

Predicting Infection Outbreak Dynamics in Toronto Healthcare Institutions: A Bayesian Approach

Respiratory Outbreaks Are More Likely in Long-Term Care Homes and Retirement Homes,
Driven by Seasonal Trends and Causative Agents Like COVID-19 and Influenza

Divya Gupta
December 3, 2024

Abstract.....	1
1. Introduction.....	2
2. Data.....	2
Measurement.....	3
Constructed Variables.....	4
Accuracy and Limitations.....	4
Understanding the Data.....	4
Outcome Variable.....	5
Predictor Variables.....	5
Visualising the Data.....	5
3. Model.....	9
3.1 Model Set-Up.....	10
3.1 Model Set-Up.....	10
3.3 Assumptions and Limitations.....	11
3.4 Model Validation.....	11
4. Results.....	12
5. Discussion.....	14
5.1 Strengths.....	14
5.2 Limitations.....	14
References.....	15

Abstract

Identifying key trends in infection outbreaks in a city can help policymakers implement better reforms to support citizens and prevent outbreaks. This study uses cleaned datasets retrieved from Open Data Toronto to develop a model that can have insightful implications for the City of Toronto policymakers, especially regarding resource allocation for outbreak prevention. This study aims to help policymakers better understand the trends of infection outbreaks in Toronto,

so that this data can be used to identify problem areas in the city (such as specific times of year). This study examines the factors influencing healthcare outbreaks in Toronto healthcare institutions, such as: type of infection (respiratory vs. gastroenteric), causative agent of infection, setting and month of infection. This study uses Bayesian logistic regression to identify the factors associated with infectious outbreaks and predict the likelihood of a specific outbreak type depending on the aforementioned variables. The findings show that COVID-19 is a significant causative agent for respiratory infections, while norovirus is a causative agent for enteric outbreaks. Moreover, there is a significant increase in the number of outbreaks in winter months and in long term care homes.

1. Introduction

Toronto has multitudes of healthcare institutions such as hospitals, long-term care homes and retirement homes that often house hundreds of people at a time and can serve as places of community for vulnerable population groups. Any healthcare-associated infectious outbreaks in such spaces pose substantial risk for all vulnerable populations in the space. Understanding the varying factors that influence and further an outbreak - such as its causes, duration, and type of outbreak - can help better understand the nature of the outbreak and plan interventions and outbreak prevention measures accordingly. Furthermore, analyzing data about infection outbreak trends in Toronto can help us predict the time, region and duration of the next outbreak occurrence which can be instrumental in helping vulnerable populations such as older people and children. This study aims to answer the research question: *How can Bayesian models help predict infection outbreak behavior, specifically based on the outbreak's healthcare setting, causative agent and month of outbreak?*

The estimand for this analysis is the probability of an infection outbreak being of a specific type (namely respiratory or enteric) as determined by the causative agent, month of outbreak, or healthcare setting.

Developing a robust Bayesian model that can predict outbreaks in the city and develop outbreak management strategies in real time, ensuring better resource allocation and healthcare management. I will be employing Bayesian Logistic Regression to predict the likelihood of respiratory versus enteric outbreaks and how they are influenced based on the healthcare institution, causative agent, duration of previous outbreaks, and time of the year. I will then be employing Bayesian Survival Analysis to model the duration of outbreaks and identify defining factors specifically in prolonged cases (duration of infection > 10 days).

The remainder of the paper is structured as follows. Section 2 details the data, measurement, along with helpful visualisations to understand the data. Section 3 outlines the modeling approach, while Section 4 presents the results. Section 5 discusses implications, limitations, and future directions. This is followed by the Appendix.

2. Data

This dataset, entitled “Outbreaks in Toronto Healthcare Institutions” was obtained from Open Data Toronto on November 25, 2024 (Open Data Toronto, 2024) and possesses data from January 2nd, 2024 until November 20, 2024.. It is published by Toronto Public Health and documents infection outbreaks reported by various healthcare institutions in the city (the dataset is updated weekly on Thursday to capture the latest information). This is done in compliance with Ontario’s Health Protection and Promotion Act, which requires healthcare facilities to continuously report all infection outbreaks to their local authorities. This level of data collection allows timely identification and potential containment and reform in affected communities.

For the purpose of this study, an infection outbreak (here on referred to as “outbreak”) is defined as a localized increase in the rate of infection above the baseline. In terms of settings, this dataset gathers information from hospitals (acute, psychiatric or chronic care), long term care homes (here on referred to as “LTCH”), retirement homes, and transitional care systems. Further, the dataset includes information about the type of outbreak (for eg., respiratory or enteric), location, the start and end date of the outbreak, and the causative agent(s). While Open Data Toronto possesses multiple other datasets regarding infections and outbreaks in Toronto, they were localised to information about homeless shelters or solely tracked COVID-19 outbreaks in the city. The chosen dataset encompasses a wider range of outbreak settings and includes various causative agents, including COVID-19, positioning it as a wider dataset - allowing us to extract greater insights and holistically understand the nature of infection outbreaks in Toronto.

