

# Number of COVID-19 Cases in Canada By Region, Age, and Sex\*

Predicting the Factors that Influence the Probability of High Risk of COVID-19

Bruce Zhang

November 26, 2024

This paper examines the demographic and regional factors influencing high-risk classifications for COVID-19 in Canada using logistic regression models, with predictors including age, sex, and region. The analysis reveals that rural regions, males, and specific age intervals exhibit higher probabilities of being classified as high risk, emphasizing disparities in healthcare access and outcomes. These findings highlight the importance of targeted public health interventions, such as resource allocation and vaccination prioritization, to mitigate these risks and can potentially guide interventions for future pandemics. By identifying key predictors and their implications, this study provides a scalable framework for improving pandemic preparedness and understanding health disparities.

---

\*Code and data are available at: [https://github.com/brucejczhang/covid\\_data](https://github.com/brucejczhang/covid_data).

# Table of contents

|          |  |           |
|----------|--|-----------|
| <b>1</b> | <b>Introduction</b>  | <b>3</b>  |
| <b>2</b> | <b>Data</b>  | <b>4</b>  |
| 2.1      | Overview . . . . .   | 4         |
| 2.2      | Measurement . . . . .  | 4         |
| 2.3      | Data Cleaning . . . . .                                      | 5         |
| 2.4      | Outcome Variables and Predictor Variables . . . . .          | 5         |
| 2.4.1    | Number of Cases and Number of Deaths . . . . .               | 6         |
| 2.4.2    | Case Fatality Rate . . . . .                                 | 6         |
| 2.4.3    | Time . . . . .   | 6         |
| 2.4.4    | Age, Sex, and Region . . . . .                               | 7         |
| <b>3</b> | <b>Model</b>   | <b>8</b>  |
| 3.1      | Model set-up . . . . .                                       | 8         |
| 3.2      | Model justification . . . . .                                | 9         |
| 3.3      | Assumptions and Limitations . . . . .                        | 10        |
| 3.4      | Model Validation . . . . .                                   | 10        |
| <b>4</b> | <b>Results</b>   | <b>11</b> |
| 4.1      | Predicted Probability of High Risk by Region . . . . .       | 11        |
| 4.2      | Predicted Probability of High Risk by Sex . . . . .          | 12        |
| 4.3      | Predicted Probability of High Risk by Age Interval . . . . . | 12        |
| <b>5</b> | <b>Discussion</b>  | <b>14</b> |
| 5.1      | Interpreting High Risk Predictions . . . . .                 | 14        |
| 5.2      | Key Future Implications . . . . .                            | 14        |
| 5.3      | Limitations . . . . .  | 14        |
| 5.4      | Next Steps . . . . .   | 15        |
| <b>A</b> | <b>Appendix</b>  | <b>16</b> |
| A.1      | Appendix A: Survey and Sampling Details . . . . .            | 16        |
| A.1.1    | Collection Process . . . . .                                 | 16        |
| A.1.2    | Sampling and Representativeness . . . . .                    | 16        |
| A.1.3    | Observational Nature and Limitations . . . . .               | 17        |
| A.2      | Appendix B: Model Selection Details . . . . .                | 17        |
| A.3      | Appendix C: Model Validation Details . . . . .               | 18        |
| A.3.1    | ROC Curve and AUC . . . . .                                  | 18        |
| A.3.2    | Confusion Matrix . . . . .                                   | 18        |
| A.3.3    | Odds Ratio with Confidence Intervals . . . . .               | 19        |
|          | <b>References</b>  | <b>20</b> |

# 1 Introduction

The COVID-19 pandemic is something that is difficult to forget. Although it is now in the past, the effects that it had on society and how we function as individuals is still profound. At the same time, the statistical information that sprouted from the pandemic has a high value for analysis. This is not only to better understand the pandemic and how it affected populations but also to form a strong idea of its patterns to suggest tactics for dealing with pandemics shall something similar happen in the future.

Studies have been done in recent years in attempts to utilize the data that was gathered during COVID times to better advise future preventions and treatment for pandemic situations. Many studies look into the relationship between demographic information of populations and their likelihood of getting diagnosed or fatality. One study found that communities on hillsides and other locations with enhanced ventilation and oxygen had reduced numbers of COVID cases (Sharma et al. 2023). Other studies have found relationships between exercise and cardiovascular function and the likelihood of getting COVID (Yang et al. 2024). Although many relationships have been examined, there has been a lack of studies focusing on a broad set of demographic traits such as age and sex. There has also not been cross comparisons of region and demographic aspects. Another gap is the lack of country-specific analysis, particularly Canada.

This analysis looks at COVID-19 data from 2022 in Canada and analyzes how the different population and individual-level variables, such as age, sex, and region of life, influenced the risk of being diagnosed with COVID and the risk of fatality. The paper uses the demographic data gathered from 2022 to predict the characteristics of populations and individuals that may be at higher risk of fatality once they have been diagnosed with COVID. This analysis may have broader implications for preventative measures for flues, other diseases, and pandemics in the future.

The analysis focuses on the outcome variable of predicted probability of high risk, which is a function of number of cases and number of deaths per datapoint and is further defined in Section 3. The model predicts the probability of high risk as a result of a series of predictor variables including region, sex, and age of the reported data.

The analysis revealed that there are disparities in the risk of severe COVID-19 outcomes across demographic and regional categories. Logistic regression models demonstrated that individuals in rural regions consistently faced higher probabilities of being classified as high risk, with males and specific age intervals further amplifying this vulnerability. These findings were supported by consistent trends in summary statistics and model visualizations, which revealed the interplay between predictors such as age, region, and sex in influencing COVID-19 outcomes.

Understanding the drivers of COVID-19 severity can craft equitable and effective public health responses. By identifying which demographic and regional groups face the greatest risks, the results of the analysis can guide resource distribution, including vaccination programs and

medical support, to areas where they are most needed. These are applicable not only to COVID-19 but also to future public health challenges.

