

# Canada\_leading\_causes\_of\_death\*

Scarlet Ruoxian Wu

March 12, 2024

This study analyzed the main causes of death in Canada between 2000 and 2022 using advanced statistical models to uncover patterns in mortality data. The research used Poisson and negative binomial regressions to identify significant trends and variances in death causes, with a focus on the over-dispersion of such data. The findings reveal a landscape of mortality and highlight the predominance of certain diseases over time and their fluctuating incidence rates. By enhancing our understanding of mortality dynamics, this analysis offers valuable insights for public health policy and prevention strategies that aim to mitigate the most common causes of death and improve overall life expectancy in Canada.

## Table of contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Data</b>	<b>2</b>
<b>3</b>	<b>Result</b>	<b>3</b>

## 1 Introduction

Mortality involves the complex interplay between lifestyle, environment, healthcare access, and genetic predisposition. Every year, mortality statistics shed light on a nation's health and the effectiveness of its healthcare system. The Canadian government has kept records of causes of death from 2000 to 2022. This study aims to examine the leading causes of death during this period, identify patterns, highlight healthcare challenges, and suggest potential improvements in public health policy. The primary estimate of this study is the annual number of deaths attributed to various causes.

---

\*Code and data are available at: [https://github.com/ScarletWu/Canada\\_leading\\_cause\\_of\\_death.git](https://github.com/ScarletWu/Canada_leading_cause_of_death.git)

This analysis aims to understand the complexity of mortality data and the variability in death counts. The study employs the simulation of datasets using the negative binomial distribution. This particular distribution is well-suited to model count data, especially in cases where the data exhibit over-dispersion. By simulating datasets, this study aims to gain insights into the potential distributions of deaths by cause and year, allowing for a deeper understanding of the trends and variability in the actual data.

This study employs both Poisson and negative binomial regression models to analyze the Canadian death data. The reason to use both models is to accommodate the variable nature of mortality counts. While Poisson regression is a conventional choice for count data, it assumes equality between the mean and variance of the data, which is often not the case in mortality statistics due to over-dispersion. In contrast, the negative binomial model provides greater flexibility by accommodating over-dispersion, making it a more realistic tool for analyzing the complex nature of mortality data.

The critical comparison between Poisson and negative binomial models in this analysis is not merely a statistical preference but a methodological necessity. It underscores the importance of selecting a model that accurately reflects the data's underlying distribution, ensuring the reliability and validity of the findings. This comparison is instrumental in identifying the model that best captures the nuances of mortality data, thereby providing a solid foundation for concluding the leading causes of death in Canada and the potential implications for public health policies and initiatives.

In summary, this paper comprehensively analyzes mortality data from Canada by using Poisson and negative binomial models to capture the nuances of death-related statistics. The primary estimand—the annual number of deaths by various causes—sets the stage for a detailed exploration of mortality patterns within Canada. Besides quantifying the burden of mortality, this study reveals underlying trends that can inform future health policies and interventions to reduce preventable deaths and improve Canadians' health and well-being.

## 2 Data

This analysis examines mortality data in Canada from 2000 to 2022 obtained from the comprehensive database maintained by Statistics Canada. The data includes the annual number of deaths categorized by different causes, offering a detailed view of mortality trends in the country. The primary variables of interest are the year of death, the cause of death according to the ICD-10 classification, and the total number of deaths attributed to each cause. These variables provide an overview of how mortality patterns have changed over the past two decades, reflecting the impact of healthcare advancements, public health initiatives, and emerging challenges.

To prepare the dataset for analysis, we made significant efforts to ensure that the data was clean and appropriately structured. This included truncating the cause of death descriptions

for readability and consistency, which helped ensure that our analyses were accurate and easy to interpret. We chose the Canadian dataset specifically due to its comprehensive coverage and the high quality of data reporting standards maintained by Statistics Canada. This allowed us to gain a nuanced understanding of mortality within the Canadian context, which may differ from individual provincial trends due to various socio-economic, environmental, and healthcare factors.

