

Interrupted Time Series and Machine Learning with Application to Effect of Influenza Vaccine

By

Cynthia Ombok Juma

144892

Master of Science in Data Science and Analytics

2024

Interrupted Time Series and Machine Learning with Application to Effect of Influenza Vaccine

By
Cynthia Ombok Juma
144892

**Submitted in Partial Fulfillment of the Requirements for the Degree of Master of Science in Data
Science and Analytics at Strathmore University**

Institute of Mathematical Sciences
Strathmore University
Nairobi, Kenya

April 2024

**This thesis is available for Library use on the understanding that it is copyright material and that
no quotation from the thesis may be published without proper acknowledgement.**

Declaration and Approval

Declaration

I declare that this work has yet to be submitted and approved for the award of a degree by this university or any other university. The dissertation contains no material previously published or written by another person except where due reference is made in the dissertation itself.

© No part of this thesis may be reproduced without the permission of the author and Strathmore University


Student's Name: **Cynthia Ombok Juma**

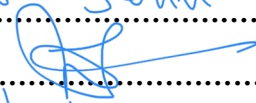
Signature: 

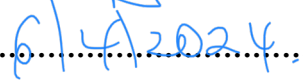
Date: **06 April 2024**

Approval

This dissertation by Cynthia Ombok Juma was reviewed and approved by:

Lecturer: 

Signature: 

Date: 

Dedication

To everyone who helped me to complete this work, I dedicate it. They include my family, who have always supported and encouraged me, and my peers and professors, who have provided thoughtful feedback throughout this project. The most significant thing is God, without whom I could not have advanced as far as I have.

Acknowledgment

I want to thank Dr. Olukuru, my supervisor, for his help with the research. I appreciate the opportunity and exposure that the Institute of Mathematical Sciences at Strathmore University and @ilabAfrica offered me throughout my study period.

Abstract

Interrupted time series analysis is being increasingly employed to assess the effects of extensive health interventions. Autocorrelation and seasonality are best captured but are not well captured by the simple implementation of the time series model like segmented regression, which is commonly used. An Autoregressive Integrated Moving Average (ARIMA) model presents an alternative approach to address these issues. In this study, the fundamental principles of ARIMA and LSTM models are expounded upon, along with their application in evaluating interventions at a population level, such as determining the effect of influenza vaccine administration. Considerations such as determining the impact shape, model selection process, transfer functions, loss functions, selection of batch sizes and epochs training the neural networks, evaluation metrics, and interpreting results are discussed. Additionally, detailed R and Python codes are provided for result replication. The application of ARIMA and LSTM predictive modeling is demonstrated through an analysis of influenza vaccination intervention to reduce the number of medically attended respiratory illnesses among children under five years. Precisely, from November 2019 to November 2021, an influenza vaccination demonstration project. In conclusion, ARIMA modeling and LSTM serve as valuable tools for assessing the effects of large-scale interventions when traditional methods are not applicable, given their ability to consider underlying patterns, autocorrelation, seasonality, and flexibility in modeling various impacts. Comparing the MAE and RMSE error results, LSTM outperformed the ARIMA model. Key terms: Interrupted time series analysis, Autoregressive integrated moving average models, LSTM, Intervention analysis

Table of Contents

Declaration and Approval.....	i
Declaration.....	i
Dedication	ii
Acknowledgment	iii
Abstract.....	iv
Table of Contents	v
List of Figures.....	viii
List of Tables	ix
List of Abbreviations	x
Chapter 1: Introduction.....	1
1.0 Background of the study	1
1.1.1 Interrupted Time Series.....	2
1.1.2 Machine Learning.....	4
1.1.2.1 Long Short-Term Memory (LSTM) Recurrent Neural Network Model.....	4
1.2 Statement of the Problem	6
1.3 Study Objectives.....	7
1.3.2 Specific objectives.....	7
1.3.3 Research Questions	7
1.4 Scope of the study	8
1.5 Relevance of the study	8
Chapter 2: Literature Review	9
2.0 Introduction	9
2.1.0 Evaluating Vaccine Effectiveness, Efficacy, and Impact.....	9
2.1.2 Effectiveness of Influenza Vaccination.....	10
2.1.3 Analyzing the effects of vaccination programs.....	11
2.2.0 Predictive Modeling for Respiratory Illness.....	12
2.2.1 Interrupted Time Series Analysis in Public Health	13
2.2.2 Machine Learning Techniques in Healthcare Forecasting	15
2.2.3 Evaluating Model Precision in Interrupted Time Series Analysis and Deep Learning Models	16
2.3 Conclusion	18