For data cleaning, analysis and visualisation, this study uses R (R Core Team 2023), along with packages such as tidyverse for visualisation (Wickham et al. 2019), dplyr for data wrangling, and lubridate for handling temporal data. Further, ggplot2 was used for creating visualizations (Wickham, Chang, et al. 2023) and forcats for working with categorical variables. For the Bayesian analysis, rstanarm was and brms were used. Readr was used to read the structured data and scales allowed further development upon the visualisations. We also used ggeffects for visualisations of the model predictions and testthat for unit testing.

Measurement

The measurement process aimed to transition real-world phenomena such as infection outbreaks in healthcare facilities into usable, numerical or categorical data that can be used for analysis. Each data entry in the data set refers to a case of infection outbreak in Toronto from January - November 2024. Variables of interest include:

Type of Outbreak: This refers to the nature of the outbreak, categorized as either respiratory or enteric - based on the symptoms of the infection. Respiratory infections are categorized by symptoms such as cough, fever, and sore throat. Enteric infections are categorized by symptoms such as diarrhea, nausea and vomiting.

Outbreak Setting: This includes variables such as hospitals (acute, psychiatric or chronic care), long-term care homes (LTCHs), retirement homes, or transitional care facilities.

Causative Agent: This refers to the variable, “Causative Agent - 1” in the dataset. This records the primary pathogen named to be responsible for the outbreak. Significant data entries included COVID-19, norovirus, and Influenza as the causative agents. This categorical variable can help one understand the nature of the infection and its instigating factors.

Date Outbreak Began: Records the first date, reported as YYYY-MM-DD, that the outbreak was reported.

Date Declared Over: Records, as a YYYY-MM-DD date, when all affected patients were declared infection-free, marking the end of the outbreak.

Constructed Variables

We also constructed some variables for the ease of data analysis from the provided variables from the raw dataset. These include:

Duration: This was calculated as the difference between Date Declared Over and Date Outbreak Began to give a numerical value for the number of days an outbreak lasted.

Month: This was constructed by monitoring the start date of an outbreak, from Date Outbreak Began, to track how seasonal changes may affect outbreak levels in the city.

Accuracy and Limitations

The Open Data Toronto portal provided the raw dataset that had some inconsistencies, namely - various fields had missing entries in the “Date Declared Over” tab, which rendered those data points useless as we could not track the duration of the outbreak. This value is essential for this analysis as our variable, Duration, was constructed as the difference between “Date Declared Over” and “Date Outbreak Began.” Further, some data points had no specific entry for “Type of Outbreak”, and had values like “other” which did not provide usable information for this analysis. For ease of analysis, these data entries were removed while constructing the analysis dataset.

Understanding the Data

Outcome Variable

Our model predicts the probability of an outbreak being labeled as respiratory or enteric based on its predictors. To that effect, our outcome variable is “Type_of_Outbreak”, with its results being in binary form of either respiratory or enteric.

Predictor Variables

Our analysis mainly focused on three predictor variables to help determine our outcome variable:

1. **Outbreak Setting:** This categorical variable utilises information about which type of healthcare institution housed an outbreak, such as hospitals or LTCH.
2. **Causative Agent:** This categorical variable houses information about the primary cause of the outbreak, such as COVID-19 or Influenza.
3. **Month:** This categorical variable represents the month of the year in which the outbreak began.

Visualising the Data

Figure 1 illustrates the difference between the levels of respiratory and enteric outbreaks in the city. As can be seen, the overwhelming amount of outbreaks reported are respiratory.