The remainder of this paper is structured as follows. The data section (Section 2) highlights the characteristics of the dataset. The section summarizes the data through a series of summary statistics (Figure 2) and represents the data in a visual way where specific trends and patterns can be observed (Figure 3, Figure 4). The model section (Section 3) creates a logistic regression model predicting the probability of high risk depending on the region, sex, and age of the data point. This section presents the mathematical basis of the model and justifies the selection for the logistic regression model over others. The results section (Section 4) includes model figures that display the predicted probabilities of high risk for the predictor variables respectively. The discussion section (Section 5) goes over how the results from the modeling can be interpreted in the context of the real world and how these findings can be utilized to advise public health measures. It also discusses limitations of the analysis and future directions.

## 2 Data

### 2.1 Overview

I use the statistical programming language R (R Core Team 2023) to analyze the data and to create graphs and models. The packages that were used include tidyverse (Wickham, Averick, et al. 2023), tidyr (Wickham, Henry, and Vaughan 2023), dplyr (Wickham, François, et al. 2023), caret (Kuhn 2023), and pROC (Robin et al. 2011). My data (Open Science Framework (OSF) 2022) was obtained from COVerAGE-DB, which was housed in Open Science Framework. The data analysis was conducted based on the guidance of Alexander (2023).

### 2.2 Measurement

The dataset used in this analysis was gathered by COVerAGE-DB, a database that focuses on COVID related data. COVerAGE-DB often gathers data relating to COVID such as the number of cases, deaths, tests, and vaccinations through governmental institutes such as health ministries and statistical offices. This data is then organized by other variables such as sex, region, country, and age and presented to the general public.

The dataset reflects a transformation of real-world phenomena—COVID-19 cases, deaths, and testing rates—into structured entries that can be analyzed. Each entry in the dataset represents aggregated information collected from health departments, governmental organizations, and research institutions worldwide. These organizations report COVID-19 statistics at varying levels of granularity, such as daily case counts or weekly summaries, which are then standardized and compiled into the dataset. This collective and comprehensive effort of data collection allows the real-world challenges of facing COVID-19 to be transformed into quantitative measures that can be analyzed and modeled.

## 2.3 Data Cleaning

The original dataset contained data from five different countries including. For the purpose to focus on data in Canada for this analysis, the datapoints for other countries were removed. Further cleaning was done to remove any problematic datapoints such as when the number of deaths exceeded the number of cases.

## 2.4 Outcome Variables and Predictor Variables

The cleaned dataset contains variables including the country, region, date, sex, and age where the number of cases, deaths, and tests were reported. Figure 1 shows a sample of the dataset, presenting the first six rows of the analysis data. The outcome variables highlighted include the total number of cases and the total number of deaths.

| Country | Region | Date       | Sex | Age | Cases | Deaths | Tests | CFR       |
|---------|--------|------------|-----|-----|-------|--------|-------|-----------|
| Canada  | Urban  | 2022-01-01 | m   | 0   | 4091  | 215    | 44602 | 5.255439  |
| Canada  | Rural  | 2022-01-01 | m   | 0   | 14555 | 269    | 42100 | 1.848162  |
| Canada  | Urban  | 2022-01-01 | m   | 0   | 4981  | 407    | 39566 | 8.171050  |
| Canada  | Rural  | 2022-01-01 | m   | 0   | 2960  | 495    | 45319 | 16.722973 |
| Canada  | Urban  | 2022-01-01 | m   | 0   | 18488 | 394    | 46584 | 2.131112  |
| Canada  | Rural  | 2022-01-01 | m   | 0   | 10762 | 391    | 11835 | 3.633154  |

Figure 1: First six rows of the dataset after cleaning

| Region | Sex | Mean Age | Total Cases | Total Deaths | Average Cases | Average Deaths | CFR (%) |
|--------|-----|----------|-------------|--------------|---------------|----------------|---------|
| Rural  | f   | 50.05    | 6763540     | 156814       | 10373.53      | 240.51         | 2.32    |
| Rural  | m   | 49.83    | 6869820     | 163430       | 10504.31      | 249.89         | 2.38    |
| Urban  | f   | 50.02    | 6769033     | 165648       | 10366.05      | 253.67         | 2.45    |
| Urban  | m   | 50.12    | 6630777     | 160191       | 10138.80      | 244.94         | 2.42    |

Figure 2: Summary Statistics of COVID-19 Cases in Canada

Figure 2 further breaks down the data summarizing the mean number of cases and deaths by region and sex. This creates an all-rounded summary of the dataset and gives a preliminary idea on the differences in cases and deaths depending on the categories of sex and region.

### 2.4.1 Number of Cases and Number of Deaths

The total and average numbers of cases and deaths are key outcome variables that can be correlated with the categorical predictors. These number can demonstrate the differences in COVID-19 susceptibility and threat level based on the characteristics of the region and the individual.

### 2.4.2 Case Fatality Rate

The case fatality rate (CFR) is an additional aspect of the summary that was calculated to standardize the deaths and cases by category. This can be used as a processed outcome variable that can better represent the dataset and aid modeling later on. The total number of deaths and cases may lead to inaccurate representations of a region, sex, or age group due to the absolute number of individuals that fall within the category. The CFR value allows better comparison across categories.

