong Neighbourhoods Strategy (TSNS) designation, which indicates whether a neighbourhood is an emerging area, an area that needs improvement, or neither. Additionally, data for each neighbourhood’s age population, income, education, and more from the 2021 Census is provided. Table 1 below offers a small preview of this dataset.

Table 1: Toronto Neighbourhood Profiles Data

Neighbourhood Name	West Humber-Clairville
Neighbourhood Number	1
TSNS 2020 Designation	Not an NIA or Emerging Neighbourhood
Total - Age groups of the population - 25% sample data	33300
0 to 14 years	4295
0 to 4 years	1460
5 to 9 years	1345

Table 1 displays the name, number, and TSNS designation along with data for the number of individuals within different age groups for a neighbourhood (West Humber-Clairville) in Toronto. By looking at Table 1, West Humber-Clairville is a neighbourhood that is not an emerging area and one that does not need improvement. It is also a region with a fairly large number of children ages 0 to 14 years old. The rest of the dataset follows the same format, displaying different characteristics (e.g. older age groups, education attainment, income levels, etc.) of all 158 neighbourhoods in the city of Toronto.

2.2 Analysis Data

For the present analysis, the variables of interest are the different types of census family sizes along with the average after-tax income for each of Toronto’s neighbourhoods. The Neighbourhood Profiles data includes four different census family sizes: (1) 2-person families, (2) 3-person families, (3) 4-person families, and (4) five or more-person families. In simple terms, Statistics Canada defines a census family as one where all family members (related by blood marriage, common-law union, adoption, or a foster relationship) live together in the same dwelling (Statistics Canada 2023). Census families can also be referred to as economic families (Statistics Canada 2021). The data further provides the average after-tax income (\$) recorded in 2020 for each neighbourhood.

Since the objective is to find out what kinds of families are driving average neighbourhood income in Toronto, the data used throughout this analysis reflects the average income level for all 158 neighbourhoods alongside their counts of different census family sizes. With minimal

data wrangling, the analysis data is simply an extraction of the larger dataset with a focus on the different family sizes and average after-tax income. Table 2 below illustrates this analysis data for three neighbourhoods and summary statistics for this data can be found in the Appendix (Section A.1).

Table 2: Toronto Neighbourhood Profiles Analysis Data

Name		2	3	4	5 or more	Average
	Number persons	persons	persons	persons	persons	Income
West Humber-Clairville	1	3635	2265	2025	805	101300
Mount Olive-Silverstone-Jamestown	2	2855	2145	1765	1290	85300
Thistletown-Beaumont Heights	3	1095	665	555	310	98100

In Table 2, the “Name” column represents each neighbourhood’s name, and the “Number” column reflects each neighbourhood’s identification number. The variables, “2 persons”, “3 persons”, “4 persons”, and “5 or more persons” represent the number of census families in each neighbourhood with those particular family sizes. Lastly, the “Average Income” variable reflects the average after-tax income for all neighbourhoods in 2020 (recorded in dollars). Within the analysis data file provided by this analysis, there are additional variables that represent the total number of census families by family size, the average census family size, and the average number of children in census families for each neighbourhood. These variables are considered to obtain additional context about how families are made up and distributed across Toronto’s neighbourhoods. It is important to note that while this additional information provided a better understanding of the variables of interest, they are not included in the model of the analysis. A detailed account of the model can be found in (Section 3).

2.3 Map Data

To further understand the type of families that drive average neighbourhood income, this analysis also uses Neighbourhoods data – a shapefile that contains geographic information of Toronto neighbourhood boundaries – to map the results of its model. The boundaries of each neighbourhood are defined using census tract information provided by Statistics Canada, and the shapefile itself is published by the Social Development, Finance & Administration from Toronto’s Open Data portal (Social Development, Finance & Administration 2024b). The map within this analysis is created using ArcGIS Pro software (Esri 2024). By joining the analysis data shown in Table 2 with the shapefile based on neighbourhood names, a map that highlights the distribution of family sizes that drive average neighbourhood income can be created. The information contained in the shapefile is shown below in Table 3.