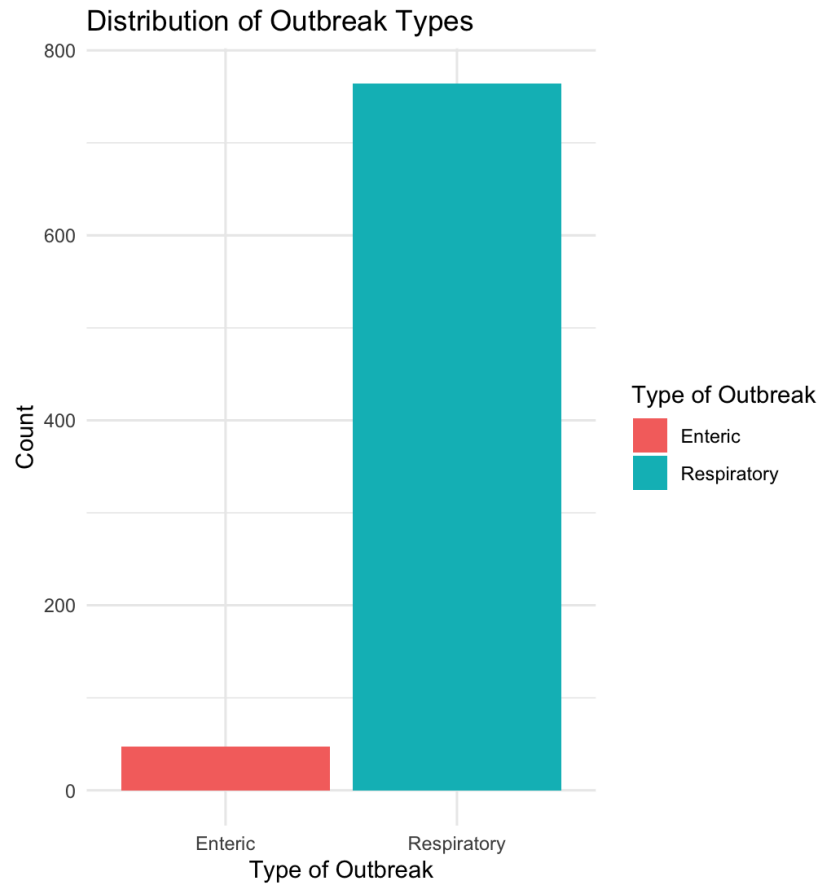


Figure 1. Distribution of Outbreak Types

Figure 2 takes this one step further and illustrates the number of outbreaks determined by healthcare settings. This shows that LTHC had an overwhelming majority in the number of outbreaks reported, followed by retirement homes.