2.5 Research Gap	18
Chapter 3: Methodology	19
3.0 Study Design	19
3.1. Data Sources and Collection Methods	19
3.2 Data Preprocessing	20
3.2.1 Data Cleaning	20
3.2.2 Data Exploration	21
3.2.3 Handling Missing Mata	21
3.2.5 Outlier Detection and Handling.....	22
3.2.6 Data Type Conversion	22
3.3 Data Transformation	22
3.3.1 Creation of Categorical Variables	22
3.4. Preprocessing Steps for ITS Analysis.....	22
3.4.1 ARIMA Model Specification.....	23
3.4.2 Determining Stationarity.....	24
3.4.3 Autocorrelation	26
3.4.3 Seasonality	26
3.4.4 Transfer Functions.....	27
3.5 Building the ARIMA Model.....	28
3.5.1 Plotting Data to Understand Underlying Patterns.....	28
3.5.2 Choosing Model Parameters (p, d, and q)	28
3.5.2 Determining the Seasonal Decomposition of the Time Series.....	29
3.5.3 Plotting the Time Series Data with Intervention.....	29
3.6.4 Fitting the ARIMA Model.....	29
3.7.5 Forecasting the ARIMA Model	29
3.6 Data Preparation for LSTM-RNN Modeling	29
3.6.1 Data preprocessing.....	29
3.6.2 Feature scaling	29
3.7 LSTM-RNN Specification	30
3.7.1. Architecture and Model Setup.....	30
3.7.2 Model Compilation:	32
3.7.2.1 Training the NN	32
3.7.2.2 Layers, Neurons and Hyperparameters.....	32
3.7.2.3 Loss Function.....	32
3.7.2.4 Training Procedure.....	32

3.7.2.5 Batch Sizes and Number of Epochs.....	33
3.8.4 Evaluation Metrics.....	33
3.9 Comparison of Traditional Time Series Analysis and LSTM -RNN algorithm.....	34
3.10 Ethical considerations	34
3.11 Limitations.....	34
Chapter 4: Results.....	35
4.0 Introduction.....	35
4.2 Time Series Plot of The Cases and Controls.....	36
4.3 Building an ARIMA Model.....	38
4.3.1 Determining Stationarity.....	38
4.3.2 Autocorrelation Plots.....	39
4.3.4 Fitting an Automated ARIMA Model.....	40
4.3.5 Residual check.....	41
4.3.6 Final ARIMA Model.....	41
Chapter 5: Discussion	44
Chapter 6: Conclusion and Recommendation.....	46
References	47
Appendices.....	62

List of Figures

Figure 1 Uses cases of data science in health care.	2
Figure 2 Use of impact models used in ITS. Source : (Barone-Adesi et al., 2011)	3
Figure 3 Example of a predictive deep learning model.	5
Figure 4 Flow chart for ARIMA model selection. Source:(8.7 ARIMA Modelling in R Forecasting, n.d.)	23
Figure 5 Description of transfer functions within interrupted time series analysis in the context of ARIMA. Source:(Schaffer et al., 2021)	28
Figure 6 Sigmoid functions for LSTM gates.....	30
Figure 7 LSTM memory cell.	31
Figure 8 Retrospective Time Series Plots of Cases Data from 2019 to 2022.....	36
Figure 9 Retrospective Time Series Plots of Control Data from 2019 to 2022.....	37
Figure 10 Time series plot illustrating the intervention period spanning from November 2019 to November 2021.	38
Figure 11 AFC and PACF plots for the cases and control time series before differencing.....	39
Figure 12 AFC and PACF plots for the cases and control time series post differencing.	40
Figure 13 Residual check for the final model, ARIMA (01,2) (0,0,1) ₁₂	41
Figure 14 Observed Values and Predicted Values in Absence of Intervention based on ARIMA Model.	42
Figure 15 The observed number of medically attended respiratory illness and values predicted by LSTM from 2019 to 2022	43