Table 3: Toronto Neighbourhood Location & Boundaries Data

Object ID	Neighbourhood Number	Neighbourhood Name	TSNS Designation	Geometry
1	174	South Eglinton-Davisville	Not an NIA or Emerging Neighbourhood	POLYGON ((-79.38635 43.6978...
2	173	North Toronto	Not an NIA or Emerging Neighbourhood	POLYGON ((-79.39744 43.7069...
3	172	Dovercourt Village	Not an NIA or Emerging Neighbourhood	POLYGON ((-79.43411 43.6601...

The Neighbourhoods shapefile data contains neighbourhood names, numbers, and TSNS designations similar to the Neighbourhood Profiles data (Table 1). It also includes geographic information (i.e. the geometry/coordinates of a neighbourhood) that defines the boundaries of each neighbourhood as shown in Table 3.

2.4 A Note on Measurement

The Neighbourhood Profiles data is a reflection of the 2021 Census for neighbourhoods in Toronto, meaning that it measures and represents the demographic and socio-economic characteristics of families across Toronto at this particular time period. This information is collected through the use of questionnaires – a short and long-form census. The short-form census is sent out to all households, and it asks questions about simple characteristics like age, gender, and household size. With this information, the short-form census attempts to enumerate all individuals within a geographic region and may impute any missing data. The long-form census is sent out to only 25% of households and is a deeper source of data that outlines information such as housing, education, and ethnicity.

So, various aspects of families’ livelihoods are measured through the Census and aggregated to create an overall picture of particular geographic regions – in this case, Toronto neighbourhoods. It is important to note that through this measurement process, there is room for error as well as a need to be aware of this potential for error. Missing data values can be imputed and the aggregation of individual household characteristics or experiences can lead to incorrect assumptions about certain regions and their populations. Thus, this analysis presents its findings with this in mind.

3 Model

To understand the types of families that drive average neighbourhood income, this analysis constructs a multiple linear regression model using R (R Core Team 2023) that considers average neighbourhood income as the response variable and the different census family sizes as predictors. As mentioned in Section 2, the average neighbourhood income variable reflects the average after-tax income of all census families for each neighbourhood in 2020 (recorded in dollars). Different family sizes (i.e. 2 to 5 or more-person families) refer the number of family members that live in the same dwelling within each individual neighbourhood. With this, the final model within this analysis is structured as follows:

$$\hat{y} = b_0 + b_1 \cdot x_1 + b_2 \cdot x_2 + b_3 \cdot x_3 + \epsilon$$

where

- \hat{y} represents the average neighbourhood after-tax income,
- b_0 represents the intercept of the regression model,
- b_1 represents the effect of two-person census families,
- x_1 represents the number of two-person census families,
- b_2 represents the effect of three-person census families,
- x_2 represents the number of three-person census families,
- b_3 represents the effect of four-person census families,
- x_3 represents the number of four-person census families,
- ϵ captures the error within the regression model

It is important to highlight that as a result of model selection and individual predictor t-tests, families with 5 or more members are not considered in the final model. The impact of 2, 3, and 4-person families on average neighbourhood income across Toronto neighbourhoods are examined through this model.

3.1 Model justification

As the objective of the analysis is to understand what might be driving the average neighbourhood income in terms of family size, a multiple linear regression will allow for an examination of the relationship between different census family sizes and average income. By striving to preserve the interpretability of the model, this analysis can further determine which family sizes are meaningful influencers of a neighbourhood's average after-tax income.

The underlying assumption of linearity that multiple linear regression models hold can be a potential limitation as these models will not be able accurately explain the influence of predictors on a response variable if these variables have a non-linear relationship with each other. The use of a multiple linear regression may not be the best choice in these cases. However, with average neighbourhood income and increasing census family sizes, the model's diagnostic plots provide evidence for linearity, meaning that the use of a multiple linear regression model can be suitable for this analysis. Further justification and validation for the model can be found in the Appendix (Section [A.2](#)).