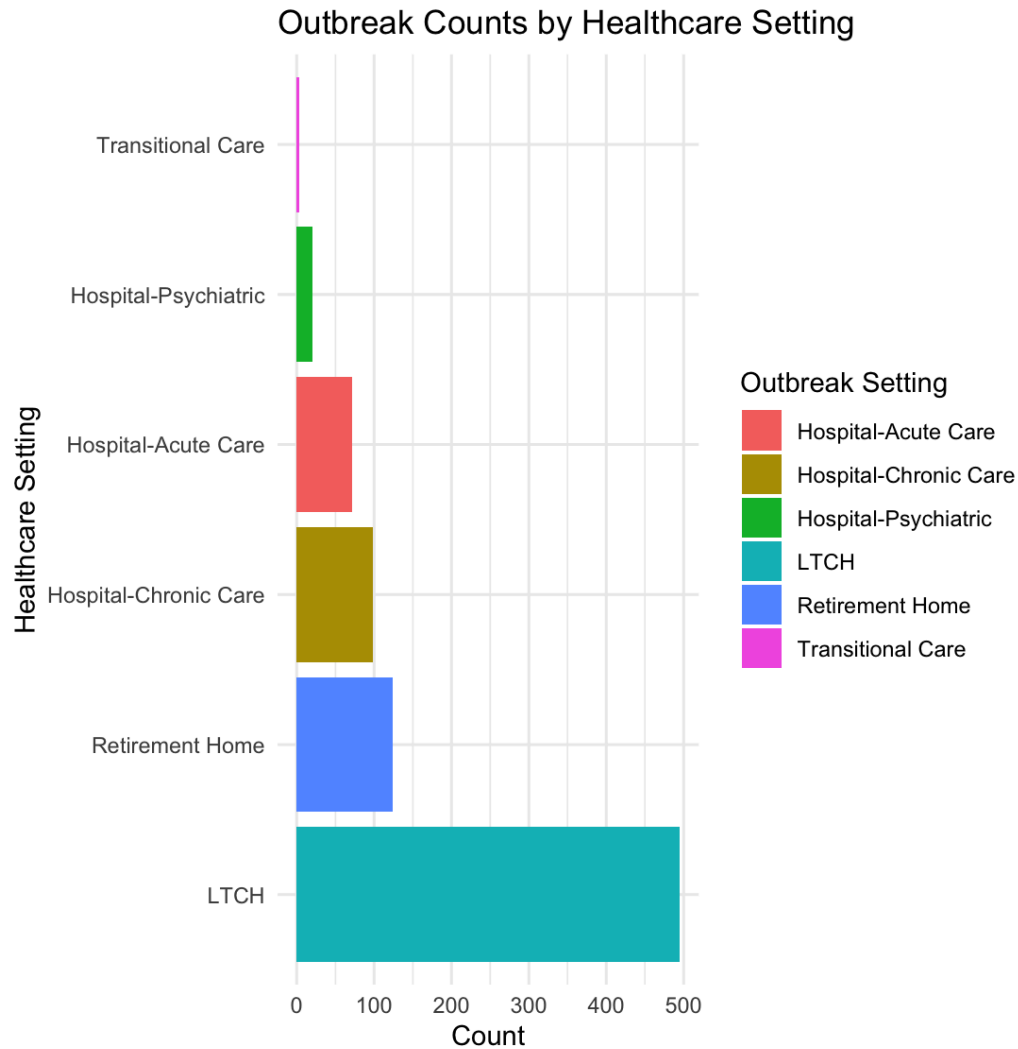


Figure 2. Outbreak Counts by Healthcare Setting

Next, we aimed to understand the interaction between different types of infection outbreaks and the setting in which they occurred. This is illustrated in Figure 3 which is a stacked bar chart. Key insights included the fact that respiratory infections dominated across all types of healthcare settings. While enteric outbreaks are less common, they are most prevalent in retirement homes and hospitals-acute care. Psychiatric care units in hospitals and transitional care systems reported no cases of enteric outbreaks in this dataset.

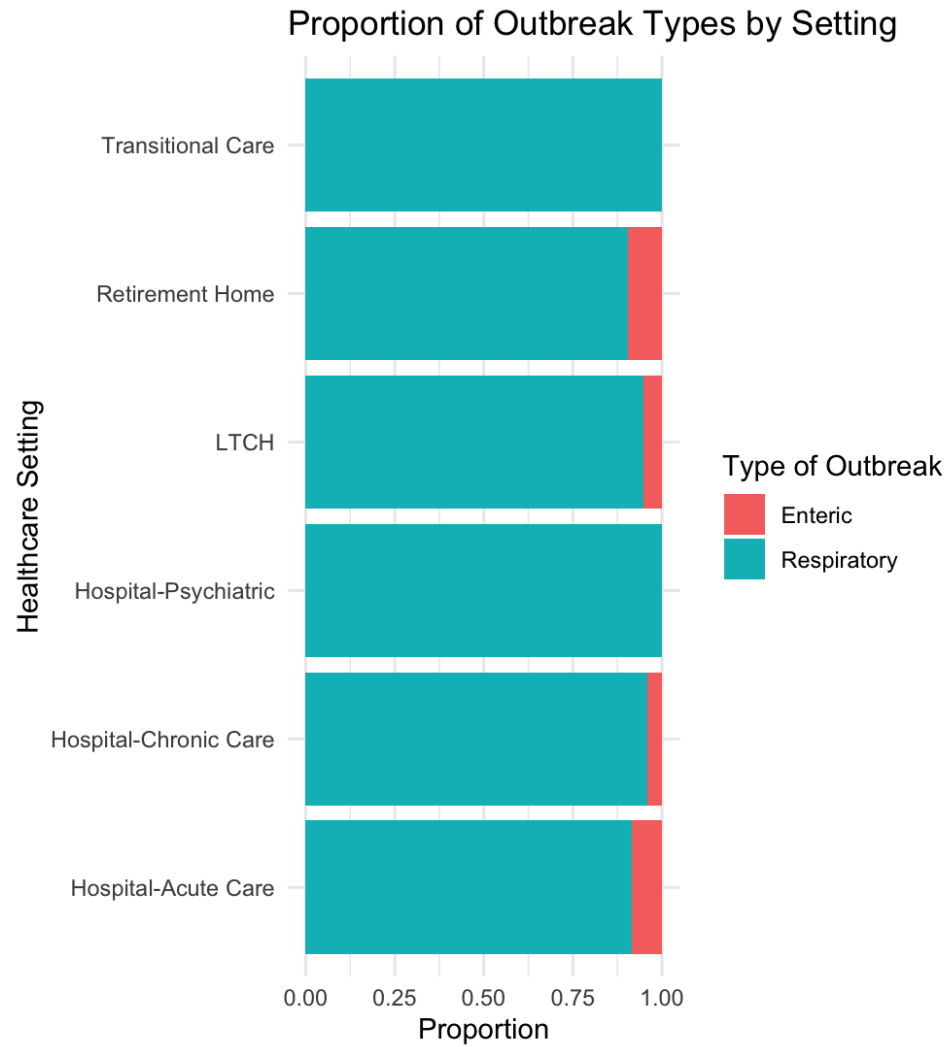


Figure 3. Proportion of Outbreak Types by Setting

Lastly, we aimed to understand the overall interaction trends between the outbreak duration, the causative agent, and the setting of the outbreak. This is illustrated in Figure 4.