List of Tables

Table 1 Medically attended respiratory illness indicators in the DHIS. 19

Table 2 Handling missing values..... 21

Table 3. Descriptive summaries of the number of medically attended respiratory illness. 35

Table 4 Augmented Dickey-Fuller Test 38

List of Abbreviations

ACF - Autocorrelation Function

AR - Autoregressive

ARIMA Model - Automated Regressive Integrated Moving Average Model

LSTM - Long Short-Term Memory

RNN - Recurrent Neural Network

ML - Machine Learning

ITS - Interrupted Time Series

ITSA - Interrupted Time Series Analysis

MA - Moving Average

PACF - Partial Autocorrelation Function

RCTSs - Randomized Controlled Trials

CDC - Centre for Disease Control

WHO - World Health Organization

ACEs - Acute Coronary Events

DHIS - District Health Information System

Chapter 1: Introduction

1.0 Background of the study

Influenza (flu) is a highly contagious respiratory illness caused by influenza viruses that infect the respiratory system, including the nose, throat, and lungs. Children, adults, and people with chronic medical conditions are at higher risk of severe flu complications. It is primarily characterized by sudden fever, cough, sore throat, headache, chills, runny or stuffy nose, body aches, and fatigue (CDC, 2022). It can cause mild to severe illness and sometimes lead to death. Influenza is a leading cause of morbidity and mortality worldwide in children under five years. Children under the age of 2 years are at higher risk of developing severe flu-related complications (CDC, 2020). There are four main types of influenza (flu) viruses: Types A, B, C, and D. However, the influenza A and B viruses that routinely spread in people (human influenza viruses) are responsible for yearly seasonal flu epidemics. Currently, the best way to reduce the risk of flu infections and their potentially serious complications is by getting vaccinated each year, as the vaccines are safe and reliable and have been used for over 60 years, as the World Health Organization (WHO) outlined. Additionally, WHO recommends annual vaccination for the following group of persons: children aged between 6 months to 5 years; pregnant women at any stage of pregnancy; elderly individuals (aged more than 65 years); individuals with chronic medical conditions and healthcare workers (*Influenza (Seasonal)*, n.d.).

Annually, the Centre for Disease Control (CDC) conducts studies to determine how well the influenza vaccine contributes to the protection against flu. Recent studies reveal that the influenza vaccine reduces the risks of flu illness by 40.0% to 60.0% among all populations, especially if the flu in circulation is well-matched with the vaccine administered (*Vaccine Effectiveness*, 2022). Between 2019 and 2020, the last flu season of the pre-COVID-19 pandemic period, it was estimated that 7.5 million influenza illnesses, 3.7 million influenza-associated medical visits, 105,000 influenza-associated hospitalizations, and 6,300 influenza-associated deaths were prevented by flu vaccination in the United States of America. A 2022 study to ascertain cause-specific mortality for neonates and children younger than five years from 2000 to 2019 found that lower respiratory illness contributed to 13.95% of the reported child deaths. At the same time, vaccine-preventable deaths, such as meningitis, measles, and lower respiratory illness, among others, were at 21.7% (Perin et al., 2022).

Global public health threats have been rapidly emerging, with the most recent notable cases being COVID-19 (*Coronavirus*, n.d.) and Monkey Pox (*Mpox (Monkeypox)*, n.d.), in turn, the use of data science methodologies has risen, especially with the increase in the amount of big data generated, fields like medicine and public health have leveraged on these techniques to draw data-driven insights (Subrahmanya et al., 2022). An evidence-based approach to better health care has been facilitated by

statistical modeling in epidemiology and public health. Biostatistics tools extensively provide an understating of diseases and allow for developing new strategies for control and prevention. Epidemiology data analysis has proven helpful in decision-making, policy changes, and the management of diseases (Matranga et al., 2021). In this study, we will leverage both methodologies to predict the effect of influenza vaccination.