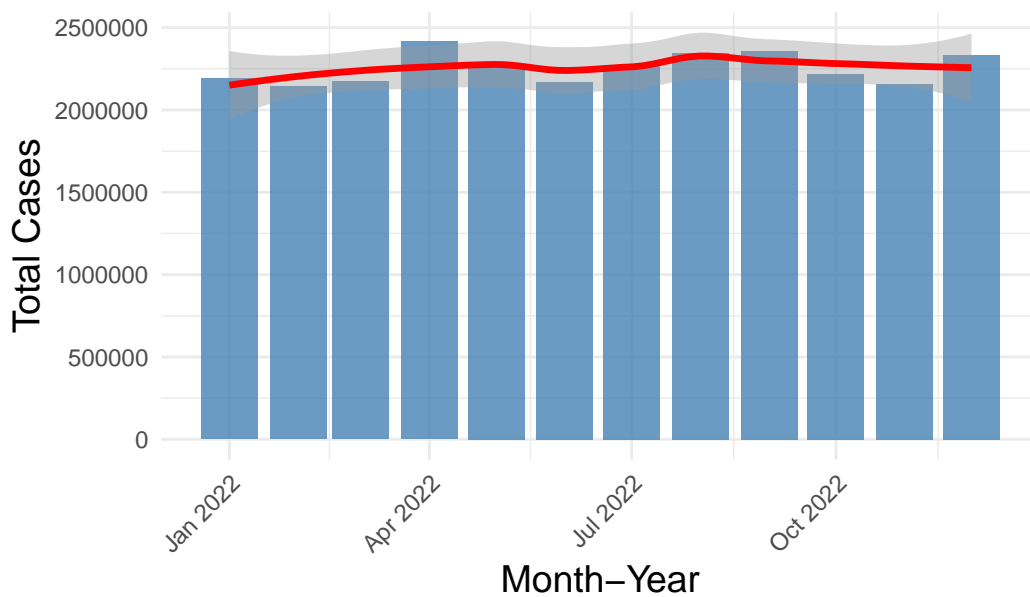


Figure 3: Number of cases per month in 2022 in Canada

### 2.4.3 Time

The time of year can also be correlated with the likelihood of getting COVID-19 independent of other predictor variables. As shown in Figure 3, there is a weak pattern of increased total

number of cases in the early summer months of May and June. The number of cases in winter and near winter months, specifically January and October, seem to be lower.

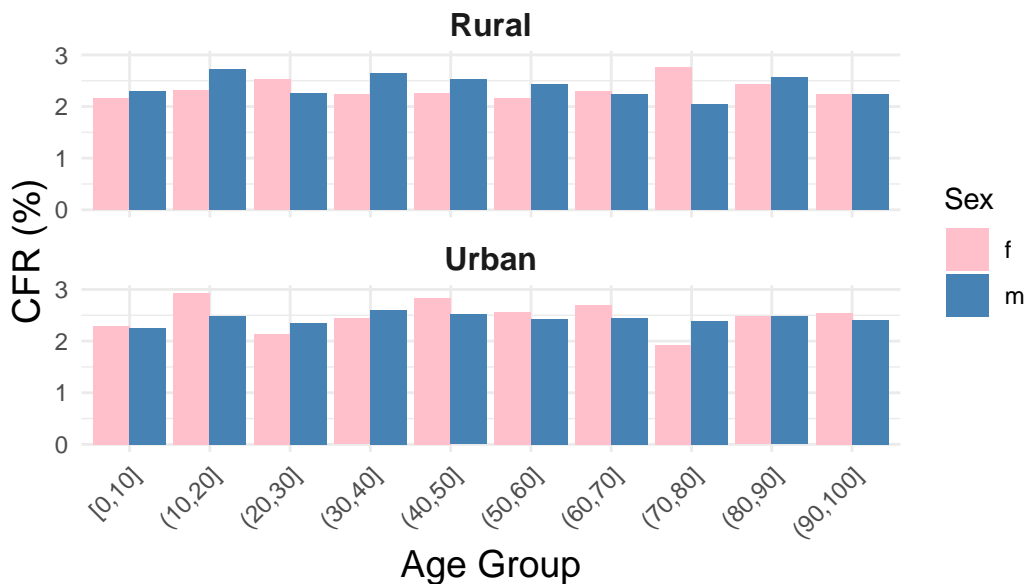


Figure 4: Case Fatality Rate (CFR) by Age Group, Region, and Sex

#### 2.4.4 Age, Sex, and Region

Age, sex, and region are key predictor variables that are examined in this analysis. In Figure 4, these variables are organized relative to the CFR, which calculates a rate based on the number of cases and number of deaths and allows standardized comparisons across different predictor variables. Figure 4 shows that the urban and rural regions have similar patterns in CFR across different age groups. The CFR for males tend to be higher more often than not compared to the CFR of females, indicating that males may have a higher likelihood of death after diagnosis of COVID in Canada. The pattern of CFR in relation to age seem to vary quite significantly between urban and rural regions and between male and females. For females, the highest CFR seems to be for individuals aged under 30 and for those around 70 to 80 regardless of region. Male CFR values peak at age interval 20 to 30 for urban regions and 30 to 40 for rural regions.

Figure 5 shows a more direct comparison of the CFR between urban and rural regions of Canada, combining different age groups and sexes. This figure shows that the CFR in rural regions is noticeably higher than that of urban. Although the difference is small, it is still a

considerable size difference when considering the number of people that a small percentage can be responsible for, based on the summary statistics from Figure 2.