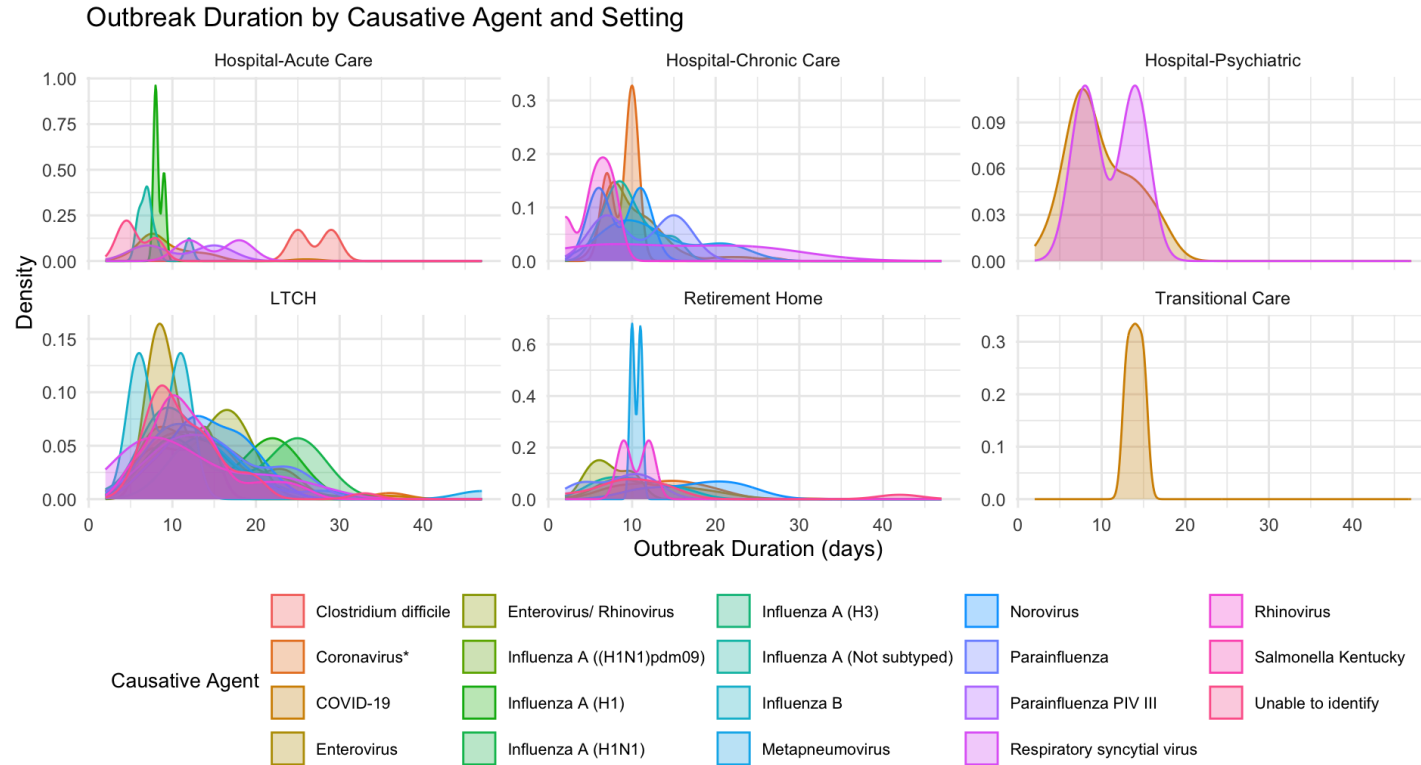


Figure 4. Outbreak Duration by Causative Agent and Setting

This series of density plots illustrates the duration of outbreaks as stratified by the causative agents and various healthcare settings included in the dataset. This shows that LTCH and Retirement homes report significantly longer outbreak durations, specifically for causative agents such as norovirus and Influenza. Further, transitional care settings have a short outbreak duration. This may be due to strict infection control practices - this raises a question about outbreak control policies in less-monitored settings such as retirement homes which have reported longer outbreak durations. These differences in durations according to healthcare settings also raises a need for setting-tailored outbreak interventions.

Overall, these visualisations help understand key parts of the data that will inform our model. Namely, they illustrate that an overwhelming majority of infections are respiratory in nature. Next, there is a difference in the duration of an outbreak according to the healthcare setting, and that there are potential impacts of the causative agent and the month in which the outbreak occurred upon the duration of the outbreak.