Figure 1 Uses cases of data science in health care. ¹

1.1.1 Interrupted Time Series

An interrupted time series (ITS) design involves consistent data collection before and after an interruption. It entails introducing and withdrawing an intervention or some part of it and observing any changes in the outcome under assessment within a study population. ITS analysis utilizes statistical methods to quantify changes in level (slope) and trend before and after intervention and assess the significance of the observed differences. The primary assumption is that observations from the baseline period predict where the future data points lie in the absence of the introduced intervention. The effect

¹ Source: <https://www.bacancytechnology.com/blog/data-science-use-cases>

size is expressed in terms of the level of change and slope change. The level change does not provide all the information in the data, while the slope change post-intervention period shows the level of impact (Fretheim & Tomic, 2015). ITS design has increasingly been recommended as a more robust design where randomization is impossible.

A 2013 study used ITS to evaluate healthcare quality improvements and reported that it was a simple but powerful approach for program evaluation. Additionally, this study highlighted that it was valuable and less expensive than comparable randomized control trials (RCTs) conducted to answer similar questions (Penfold & Zhang, 2013). A similar 2017 study found that it allowed for a more detailed assessment of the longitudinal impact of an intervention, unlike RCTs. What stood out was how the graphical and numerical presentation of the outputs allowed for easy understanding by the knowledge of epidemiological and statistical methods (Bernal et al., 2017).

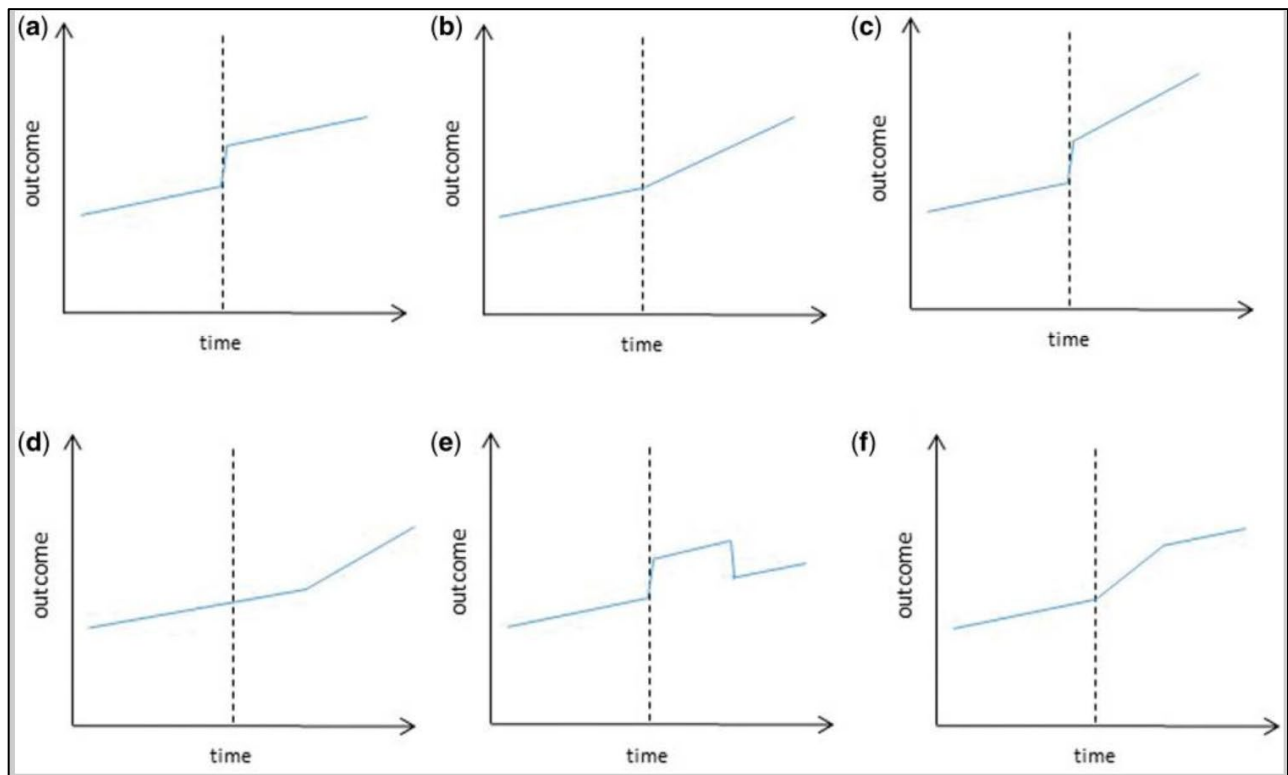


Figure 2 Use of impact models used in ITS. Source : (Barone-Adesi et al., 2011)

(a) Level change; (b) Slope change; (c) Level and slope change; (d) Slope change following a lag; (e) Temporary level change; (f) Temporary slope change leading to a level change.