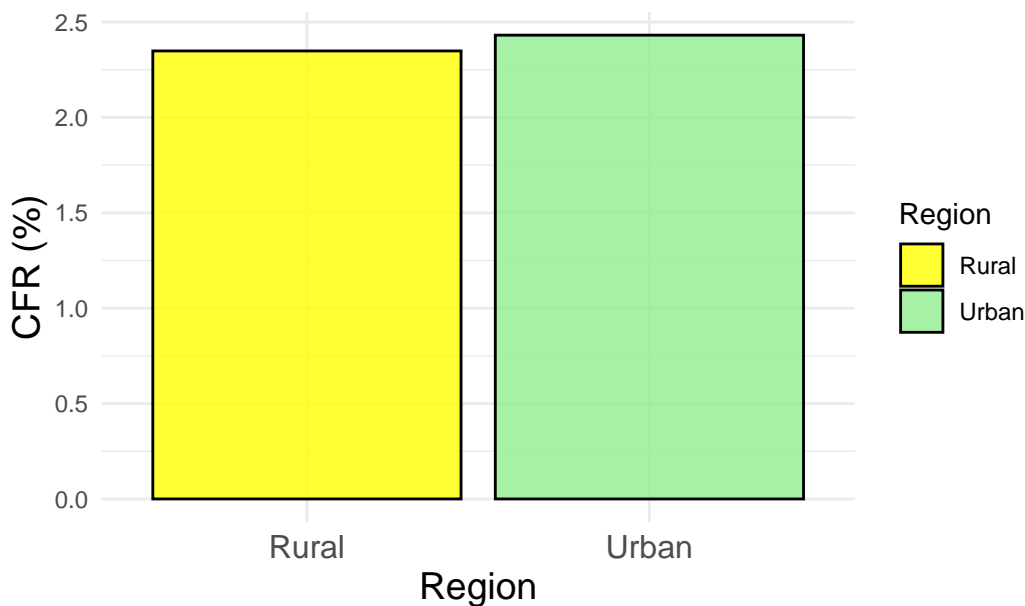


Figure 5: Case Fatality Rate (CFR) by Region in 2022 in Canada

### 3 Model

The goal of our modelling strategy is twofold. Firstly,...

Here we briefly describe the Bayesian analysis model used to investigate... Background details and diagnostics are included in [?@sec-model-details](#).

#### 3.1 Model set-up

A logistic model has been selected as part of this analysis to examine the relationship between region, sex, age, and the number of cases and deaths of COVID-19.

The logistic regression model predicts the probability of a region or individual being classified as “high risk” based on key predictors. In this case, “high risk” is defined as regions where the death rate (CFR) exceeds the median death rate across the dataset.



The CFR is calculated as:

$$\text{CFR} = \frac{\text{Deaths}}{\text{Cases} + 10^{-6}} \times 100 \quad (1)$$

The box-and-whisker plot visualizes the distribution of predicted probabilities ( $p$ ) across different age intervals. Each box represents the interquartile range (IQR) of predicted probabilities, with the line inside the box indicating the median predicted probability for that age interval. Whiskers extend to 1.5 times the IQR, and any points outside this range are considered outliers. This visualization highlights the variability of predicted probabilities within each age group.

The logistic regression model equation is:

$$\text{logit}(p) = \ln\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 \cdot \text{Tests} + \beta_2 \cdot \text{Age} + \beta_3 \cdot \text{Region} + \epsilon \quad (2)$$

where:

- $p$ : Predicted probability of being high risk ( $p$  = Predicted Probability).
- $\text{logit}(p)$ : Log-odds of the high-risk classification.
- $\beta_0$ : Intercept of the model.
- $\beta_1, \beta_2$ : Coefficients for the predictors (e.g., number of tests, age interval).
- $\epsilon$ : Error term.

This model explains how testing rates and age intervals affect the likelihood of a region or demographic being classified as high risk. The box-and-whisker plot shows the variability of predicted probabilities across age groups, helping to identify specific age intervals with consistently higher or lower risks.

We run the model in R (R Core Team 2023) using the ‘pROC’ (Robin et al. 2011) and ‘caret’ (Kuhn 2023) packages to construct and test the model.

## 3.2 Model justification

For the analysis, a logistic regression model was chosen to predict the probability of a region or individual being classified as “high risk” based on key predictors such as age, sex, region, and the number of COVID-19 tests conducted. Logistic regression is particularly well-suited for this dataset because the outcome variable, high risk, is binary (e.g., classified as high risk or not). This makes logistic regression the most appropriate modeling technique as it directly estimates the probability of the binary outcome through the log-odds transformation. Compared to other models like linear regression, logistic regression avoids the issue of nonsensical

predictions, such as probabilities below 0 or above 1, ensuring the interpretability of the results. Most importantly, logistic regression can handle both categorical and continuous predictors effectively. For this particular dataset, the key predictor variables are all categorical, even age is categorized by intervals. Therefore, it was extremely necessary for the model to be able to handle categorical data, which the logistic regression model is capable of.

Another key advantage of the logistic regression model in this context is that it provides interpretable coefficients, which allow for an understanding of how specific predictors influence the odds of being high risk. This interpretability is crucial for public health applications, where stakeholders need clear and actionable insights. Furthermore, logistic regression is computationally efficient, making it ideal for datasets with relatively straightforward structures like this one. Thus, the logistic regression model strikes an appropriate balance between accuracy, interpretability, and simplicity, making it the most suitable choice for this analysis.

The details on model selection can be found in the appendix in [Section A.2](#)

### 3.3 Assumptions and Limitations

The logistic regression model assumes that observations within the dataset are independent, meaning that each data point does not influence others. However, this assumption could be violated if individuals within the same region share similar environmental, healthcare, or socioeconomic factors that contribute to their risk classification. Additionally, the model assumes that the predictors, such as age, sex, and region, are not highly correlated. Significant multicollinearity between these variables could undermine the stability of the model coefficients and reduce interpretability. Finally, the classification of “high risk” as regions or individuals exceeding the median death rate simplifies a complex phenomenon. This threshold may fail to capture subtler variations in risk factors across different demographic and geographic contexts.

One key limitation of logistic regression is its inability to account for nuanced interactions between predictors without explicitly adding interaction terms. For instance, the combined effects of age and sex on risk may not be fully captured in this model. Additionally, the model assumes a predefined functional form, which limits its flexibility in identifying nonlinear relationships between predictors and outcomes. As a result, more complex patterns in the data may remain undetected, suggesting the need for alternative modeling approaches such as random forests or decision trees in future studies.

### 3.4 Model Validation

For the model validation process of the selected logistic regression model, multiple validation techniques were used. Firstly, a ROC curve was constructed to visualize the trade-off between sensitivity and specificity. The ROC curve showed a high curve that deviated from the linear

line representing random chance between sensitivity and specificity. The AUC value accompanying the ROC curve was above 0.8 for all three prediction models constructed for different variables.

A confusion matrix was also constructed to assess classification accuracy, precision, and recall by comparing the predicted classes to the actual outcomes. Values of above 0.75 were yielded for the accuracy rating. Lastly, the odds ratios with confidence intervals provided insights into the strength and direction of the predictors' effects.