4 Results

5 Discussion

5.1 Limitations & Future Directions

A Appendix

A.1 Analysis Data Summary Statistics

Table 4: Toronto Neighbourhood Profiles Data Summary Statistics

2-person Families	3-person Families	4-person Families	5 or more-person Families	Average Income
Min. : 815	Min. : 260.0	Min. : 220.0	Min. : 40.0	Min. : 76800
1st Qu.:1596	1st Qu.: 716.2	1st Qu.: 590.0	1st Qu.: 185.0	1st Qu.: 93450
Median :2172	Median :1030.0	Median : 867.5	Median : 310.0	Median :108150
Mean :2288	Mean :1099.4	Mean : 900.2	Mean : 351.4	Mean :121582
3rd Qu.:2879	3rd Qu.:1393.8	3rd Qu.:1180.0	3rd Qu.: 460.0	3rd Qu.:129500
Max. :5435	Max. :2265.0	Max. :2025.0	Max. :1290.0	Max. :351600

Table 4 presents the summary statistics for each census family size and the average after-tax income across all 158 neighbourhoods in Toronto.

A.2 Model Validation

The construction of a model for this analysis began by fitting an initial model with average neighbourhood income as the response variable and 2, 3, 4, and 5 or more-person families as predictors. As seen by Table 5 and the individual t-tests, 5 or more-person families do not appear to significantly impact a neighbourhood's average income with the largest p-value among all predictors ($p > 0.001$).

So, another model was created without this variable (i.e. a reduced model), and all predictors seem to have significant impacts ($p < 0.001$) on neighbourhood income as shown in Table 6.

Model selection was further conducted by comparing the BIC scores of both models (shown in Table 7). As the new model has a lower BIC score, this is the one that is used within the analysis and outlined in Section 3.

Table 5: Summary Table of the Initial Model

	(1)
(Intercept)	125 502.083
	8644.967 (<0.001)
size_two	26.444
	4.960 (<0.001)
size_three	−178.177
	18.084 (<0.001)
size_four	157.742
	17.431 (<0.001)
size_five_plus	−29.981
	20.206 (0.140)
Num.Obs.	158
R2	0.433
R2 Adj.	0.418
AIC	3759.2
BIC	3777.6
Log.Lik.	−1873.608
RMSE	34 178.04

Table 6: Summary Table of the Reduced Model

	(1)
(Intercept)	125 318.176
	8677.734 (<0.001)
size_two	28.593
	4.763 (<0.001)
size_three	−184.416
	17.657 (<0.001)
size_four	148.403
	16.318 (<0.001)
Num.Obs.	158
R2	0.424
R2 Adj.	0.413
AIC	3759.5
BIC	3774.8
Log.Lik.	−1874.736
RMSE	34 423.06

Table 7: BIC Scores for Initial and Reduced Models

Model	BIC Score
Initial	3777.6
Reduced	3774.8

The purpose of the model is to be able to interpret the coefficient estimates and what they mean in terms of influencing average neighbourhood income. For this reason, transformations (e.g. log transformations on predictor variables) were not applied to the data even with the presence of some evidence against the linear regression model assumptions.

Looking at the diagnostic plots (Figure 1), the lack of observational patterns in the Residuals vs Predictor plots for 2, 3, and 4-person families points to evidence for the linearity assumption of linear regression models.