To explore and visualize the data, we used a suite of R (R Core Team (2022)) packages, each chosen for its specific capabilities. The `dplyr` (Wickham et al. (2023)) package facilitated efficient data manipulation, while `ggplot2` (Wickham (2023)) helped us visualize trends and patterns in mortality. We used Poisson and negative binomial regression techniques to model the data, and `rstanarm` (Team (2023)) was chosen for its advanced Bayesian modeling capabilities. The `modelsummary` (Arel-Bundock (2023)) package efficiently summarized the results, providing clear and concise insights into the findings. We also utilized the `broom` (Robinson et al. (2023)) and `broom.mixed` (Bolker et al. (2023)) packages to tidy the outputs of our statistical models, making the results more accessible and interpretable.

By systematically examining the variables within the dataset and using rigorous statistical modeling, this analysis provides a deeper understanding of mortality trends in Canada, underscoring the critical role of data-driven approaches in public health planning and evaluation.

### 3 Result

The analysis comprises various elements that are meticulously designed to provide a comprehensive insight into the mortality trends prevailing in Canada. It involves statistical models to analyze these trends and evaluate the model performance.

The analysis begins with the creation of a simulated death dataset that replicates the structure of real-world data, where causes of death over the years 2000 to 2022 are distributed according to a negative binomial distribution. Figure 1 shows the simulated distributions of death for each cause.

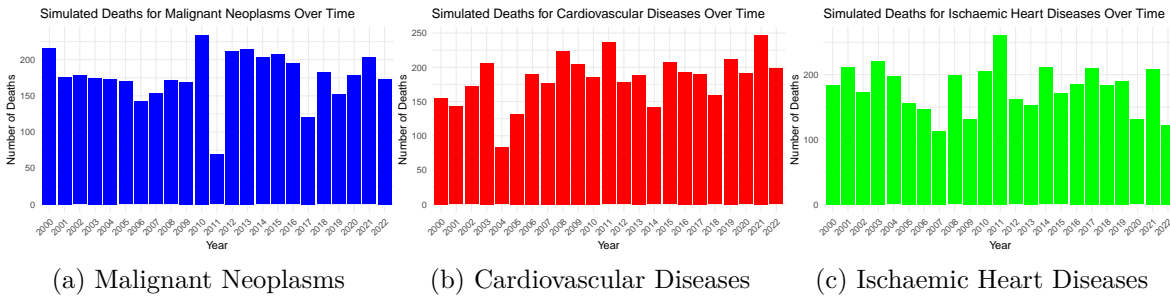


Figure 1: Negative binomial death simulation for each cause of death

Table 2: Summary statistics of the number of yearly deaths, by cause, in Canada

	Min	Mean	Max	SD	Var	N
value	14 466	42 831	82 822	22 480	505 344 387	153

Table 1 ranks the ten leading causes of death in Canada for the year 2022. It lists the specific causes, the number of deaths attributed to each cause, and the rank based on the number of deaths. This snapshot is crucial for understanding the most pressing health threats faced by Canadians in the latest year of the study.

Table 1: Top-ten causes of death in Canada in 2022

Year	Cause	Deaths	Rank	Years
2022	Malignant neoplasms [C00-C97]	82412	1	23
2022	Major cardiovascular diseases...	76639	2	23
2022	Diseases of heart [I00-I09,...]	57357	3	23
2022	Ischaemic heart diseases [I...	34830	4	23
2022	Dementia [F010-F019, F03]	25994	5	6
2022	Unspecified dementia [F03]	23896	6	6
2022	Other forms of chronic isch...	20126	7	23
2022	COVID-19 [U07.1, U07.2, U10.9]	19716	8	3
2022	Malignant neoplasms of trac...	19151	9	23
2022	Other heart diseases [I26-I51]	18913	10	23

Figure 2 shows a series of line graphs, each representing a different cause of death in Canada, such as heart disease, cancer, and dementia. Each line displays the change in the number of deaths from 2000 to 2022, enabling viewers to track how each cause of death has fluctuated over time. One noticeable line represents COVID-19, indicating its emergence and impact in recent years.