In general, these results demonstrated the suitability and strong effect that the logistic regression model has to predict probability of high risk based on the data.

Details of the model validation are provided in the appendix (Section [A.3](#)).

## 4 Results

The results of this analysis are summarized in Figure [6](#), Figure [7](#), and Figure [8](#) to give an all-rounded picture of the predicted probability of high risk based on the different predictor variables.

### 4.1 Predicted Probability of High Risk by Region

Figure [6](#) models the predicted probability of high risk for different regions, urban or rural. This is represented by a faceted density plot.

Based on these results, the predicted probability of high risk is higher in rural areas than urban areas and the predicted probability of low risk is higher in urban areas than rural. This is indicated by the height of the peaks at the higher and lower end of the scale, where the peaks for urban regions are taller at the lower end of the probability of high risk scale and shorter at the higher end. The height of the peaks indicate the density of datapoints that fall around the respective predicted value.

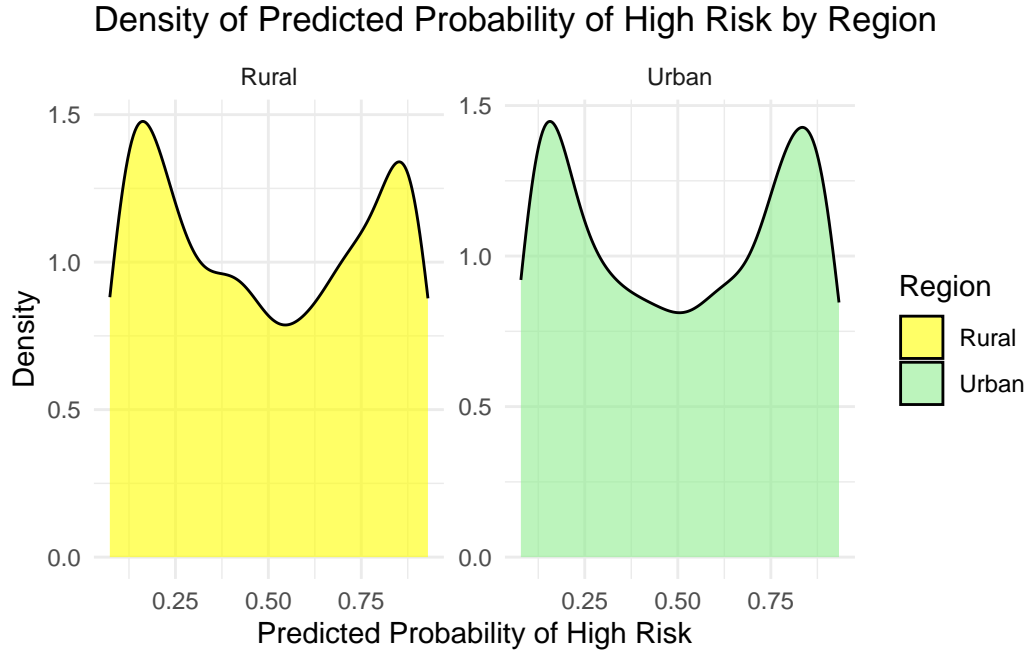


Figure 6: Predicted probability of high risk of COVID-19 by region (Urban vs. Rural)

## 4.2 Predicted Probability of High Risk by Sex

Figure 7 compares the predicted probability of high risk based on biological sex, representing the model in a similar box-and-whiskers plot method. As shown in the figure, males (represented by 'm') have a noticeably higher predicted probability of high risk than females (represented by 'f'). The predicted probability of a male qualifying as high risk for COVID-19 is above 50%, while the number is below 50% for females.

## 4.3 Predicted Probability of High Risk by Age Interval

Figure 8 shows the predicted probability of high risk for different age groups. The age groups are intervals of 10 from 0 to 100 as classified by the original dataset. As shown in the figure, age intervals 10 to 20 and 30 to 40 have a significantly higher predicted probability of high risk than other age groups. This means that it is more likely (around 55% and 57% chance) that someone within these age ranges will qualify as a high risk than the other age groups. The age group with the lowest predicted probability of high risk is 0 to 10 year-olds, yielding a 38% chance median of being predicted as high risk.