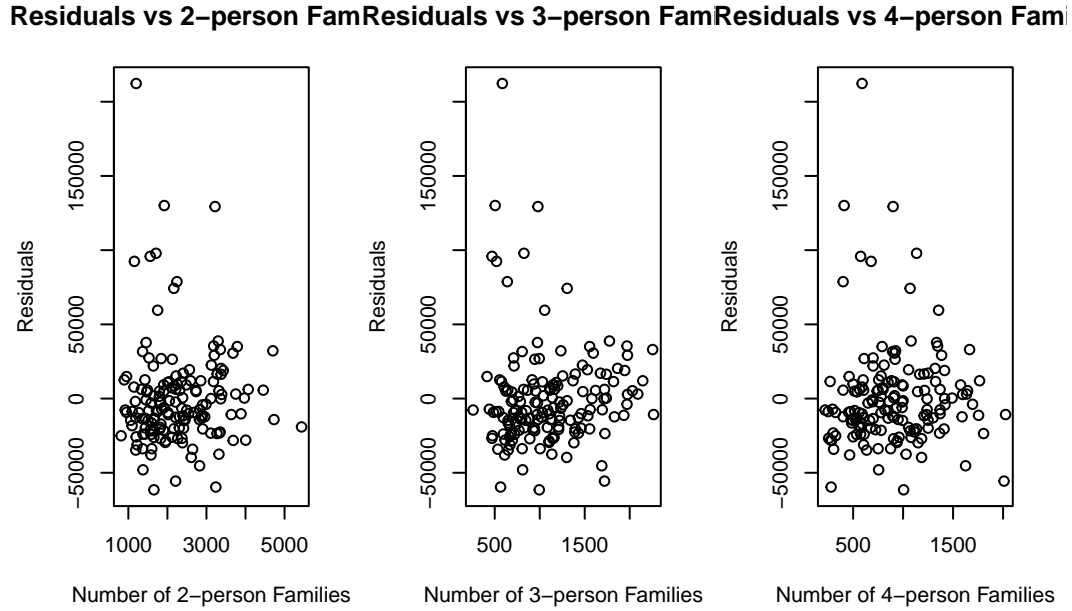


Figure 1: Showing the Residual vs. Predictors

The Residual vs Fitted Values plot (Figure 2) seems to have some clustering of points, indicating that the constant variance and independence assumptions are somewhat satisfied. However, because the overall pattern of the plot is still relatively spread out with no observable patterns, and as the objective of the model is interpretability, transformations were not applied to the variables in the model to better these assumptions.

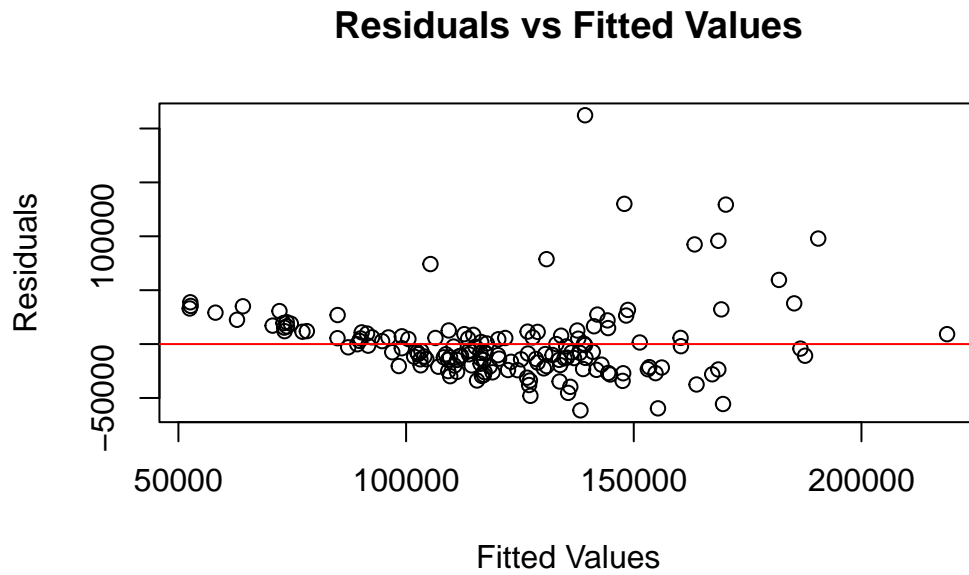


Figure 2: Showing the Residual vs. Fitted Values

Lastly, though the Normal Q-Q plot (Figure 3) shows some deviation from the Q-Q line, the normality assumption is somewhat satisfied with little evidence to suggest an alarming violation of this assumption. Putting all this together, while improvements can be made to the model for better, more accurate predictions, it is sufficient for the purpose of interpretability.