Table 2 provides a broad overview of the number of deaths from the top causes over the years. It includes the minimum, maximum, and average number of deaths, standard deviation, variance, and the total count of data points (N). These statistics offer a foundational understanding of the distribution and variability of mortality data.

The analysis uses two statistical methods to study the causes of death. The Poisson model was first applied. However, this model can sometimes be limited due to its assumption of equal mean and variance. To address this limitation and consider over-dispersion in the data, the negative binomial model was also used. This model allows for greater variance than the mean, providing a more detailed understanding of the mortality data. These models were useful in clarifying the relationship between various causes of death and their frequencies over time.

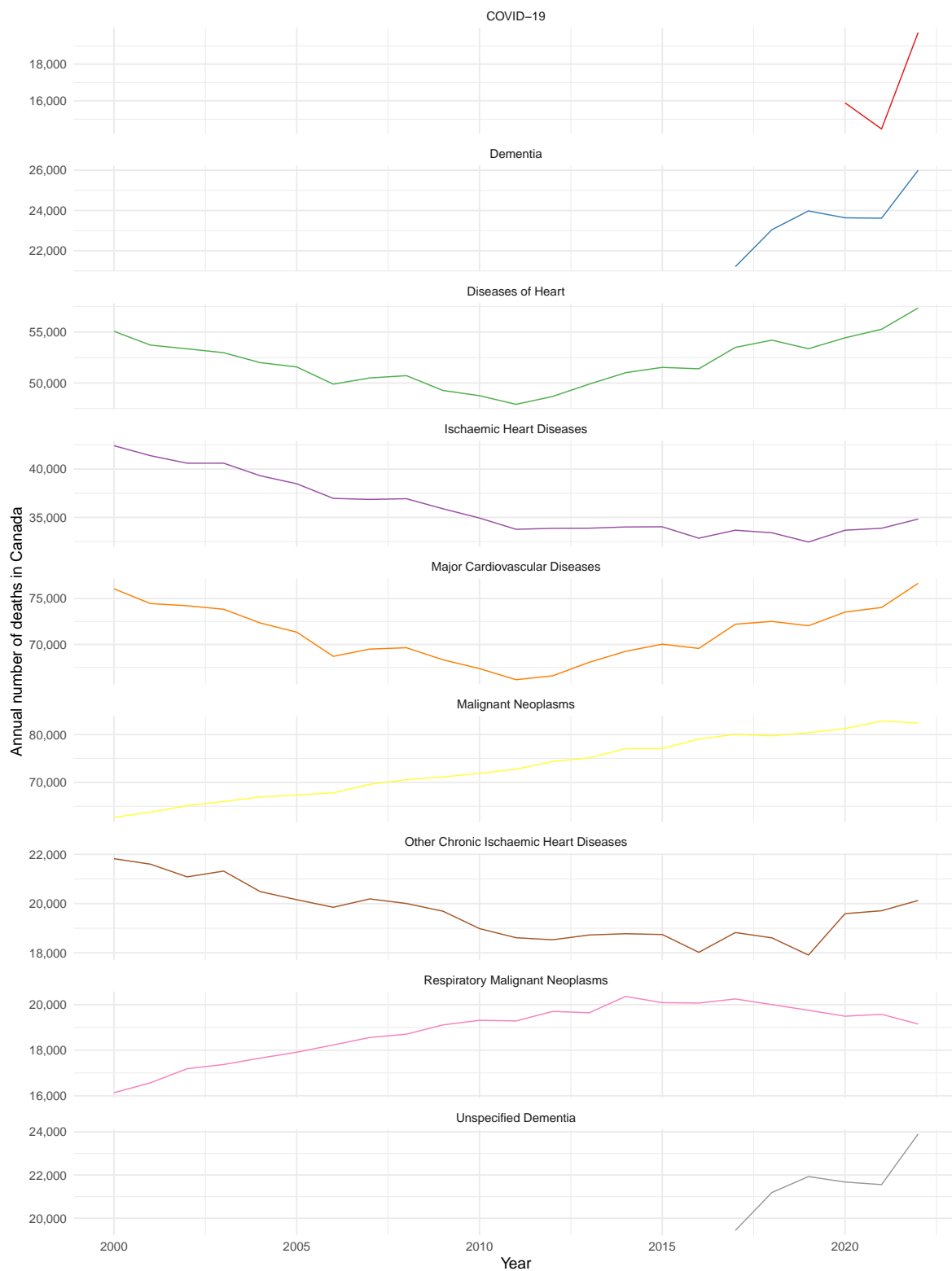


Figure 2

Table 3: Modeling the most prevalent cause of deaths in Canada, 2001-2020

	Poisson	Negative Binomial
Malignant Neoplasms	1.479	1.472 (0.082)
Diseases of Heart	1.137	1.129 (0.083)
Respiratory Malignant Neoplasms	0.123	0.116 (0.082)
Major Cardiovascular Diseases	1.450	1.443 (0.085)
Ischaemic Heart Diseases	0.770	0.761 (0.083)