3. Model

3.1 Model Set-Up

3.1 Model Set-Up

This study aims to build a Bayesian Logistic Regression model using R (R Core Team, 2023) to predict the probability of a respiratory outbreaks in Toronto. Based on our exploratory data analysis, we can see that predictor variables such as causative agent, healthcare setting, and month/duration of the outbreak are key predictors in helping model this. The model is, thus, as follows:

$$Y_i \sim \text{Bernoulli}(\pi_i)$$

$$\text{logit}(\pi_i) = \beta_0 + \beta_1 \cdot \text{Outbreak_Setting}_i + \beta_2 \cdot \text{Causative_Agent_1}_i + \beta_3 \cdot \text{Month}_i + \epsilon_i$$

Where:

1. π_i is a predicted probability of a respiratory outbreak for the i th observation.
2. β_0 is the intercept, that captures the baseline log-odds of a respiratory outbreak.
3. Outbreak_Setting helps determine the role of healthcare settings in the determination of an outbreak.
4. Causative_agent captures the role of the primary pathogen that causes the outbreak.
5. Month captures the month in which an outbreak started, aiming to capture any seasonal aspects of outbreaks.
6. Lastly, ϵ_i is an error term capturing any noise in the model.

Priors

This model incorporates the following priors:

- $\beta \sim N(0, 2.5)$ for the coefficients of predictors.
- $\beta_0 \sim N(0, 5)$ for the intercept.

These are weakly informative priors that help constrain extreme values of the variables without being extremely restrictive.

3.2 Model Justification

Since this paper deals with a binary variable such as Respiratory infection outbreaks, logistic Bernoulli regression was used. This is so that prior knowledge and uncertainty (through the priors and error term) can be accommodated and binary outcome variables can be handled effectively.

The logit link function in the model ensures that the predicted probability of a respiratory outbreak is between 0 and 1. A Bayesian approach was taken to the model to allow full posterior distributions for each predictive parameter. Further, it allows for the inclusion of prior knowledge, which can be helpful in understanding outbreak trends over a long period of time. Further, the priors help regularize parameter estimates and prevent overfitting to certain parameters such as outbreak settings or causative agents.

The predictors were chosen for this model based on our exploratory data analysis, where they were shown to be significant in understanding the type of outbreak. Certain alternative model choices such as regular logistic regression were considered but not used due to their inability to incorporate prior knowledge. The Bayesian model, overall, includes prior knowledge along with a nuance to the convoluted nature of our chosen predictor variable, making it a robust choice for this analysis.

3.3 Assumptions and Limitations

There are several key assumptions of this model. It assumes that the outcome variable, Type of Outbreak, is a binary outcome (respiratory or enteric). Further, it assumes that each outbreak is an independent event influenced by its own set of predictive variables. However, this assumption may not hold for outbreaks that occurred in the same healthcare setting, or occurred very closely together, or overlapped in the duration or month variables. This model does not account for these potential interactions. Lastly, the weakly informative priors are assumed to appropriately reflect the healthcare setting due to their unrestrictive flexibility.

This raises certain limitations to the scope of the model. The model assumes the dataset it was developed on was accurately and consistently reported on by all healthcare institutions. It is possible that the causative agent or the duration of the outbreak, for example, could be misreported which would skew the results. Further, nuanced aspects of healthcare settings such as quality of care provided, healthcare policies, severity of symptoms, and staff-to-patient ratios can vary the severity and occurrence of the outbreak but these were not captured by this model. While the month variable was included to capture seasonal trends in outbreaks, it is possible that this may vary year to year. Our dataset only captures information for 2024 and may not be generalizable to a wider range of time without proper adjustment.

3.4 Model Validation

Validation is critical to ensure the models are reliable and provide accurate predictions. As a part of the model validation process, Posterior Predictive Checks were conducted to assess the fit of the model. Figure 5 shows the Posterior Predictive Check that compares the observed data (illustrated by y) and the predicted data (illustrated by y_{rep}). The observed outcome is overlaid

with the densities of the simulated outcome to see how well the model's prediction matches with the real data.

There is a close alignment between the two types of data, which shows that the model captures the data distribution effectively. There are slight deviations in the lower probability regions which suggests that there are areas for improvement in the model. This could be tackled in the future by including more predictor variables.