In the above model, Barone-Adesi (Barone-Adesi et al., 2011) assumed no lags in the change in levels of

Acute Coronary Events (ACEs). This was supported by past evidence showing the short-time disappearance of acute cardiovascular risks derived from passive smoking (Barone-Adesi et al., 2011).

Influenza vaccination among the population has shown evidence of reducing the number and severity of infections. Data collected before and after the immunization helps predict the effect using a time series model. The graphical components of the time series allow for a visual presentation of the level of change and slope, which will be utilized in the subsequent chapters.

1.1.2 Machine Learning

Machine learning (ML) is a branch of computer science and artificial intelligence whose primary focus is using data and algorithms to imitate the learning process of humans and gradually improve the accuracy of learning. It works in three ways: a decision process, an error function, and a model optimization process. ML algorithms are used for two main processes: classification or prediction based on unlabeled or labeled data. The error function is used to evaluate the model's prediction, mainly where known examples are used for comparison and, eventually, accuracy assessment. Model optimization is carried out in the following fashion: continuous weight adjustments are made to reduce the discrepancy between known examples and the model to achieve an accuracy threshold ("What Is Machine Learning (ML)?" 2020).

Numerous innovations are made as the healthcare industry expands into the modern technological era. The advancement of the discipline depends on methods and applications based on artificial intelligence. In the fight against COVID-19, machine learning applications have lately made it possible to speed up testing and hospital response. Using GE's Clinical Command Center deep learning system, hospitals have been able to share, track, and arrange patients, beds, rooms, ventilators, EHRs, and even personnel throughout the pandemic and machine learning, which will lead to faster, more accurate, and simpler diagnosis processes (HealthLeaders, n.d.). Researchers have also utilized artificial intelligence to monitor their progress, identify genetic sequences of SARS-CoV2, and create vaccines (Malone et al., 2020). Support vector machines, linear and logistic regression, naive Bayes, decision trees (random forest, ETC), K-nearest neighbor, and neural networks (multilayer perceptron) are a few ML approaches used for prediction (Sharmila et al., 2017; Er et al., 2010).

1.1.2.1 Long Short-Term Memory (LSTM) Recurrent Neural Network Model

This neural network is recurrent and has a chain-like structure. Instead of having a single neural network layer, it has four layers, each carrying out a unique network function. It is beneficial for projecting the number of new cases for a specific time frame and creating a reasonable long-term projection (Shakeel et al., 2021). Additionally, this model was applied in Canada in March 2020, where it successfully predicted that the first wave of COVID-19 would cease in June 2020—predicting that the pandemic would end in December 2020 and that it would not endure as long as the Spanish flu of 1918 was inaccurate though (Chimmula & Zhang, 2020). In the subsequent chapters, we describe how this deep learning model will aid in predicting the effect of influenza vaccination compared to the traditional time series model.