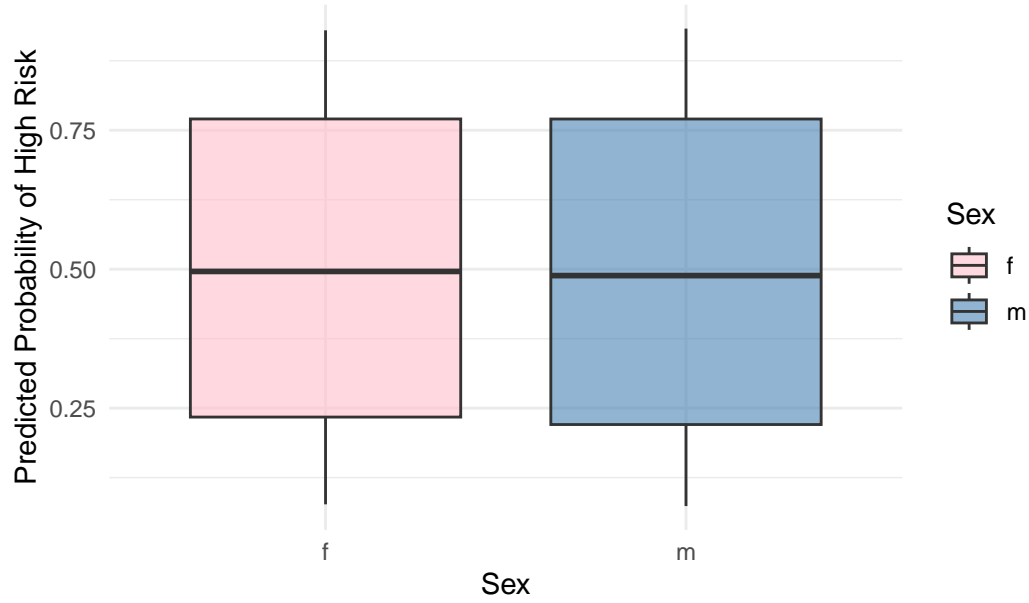


Figure 7: Predicted probability of high risk based on sex

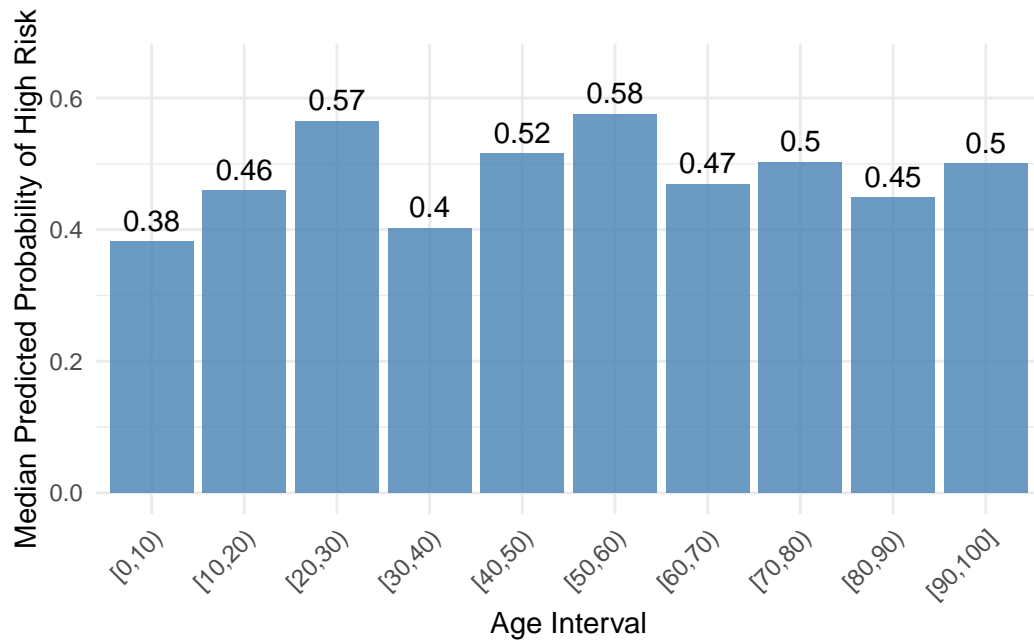


Figure 8: Predicted probability of high risk based on age intervals

## 5 Discussion

### 5.1 Interpreting High Risk Predictions

The results of this analysis reveal critical insights into how demographic and regional factors influence the likelihood of being classified as high risk for COVID-19. The logistic regression model, designed to predict the probability of high risk, effectively highlighted the relative importance of age, sex, and region in shaping these probabilities. Notably, rural regions consistently showed higher predicted probabilities of high risk compared to urban regions. This disparity underscores the challenges faced by rural areas, such as limited access to healthcare infrastructure, which may contribute to higher case fatality rates (CFR).

Similarly, the analysis revealed that males are more likely to be classified as high risk than females, aligning with previous research suggesting that biological differences, comorbidities, and health-seeking behaviors could account for this discrepancy. Age also played a pivotal role, with certain age intervals, such as 10–20 and 30–40, showing significantly higher probabilities of high risk. This pattern may reflect complex interactions between behavioral, immunological, and societal factors influencing COVID-19 outcomes. Understanding these relationships is crucial for tailoring public health interventions and resource allocation strategies.

### 5.2 Key Future Implications

The findings of this study have implications for future pandemic preparedness and response. By identifying demographic and regional factors associated with higher probabilities of high risk, policymakers can develop targeted interventions to mitigate disparities. For example, enhancing healthcare infrastructure in rural areas and prioritizing vaccination and testing campaigns for high-risk groups can help reduce the impact of future pandemics.

Additionally, the use of logistic regression models to predict high-risk classifications provides a robust framework for decision-making. This approach could be extended to other diseases or health outcomes, offering a scalable and interpretable method for public health planning. The analysis also highlights the importance of maintaining high-quality, standardized datasets to facilitate cross-regional and cross-demographic comparisons.

### 5.3 Limitations

Despite its strengths, this study has several limitations that should be addressed in future research. First, the analysis is limited to Canada, which may restrict the generalizability of the findings to other countries with different healthcare systems, cultural contexts, or pandemic responses. Second, the dataset used for this analysis only includes data from 2022, which may not capture temporal variations in COVID-19 patterns. Future studies could incorporate longitudinal data to examine how the predictors of high risk evolve over time.

Another limitation is the reliance on aggregate data, which may obscure individual-level variability and interactions. For example, the interplay between age, sex, and comorbidities at the individual level could yield more nuanced insights. Additionally, the classification of rural and urban regions may oversimplify the diversity within these categories, such as the distinction between suburban and remote rural areas.

Thirdly, the data obtained for COVID-19 is obviously in the past and no longer an issue that is of primary concern. This issue is in line with the generalizability issue due to limited surveying countries. The patterns and models may be limited in generalizability due to the sole focus on COVID-19.

## **5.4 Next Steps**

Building on the findings of this study, future research should explore several directions. First, expanding the analysis to include additional countries or regions would provide a more comprehensive understanding of the factors influencing COVID-19 outcomes. Second, incorporating additional variables, such as socioeconomic status, vaccination rates, and comorbidities, could enhance the predictive power of the models.

Another next step that could be taken is to use the analysis and modeling on a broader series of pandemics and diseases such as seasonal flues, SARS outbreak in the 2000s, and others that may show similar trends. This will allow a broader generalizability of the demographic patterns and allow the model to better predict probability of high risk based on demographic data. Finally, public health practitioners and researchers should collaborate to ensure that findings are translated into actionable policies and interventions, bridging the gap between data analysis and real-world impact.