Unspecified Dementia	0.259	0.253 (0.096)
Num.Obs.	153	153
Log.Lik.	-14 947.833	-1463.624
ELPD	-15 593.5	-1468.3
ELPD s.e.	1752.0	7.1
LOOIC	31 186.9	2936.6
LOOIC s.e.	3504.0	14.3
WAIC	31 838.1	2936.3
RMSE	3153.30	3153.84

The model summary compares the coefficient estimates from both Poisson and negative binomial models for the leading causes of death. The estimates measure the impact each cause has on mortality rates. A higher value suggests a greater influence on increasing death rates. This table is crucial in discerning the relative importance of each cause in the context of mortality (Table 3).

Figure 3 demonstrates how well the Poisson and negative binomial models fit the actual data. The dark line represents the observed data, while the lighter lines or bands show the range of outcomes predicted by the models. A close match between the two indicates a good fit, meaning the model's predictions closely align with the real-world data.

```
poisson <- loo(cause_of_death_poisson, cores = 1)
neg_binomial <- loo(cause_of_death_neg_binomial, cores = 1)
loo_compare(poisson, neg_binomial)
```

```
elpd_diff se_diff
```

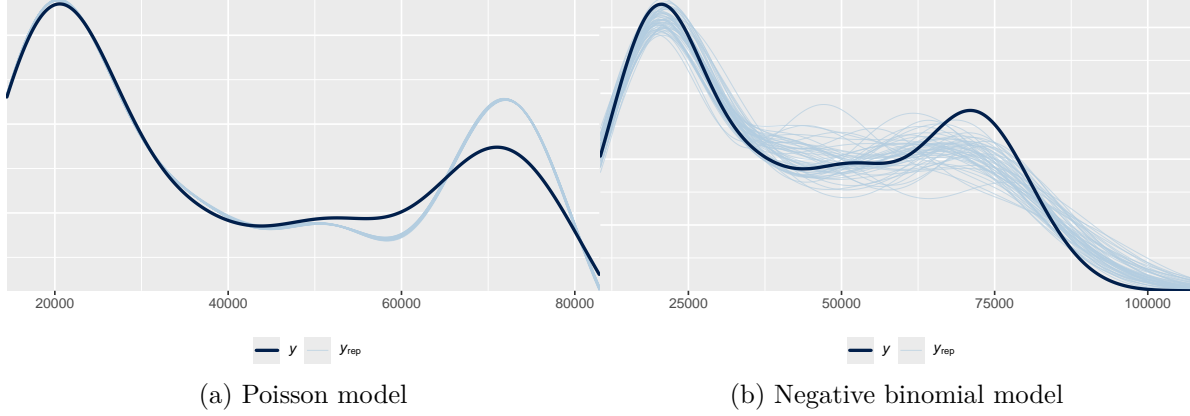


Figure 3: Comparing posterior prediction checks for Poisson and negative binomial models