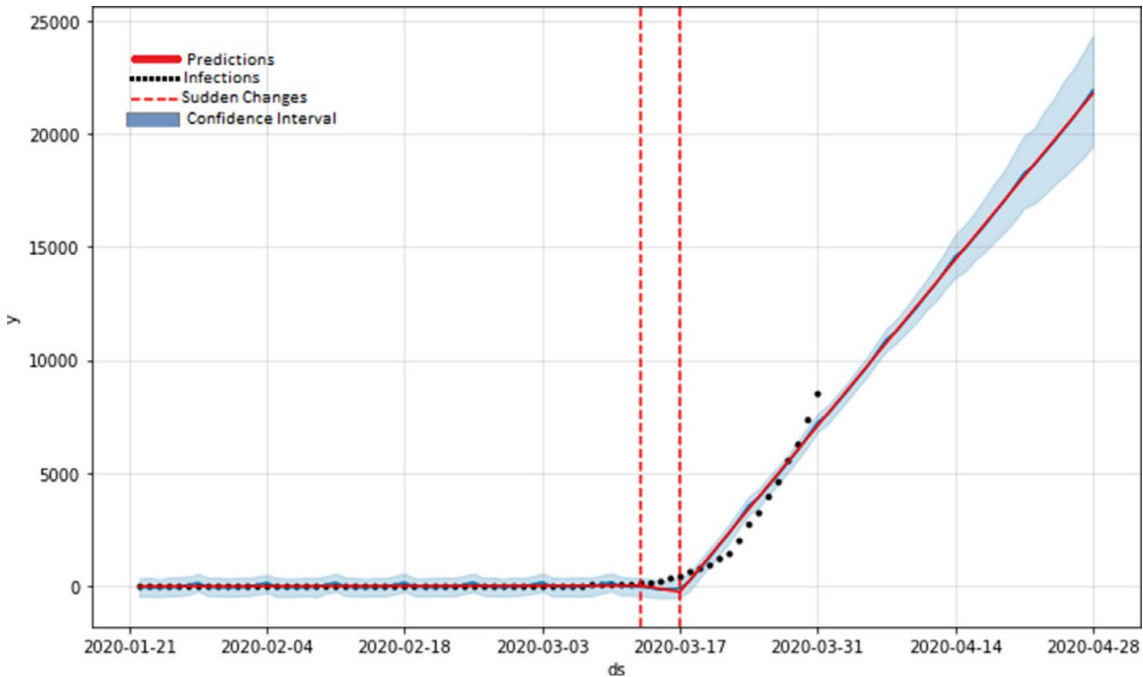


Figure 3 Example of a predictive deep learning model.

“Predictions of the LSTM model on current exposed and infectious cases (Red solid line). The red dotted lines represent the sudden changes from where the infections started following an exponential trend. The black dotted lines in the figure represent the training data or available confirmed cases”. Source : (Chimmula & Zhang, 2020).

The combination of interrupted time series analysis and machine learning techniques provides powerful tools for assessing the impact of interventions, such as the influenza vaccine and making informed predictions about their future effects. This comprehensive approach can improve public health strategies

and decision-making processes related to the influenza vaccine.

1.2 Statement of the Problem

Kenya is classified a Lower Middle-Income Country (LMIC) according to the World Bank classification. In Kenya, low-income communities face significant health challenges, particularly among young children aged five years and below. A 2023 study on causes of death among infants and children in the Child Health and Mortality Prevention Surveillance (CHAMPS) Network in sub-Saharan Africa and South Asia reported that the common causes of death were malnutrition, HIV, malaria, and diarrheal disease. However, when considering the immediate causes of death only, sepsis and lower respiratory tract infection emerged as the top cause of death (Bassat et al., 2023). Another 2023 study for estimating the national burden of respiratory syncytial virus (RSV) in children under five years concluded that their calculations pointed to a significant burden of RSV-related disease and fatalities in Kenya. The data from the study aligned with those from other research that had documented the most significant incidence of illnesses linked to RSV in the first six months of life. More severe consequences (hospitalization and mortality) occur in the first three months (Nyawanda et al., 2023). There is evidence that influenza is a significant cause of respiratory illness in Kenya, particularly among children under five (McMorrow et al., 2015; J. Dawa et al., 2020). One potential solution is the administration of influenza vaccines, which have been shown to reduce the incidence and severity of influenza among young children (Dawa et al., 2020).

In the most recent case where influenza vaccination among children aged five years and below was piloted in Nakuru and Mombasa counties in Kenya, what follows is the need to evaluate the impact of the immunization rigorously. However, intervention such as mass vaccination among children 5 years is a resource intensive undertaking. Evidence of the impact of this type of intervention is significant in the decision making, resource allocation and policy changes to the inclusion of influenza vaccine as part of the Kenya Expanded Program on Immunization (KEPI) vaccines. The KEPI was established by the Kenyan government to provide essential vaccines to children to protect them against various preventable diseases. KEPI vaccines include a range of vaccines recommended by the World Health Organization (WHO) for routine immunization programs. These vaccines are provided free of charge or at a subsidized cost through public health facilities in Kenya as part of the government's efforts to ensure widespread immunization coverage and prevent the spread of vaccine-preventable diseases.