## A Appendix

### A.1 Appendix A: Survey and Sampling Details

The dataset used in this analysis is sourced from the COVerAGE-DB, a comprehensive repository of COVID-19 demographic data. The database compiles data from a diverse set of sources, including government health agencies, international organizations, and other reputable institutions, ensuring a broad representation of regions and demographics. Each entry in the dataset represents an aggregation of cases, deaths, and demographic information categorized by age, sex, and region.

#### A.1.1 Collection Process

Data collection follows a systematic approach. For primary sources, COVerAGE-DB retrieves data directly from official government dashboards, public health reports, and institutional publications ((cite-dowd2020?)). This ensures that the dataset reflects real-world phenomena with minimal distortion.

The database ensures standardization by processing the raw data from various sources through rigorous standardization to ensure comparability. For example, age intervals are harmonized across sources to create a unified structure suitable for cross-regional analysis.

The database also ensures quality Control through running automated checks and manual reviews to detect inconsistencies, missing values, and potential errors in data entry. Such measures are crucial given the heterogeneity in reporting standards across countries and regions.

Lastly, the update frequency is ensured so the data is as relevant and recent as possible. The database is updated regularly to reflect new data releases, ensuring timeliness and relevance for ongoing research.

#### A.1.2 Sampling and Representativeness

While the dataset covers a broad range of countries and regions, its representativeness depends on the completeness and accuracy of the underlying sources. For instance, regions with underdeveloped health reporting infrastructure might have lower data quality or incomplete records. Additionally, variations in testing capacity and reporting standards can introduce biases, particularly in the calculation of derived measures such as case fatality rates (CFR).



### A.1.3 Observational Nature and Limitations

The data is observational, collected without experimental controls. As a result, confounding factors such as healthcare system differences, socioeconomic disparities, and public health policies must be carefully considered when interpreting results ((cite-ioannidis2021?)). For instance, differences in testing rates across regions can influence the apparent number of cases and consequently impact measures like CFR.

## A.2 Appendix B: Model Selection Details

During the model selection process, model figures were generated for the linear regression model. As shown in Figure 9, the data was difficult to represent and relationships were difficult to identify due to the categorical nature of the variables. This led to attempts to model the data in different ways. The second method that was attempted was using a logistic regression model, which was able to represent the data well and make predictions about high risk. These figures are shown within the Section 4 of the paper.