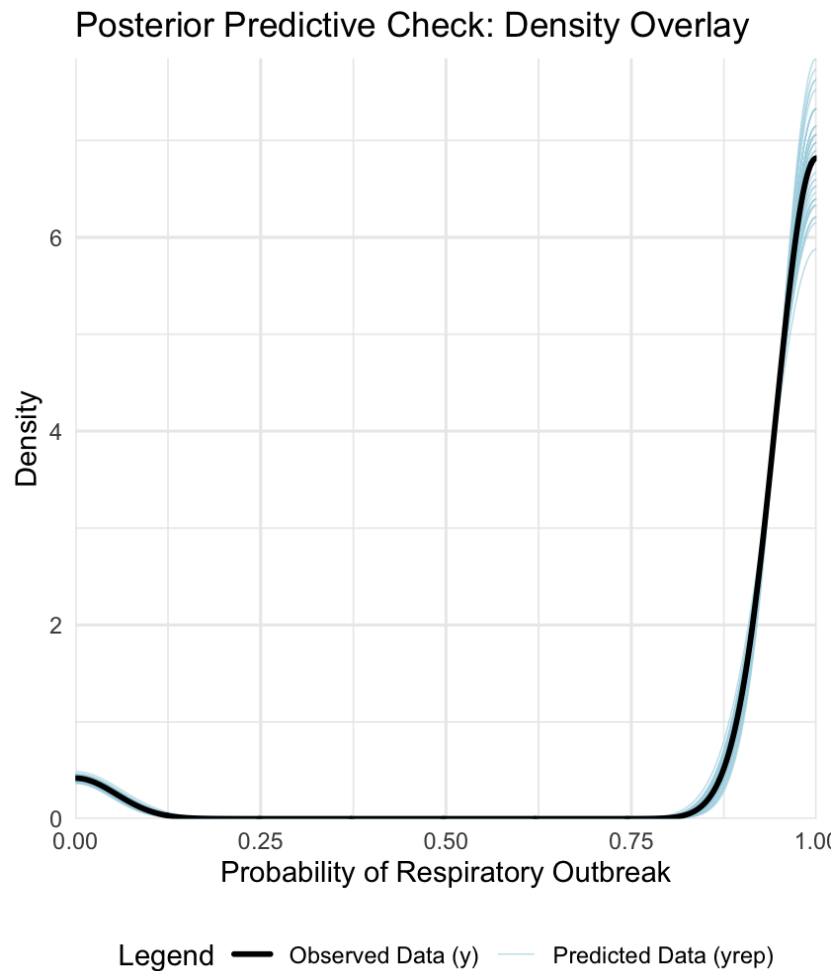


Figure 5. Posterior Predictive Check with Density Overlay

4. Results

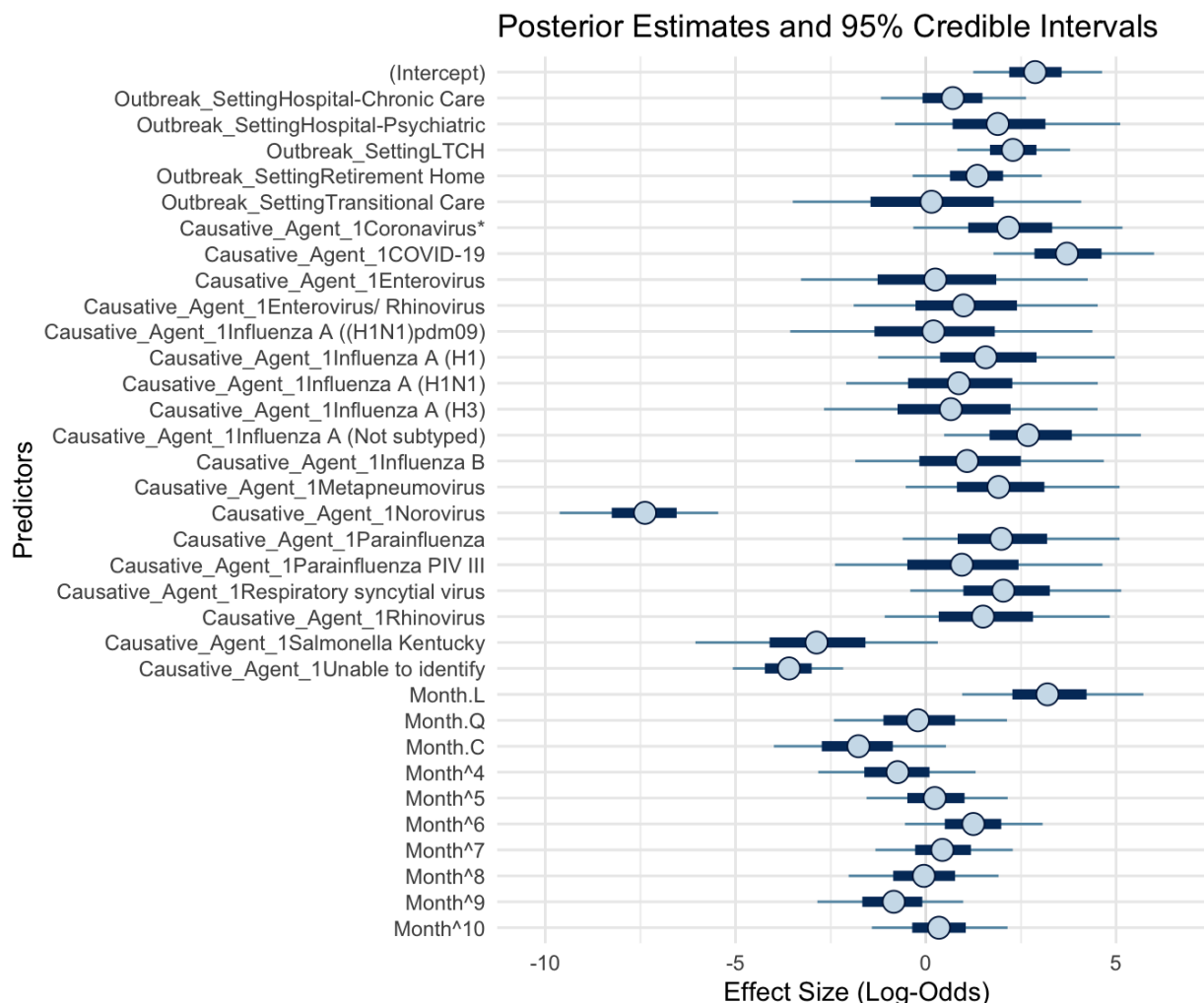


Figure 6. Posterior Estimates and 95% Credible Intervals

Figure 6 displays the posterior estimates and their corresponding 95% credible intervals for each predictor in the model. The x-axis represents the effect size as log-odds while the y-axis shows the predictors. This figure shows that LTCH have a high number of respiratory outbreaks, given that their credible intervals are above zero. Hospital-Chronic Care, Hospital-Psychiatric, and Retirement Homes also have positive effects but have differing magnitudes and credible intervals. Further, causative agents such as *COVID-19*, *Coronavirus*, and *Respiratory Syncytial Virus* have positive associations with respiratory outbreaks as their credible intervals are above zero. *Influenza A* and *Metapneumovirus* also have an association with respiratory outbreaks but have wider intervals, showing uncertainty. On the other hand, *norovirus* is heavily associated with enteric outbreaks as they have a negative association with respiratory outbreaks. There is a

wide variability in the month and its association with different outbreaks, illustrating uncertainty in this predictor.

Overall, we see that LTCH settings ($\beta=2.3$) are most associated with respiratory outbreaks, followed by Retirement homes and Hospital- Chronic care. Further, COVID-19, Coronavirus, and Respiratory Syncytial Virus are the strongest predictors of respiratory outbreaks for causative agents. Causative agents like COVID-19 ($\beta=3.8$) strongly increase the likelihood of respiratory outbreaks while norovirus ($\beta=-7.4$) strongly increase the likelihood of enteric outbreaks.

5. Discussion

Overall, this study aimed to analyse the trends and nature of infection outbreaks in Toronto as determined by various predictors such as setting, causative agents, and month of outbreak occurrence. This data was collected over the year of 2024 in healthcare settings and was reported by Toronto Public Health, available on Open Data Toronto. By employing a Bayesian logistic regression model, this study