Evidence from other researchers has shown that the most effective method for assessing the effects of such interventions would be RCTs (Hariton & Locascio, 2018). However, they can be expensive and time-consuming to conduct. They often require significant financial resources, skilled personnel, and lengthy follow-up periods despite being recommended as the gold standard measure for assessing the impact of

interventions.

This study seeks to understand better the effect of influenza vaccine administration on children under five years in Kenya. Exploring alternative methodologies, such as interrupted time series analysis and machine learning, is imperative for evaluating the impact of influenza vaccine interventions in Kenya. Leveraging these approaches on available data not only unveils valuable insights into the effects of such interventions but also empowers us to forecast their future trends. Moreover, employing interrupted time series analysis and machine learning techniques offers a cost-effective and streamlined approach to assess the effectiveness of influenza vaccine interventions.

With the rapid advancements in computational hardware and the sophistication of machine learning algorithms, intense Learning, which excels in unraveling system complexities, these tools present a promising avenue for analyzing and forecasting the effectiveness of influenza vaccines in this context. Harnessing interrupted time series analysis and machine learning techniques significantly enhances the accuracy of predicting the efficacy of influenza vaccine interventions.

1.3 Study Objectives

1.3.1 Overall objective

This study aims to compare the effectiveness of two methods for evaluating the impact of interventions in healthcare settings: the ITS design, which is a sturdy quasi-experimental design preferred in a healthcare setting, and the LSTM-RNN, which is a deep learning algorithm motivated by the fact that there has been a rise of the use of ML algorithms to predict outcomes in healthcare settings accurately.

1.3.2 Specific objectives

1. To develop predictive models using interrupted time series (ITS) analysis and machine learning (ML) techniques to forecast the number of medically attended respiratory illness cases following influenza vaccination.
2. To evaluate the effectiveness of influenza vaccination by comparing the number of cases of medically attended respiratory illness before and after vaccination using ITS and ML methodologies.
3. To assess and compare the accuracy of interrupted time series analysis and machine learning in predicting the effectiveness of influenza vaccine interventions.

1.3.3 Research Questions

1. How accurately can predictive models based on interrupted time series analysis and machine

learning forecast the number of medically attended respiratory illness cases following influenza vaccination?

2. What differences exist in the rates of medically attended respiratory illness before and after influenza vaccination, and how do these differences compare between interrupted time series analysis and machine learning techniques?
3. What are the respective accuracies of interrupted time series analysis and machine learning in predicting the effectiveness of influenza vaccination?

1.4 Scope of the study

The research aims to leverage alternative methodologies to evaluate the impact of influenza vaccination in resource-constrained settings in Kenya. It will include a review of predictive models used in similar settings, the development and deployment of the ITS and LSTM-RNN model, and finally, the development of a dashboard that is significant in monitoring the number of cases of medically attended respiratory illness over time.

1.5 Relevance of the study

The study will provide valuable insights into the effectiveness of these two methods for evaluating the effect of influenza vaccine administration in low-income communities in Kenya. The findings will inform the development of more effective methods for assessing the effects of health interventions in similar settings, ultimately contributing to improving vulnerable populations' health outcomes.

Chapter 2: Literature Review

2.0 Introduction

Over the past decades, public health interventions have involved developing and introducing many vaccines to support preventive care against certain illnesses. Vaccine developers utilize randomized controlled studies to illustrate that their products meet the safety and efficacy standards required for licensing. A critical need exists for vaccine assessments to evaluate their performance in real-world conditions, including effectiveness among population subgroups not involved in pre-licensure trials. The evaluations reveal the impacts of the vaccine on healthcare outcomes excluded from pre-licensure trials (Verani et al., 2017). This dissertation aims to establish the extent to which the administration of influenza vaccine to children five years and below influenced the number of medically attended respiratory illnesses in the pilot counties in Kenya.