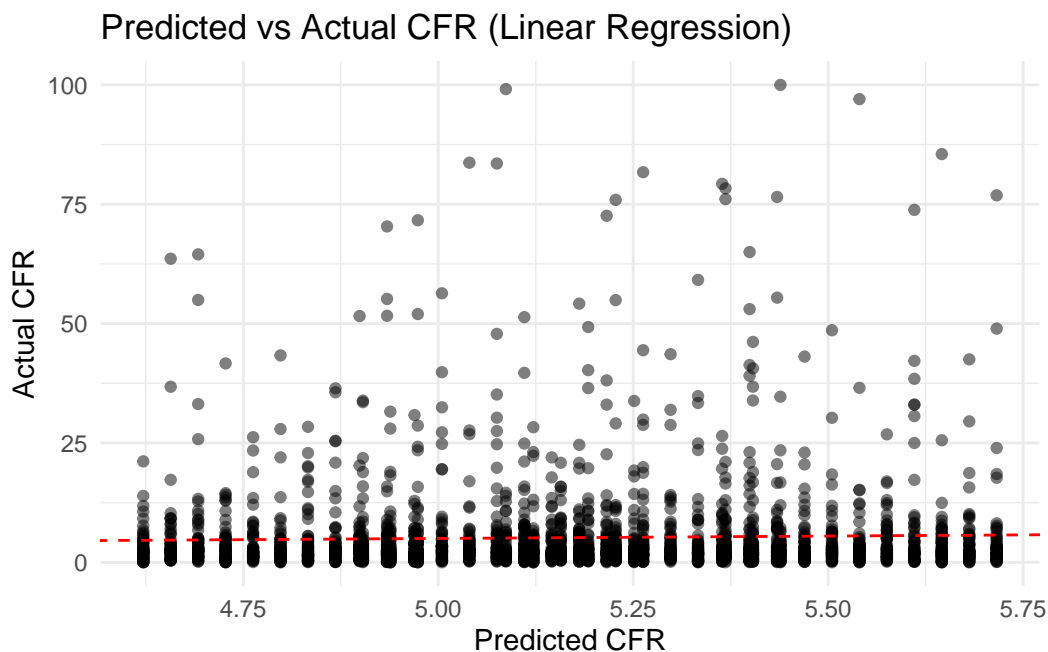


Figure 9: Linear Regression Model for dataset

## A.3 Appendix C: Model Validation Details

### A.3.1 ROC Curve and AUC

AUC: 0.8202962

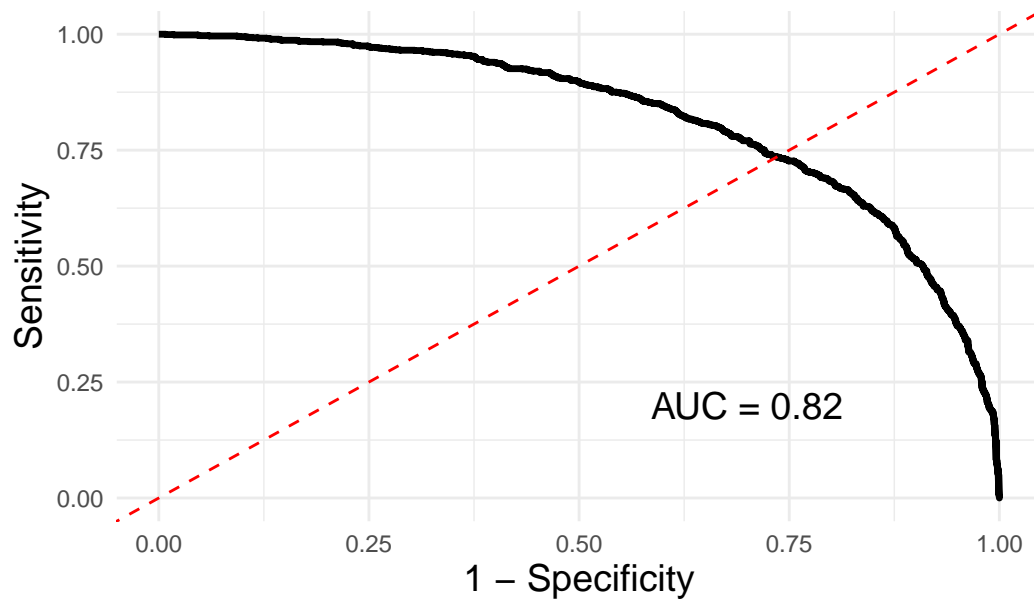


Figure 10: ROC Curve and AUC Value for Logistic Regression Model

### A.3.2 Confusion Matrix

Confusion Matrix and Statistics

|            | Reference |     |
|------------|-----------|-----|
| Prediction | 0         | 1   |
| 0          | 975       | 353 |
| 1          | 332       | 953 |

Accuracy : 0.7378  
95% CI : (0.7205, 0.7546)  
No Information Rate : 0.5002  
P-Value [Acc > NIR] : <2e-16

Kappa : 0.4757  
 Mcnemar's Test P-Value : 0.4448  
 Sensitivity : 0.7460  
 Specificity : 0.7297  
 Pos Pred Value : 0.7342  
 Neg Pred Value : 0.7416  
 Prevalence : 0.5002  
 Detection Rate : 0.3731  
 Detection Prevalence : 0.5082  
 Balanced Accuracy : 0.7378  
 'Positive' Class : 0

### A.3.3 Odds Ratio with Confidence Intervals

| Predictor            | Odds Ratio | Standard Error | Statistic | P Value | Lower 95% CI | Upper 95% CI |
|----------------------|------------|----------------|-----------|---------|--------------|--------------|
| (Intercept)          | 10.324     | 0.199          | 11.715    | 0.000   | 7.007        | 15.311       |
| Sexm                 | 0.975      | 0.094          | -0.273    | 0.785   | 0.810        | 1.173        |
| RegionUrban          | 1.008      | 0.094          | 0.087     | 0.930   | 0.838        | 1.213        |
| Age_Interval[10,20)  | 1.035      | 0.222          | 0.156     | 0.876   | 0.671        | 1.599        |
| Age_Interval[20,30)  | 1.257      | 0.220          | 1.041     | 0.298   | 0.818        | 1.936        |
| Age_Interval[30,40)  | 0.936      | 0.223          | -0.297    | 0.766   | 0.604        | 1.449        |
| Age_Interval[40,50)  | 1.320      | 0.222          | 1.248     | 0.212   | 0.854        | 2.042        |
| Age_Interval[50,60)  | 1.614      | 0.222          | 2.155     | 0.031   | 1.045        | 2.499        |
| Age_Interval[60,70)  | 1.261      | 0.222          | 1.043     | 0.297   | 0.816        | 1.950        |
| Age_Interval[70,80)  | 1.380      | 0.222          | 1.454     | 0.146   | 0.894        | 2.134        |
| Age_Interval[80,90)  | 1.332      | 0.224          | 1.277     | 0.202   | 0.858        | 2.068        |
| Age_Interval[90,100] | 1.358      | 0.192          | 1.596     | 0.110   | 0.933        | 1.979        |
| Cases                | 1.000      | 0.000          | -25.022   | 0.000   | 1.000        | 1.000        |

Figure 11: Odds Ratio table with Confidence Intervals

## References

- Alexander, Rohan. 2023. *Telling Stories with Data*. Chapman; Hall/CRC. <https://tellingstorieswithdata.com/>.
- Kuhn, Max. 2023. *Caret: Classification and Regression Training*. <https://CRAN.R-project.org/package=caret>.
- Open Science Framework (OSF). 2022. “COVID-19 Cases, Deaths, and Tests by Age, Gender, and Region.” Open Science Framework (OSF). <https://osf.io/43ucn>.
- R Core Team. 2023. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.
- Robin, Xavier, Natacha Turck, Alexandre Hainard, Natalia Tiberti, Frédérique Lisacek, Jean-Charles Sanchez, and Markus Müller. 2011. “pROC: An Open-Source Package for r and s+ to Analyze and Compare ROC Curves.” *BMC Bioinformatics* 12: 77. <https://doi.org/10.1186/1471-2105-12-77>.
- Sharma, Rajesh, Anjali Singh, Poonam Verma, and Suresh Kumar. 2023. “Analyzing Uttarakhand’s COVID-19 Outbreak: Demographic Insights and Strategies for Future Pandemic Prevention.” *Current Medical Research and Opinion* 39 (1): 45–52.
- Wickham, Hadley, Mara Averick, Jennifer Bryan, Winston Chang, Lucy D’Agostino McGowan, Romain François, Garrett Grolemund, et al. 2023. *Tidyverse: Easily Install and Load the ‘Tidyverse’*. <https://CRAN.R-project.org/package=tidyverse>.
- Wickham, Hadley, Romain François, Lionel Henry, and Kirill Müller. 2023. *Dplyr: A Grammar of Data Manipulation*. <https://CRAN.R-project.org/package=dplyr>.
- Wickham, Hadley, Lionel Henry, and Davis Vaughan. 2023. *Tidyr: Tidy Messy Data*. <https://CRAN.R-project.org/package=tidyr>.
- Yang, Ziang, Xieraili Tiemuerniyazi, Fei Xu, Yang Wang, Yang Sun, Peng Yan, Liangxin Tian, et al. 2024. “Partial Cardiac Denervation to Prevent Postoperative Atrial Fibrillation After Coronary Artery Bypass Grafting: The pCAD-POAF Randomized Clinical Trial.” *Progress in Cardiovascular Diseases* 80: 1–9.