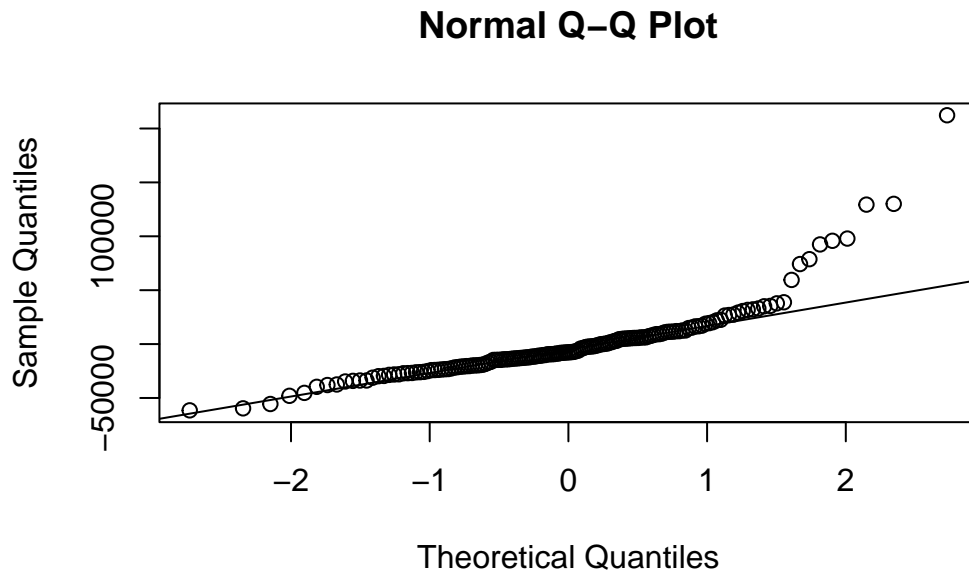


Figure 3: Showing the Normal Q-Q Plot

References

- Arel-Bundock, Vincent. 2022. *modelsummary: Data and Model Summaries in R*. *Journal of Statistical Software*. Vol. 103. <https://doi.org/doi:10.18637/jss.v103.i01>.
- City of Toronto. 2024. “City of Toronto’s Open Data Portal.” <https://open.toronto.ca/>.
- Esri. 2024. “ArcGIS Pro.” <https://www.esri.ca/en-ca/products/gis-mapping-products/arcgis-pro/overview>.
- Gelfand, Sharla. 2022. “opendatatoronto: Access the City of Toronto Open Data Portal.” <https://cran.r-project.org/web/packages/opendatatoronto/index.html>.
- Müller, Kirill, and Jennifer Bryan. 2020. *here: A Simpler Way to Find Your Files*. <https://here.r-lib.org/>.
- Pebesma, Edzer. 2018. “Simple Features for R: Standardized Support for Spatial Vector Data.” *The R Journal* 10 (1): 439–46. <https://doi.org/10.32614/RJ-2018-009>.
- R Core Team. 2023. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.
- Richardson, Neal, Ian Cook, Nic Crane, Dewey Dunnington, Romain François, Jonathan Keane, Dragoş Moldovan-Grünfeld, Jeroen Ooms, Jacob Wujciak-Jens, and Apache Arrow. 2024. *arrow: Integration to ‘Apache’ ‘Arrow’*. <https://github.com/apache/arrow/>.
- Social Development, Finance & Administration. 2024a. “Neighbourhood Profiles.” <https://open.toronto.ca/dataset/neighbourhood-profiles/>.

- . 2024b. “Neighbourhoods.” <https://open.toronto.ca/dataset/neighbourhoods/>.
- Statistics Canada. 2021. “Economic family.” <https://www23.statcan.gc.ca/imdb/p3Var.pl?Function=Unit&Id=33863>.
- . 2023. “Census family.” <https://www23.statcan.gc.ca/imdb/p3Var.pl?Function=Unit&Id=32746>.
- . 2024. “Reference materials, 2021 Census.” <https://www12.statcan.gc.ca/census-recensement/2021/ref/index-eng.cfm>.
- Wickham, Hadley. 2011. *testthat: Get Started with Testing*. *The R Journal*. Vol. 3. https://journal.r-project.org/archive/2011-1/RJournal_2011-1_Wickham.pdf.
- . 2016. “ggplot2: Elegant Graphics for Data Analysis.” <https://ggplot2.tidyverse.org>.
- Wickham, Hadley, Mara Averick, Jennifer Bryan, Winston Chang, Lucy D’Agostino McGowan, Romain François, Garrett Golemund, et al. 2019. “Welcome to the Tidyverse.” *Journal of Open Source Software* 4 (43): 1686. <https://doi.org/10.21105/joss.01686>.
- Xie, Yihui. 2015. *Dynamic Documents with R and Knitr*. 2nd ed. Boca Raton, Florida: Chapman; Hall/CRC. <https://yihui.org/knitr/>.
- . 2019. “TinyTeX: A Lightweight, Cross-Platform, and Easy-to-Maintain LaTeX Distribution Based on TeX Live.” *TUGboat* 40 (1): 30–32. <https://tug.org/TUGboat/Contents/contents40-1.html>.