quantified the association between different outbreak types (respiratory or enteric) and their predictor variables. This information can be used to predict the occurrence of an outbreak based on these variables. These outcomes help determine trouble zones in the city, such as LTCH, and can help policymakers determine key areas for further reform.

5.1 Strengths

The inclusion of an array of predictors allowed the model to account for a high level of variability in the data, and understand the effect of different factors on outbreaks. The Bayesian approach allowed the quantification of categorical variables as well as the inclusion of weakly informative priors that allow the incorporation of prior knowledge from healthcare settings. This allowed for key insights, as discussed in the Results section, that showed that LTCH and Retirement Homes are key healthcare settings to be targeted to tackle respiratory outbreaks.

5.2 Limitations

This study has several limitations. The dataset relies on accurate reporting by healthcare institutions, which tend to be overworked and can misreport data. Further, multiple data entries were excluded due to incomplete data points, which skews the robustness of the model. The dataset only incorporated information from 2024 and it is unclear if this can be extrapolated to different years. A wider dataset, with more information ranging a 5 or 10 year period, may allow for a more robust model. While some seasonal trends were seen in the analysis, it is still important to note that these are seasonal trends over the span of only one year, and they can vary significantly. More precise week-by-week analysis may be favored over monthly analysis to

understand outbreak trends. Further, our outcome variable is binary - respiratory or enteric - and can overlook the aspect of interaction between various types of infections occurring together.

References

R Core Team (2021). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. <https://www.R-project.org/>.

Toronto Public Health. (2024b, November). Outbreaks in Toronto Healthcare Institutions. Toronto; Open Data Toronto.

<https://open.toronto.ca/dataset/outbreaks-in-toronto-healthcare-institutions/>