cause_of_death_neg_binomial	0.0	0.0
cause_of_death_poisson	-14125.2	1748.0

The Leave-One-Out Cross-Validation (LOO-CV) comparison provides a statistical method to determine which model, Poisson or negative binomial, predicts the data more accurately. The result shows the difference in the expected log predictive density (ELPD) and its standard error (SE). In this scenario, the negative binomial model proves to be a superior choice compared to the Poisson model due to its higher ELPD.

**#Discussion** One important lesson we can learn from this analysis is how different diseases affect mortality rates. The comparison of different models shows that certain causes, such as heart diseases and cancer, consistently rank as the leading causes of death. This information emphasizes the chronic nature of these conditions and may also reflect the impact of an aging population on healthcare systems. The inclusion of COVID-19 as a separate category in recent years highlights how emerging health threats can quickly alter mortality rates.

Another valuable insight we gain from this investigation is how mortality causes change over time. The graph that shows the annual number of deaths per cause reveals the dynamics of how these causes evolve. Some causes, such as cancer and heart diseases, show a steady increase that aligns with population growth and aging. On the other hand, causes such as COVID-19 show a sudden spike due to the pandemic. This trend information can inform healthcare planning, such as resource allocation, to address these evolving challenges.

However, this study has limitations. While the models used are adept at handling the data, they may not fully account for all the complexities of real-world scenarios, such as socio-economic factors, healthcare access, and lifestyle changes over time. The Poisson model, in particular, with its assumption of equal mean and variance, may oversimplify the complexity of mortality data. Furthermore, our understanding is limited to the data provided by Statistics Canada, which may have its own reporting biases and gaps.

Moving forward, it is important to continue this research by incorporating more nuanced data that considers the broader determinants of health. Future studies could benefit from a more granular approach, such as analyzing sub-populations, to understand disparities in mortality rates. Additionally, as new health threats emerge and societal factors evolve, continuous updating and refining of models will be crucial. Longitudinal studies could also shed light on the lifetime risks of various demographics, providing a more comprehensive picture of mortality in Canada.

In conclusion, this study is a starting point for a deeper understanding of mortality and its causes. It highlights the persistent and emerging health threats that dominate mortality statistics and emphasizes the importance of tailored health interventions. Further research can provide more detailed insights, ultimately guiding better health policies and outcomes.

#### #Reference

- Arel-Bundock, Vincent. 2023. *Modelsummary: Summary Tables and Plots for Statistical Models and Data: Beautiful, Customizable, and Publication-Ready*. <https://vincentarelbundock.github.io/modelsummary/>.
- Bolker, Ben et al. 2023. *Broom.mixed: Tidying Methods for Mixed Models*. <https://github.com/bbolker/broom.mixed>.
- R Core Team. 2022. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.
- Robinson, David, Alex Hayes, Simon Jackson, Max Kuhn, and Hadley Wickham. 2023. *Broom: Convert Statistical Analysis Objects into Tidy Tibbles*. <https://broom.tidyverse.org>.
- Team, Stan Development. 2023. *Rstanarm: Bayesian Applied Regression Modeling via Stan*. <https://mc-stan.org/users/interfaces/rstanarm>.
- Wickham, Hadley. 2023. *Ggplot2: Create Elegant Data Visualisations Using the Grammar of Graphics*. <https://ggplot2.tidyverse.org>.
- Wickham, Hadley, Romain François, Lionel Henry, and Kirill Müller. 2023. *Dplyr: A Grammar of Data Manipulation*. <https://dplyr.tidyverse.org>.