2.1.0 Evaluating Vaccine Effectiveness, Efficacy, and Impact

Terms like effectiveness, efficacy, and impact have distinctly different meanings in the context of vaccine studies. Efficacy refers to the percentage by which the influenza vaccine reduces the target disease rate among the vaccinated individuals compared to the rate of target disease prevalence among the unvaccinated. The purpose of efficacy is to establish vaccine performance under optimal conditions. Effectiveness measures the disease rate percent reduction in the context of real-world vaccine use. Effectiveness may be the same as efficacy, but its use of the vaccinated population in real-world conditions with multiple influences makes it differ in magnitude. On the other hand, impact refers to the quantified disease reduction at the population level after vaccine introduction. The metric can be expressed as the absolute change or percentage decline in the disease rate. The impact is determined by the combined contribution of vaccine coverage, vaccine effectiveness, and herd effects, including lowered transmission rates and reduced population disease risk (Verani et al., 2017).

The efficacy of immunization programs that have been put into place is evaluated using four primary research designs. (Newall et al., 2014). The most popular designs are observational or ecologic studies, which examine changes in disease burden over time (*Reduction in Acute Gastroenteritis Hospitalizations among US Children After Introduction of Rotavirus Vaccine: Analysis of Hospital Discharge Data from 18 US States* | *The Journal of Infectious Diseases* | Oxford Academic, n.d.), and case-control studies, which compare the vaccination status of infectious cases and their controls across time (*Scopus Preview - Scopus - Document Details - Resolving the Pneumococcal Vaccine Controversy*, n.d.).

Another design uses routinely collected medical databases to estimate rates of infectious diseases in a population, mainly when the vaccination status is unknown, as it is rarely linked to these databases. Monitoring a program's impact requires identifying cases using regularly collected health statistics (such as hospitalizations, emergency department presentations, notifications from infectious disease laboratories, and death notifications). Due to data shortages or constraints, detecting every case using these datasets might not always be possible. For instance, not every case of an infectious disease results in medical attention, so it is not entered into a system regularly. In other cases, it may be impossible to estimate the incidence of new infections due to several unrelated admissions for chronic diseases, such as hepatitis B (Mealing et al., 2016). In this study, we shall use a similar design where data collected routinely and reported to the Ministry of Health databases about the infection rates of medically attended respiratory illness will be used to assess whether the vaccination of children under five years had any effect.

2.1.2 Effectiveness of Influenza Vaccination

Influenza vaccination in children is effective in preventing influenza illness and reducing hospitalizations. Studies have reported vaccine effectiveness (VE) ranging from 25.6% to 78.8% and efficacy values for the influenza vaccine between 25.6% and 74.2%. The researchers in these studies reviewed the evidence on the safety and effectiveness of influenza vaccination in young children. They also highlighted the use of quadrivalent vaccines for influenza prevention. They discussed real-time, reverse-transcriptase polymerase chain reaction and the test-negative design. In summary, influenza vaccination in young children was safe and effective, and quadrivalent vaccines provided a stable immune response against key strains of influenza. In addition, the influenza vaccine showed substantial benefits against outpatient illness in children, and vaccine effectiveness was highest in the 6- to 59-month age group (Булракова et al., 2022; Hood et al., 2023).

A systematic review of multiple seasons found that the pooled VE for any influenza was 46%. The researchers conducted a non-systematic search of studies between 2010 and 2020. Influenza vaccination was found to be effective in preventing influenza in healthy children. The efficacy values ranged from 25.6% to 74.2%, and effectiveness from 26% to 78.8%. The limitation highlighted in the research was the high variability in results due to differences in design, vaccine type, and season included in the studies (Orrico-Sánchez et al., 2023). VE was highest among children in the 6- to 59-month age group compared with older pediatric age groups, as shown by another 2021 study. The researchers did a systematic review and meta-analysis of the literature and found that a test-negative design was used in 37 studies. Ultimately, influenza

vaccination provided moderate protection against hospitalization in children. Although, effectiveness varied by influenza subtype and vaccine type (Boddington et al., 2021).

The vaccine's effectiveness varied by influenza subtype, with higher VE against influenza A/H1N1pdm09 and lower VE against influenza A/H3N2, as reported by a 2022 study. The researchers conducted a non-systematic search of studies conducted between 2010 and 2020. In conclusion, influenza vaccination was effective in preventing influenza in healthy children. Lastly, efficacy values for the influenza vaccine ranged from 25.6% to 74.2%, and effectiveness from 26% to 78.8% (Orrico-Sánchez et al., 2023). Influenza vaccination also provides moderate overall protection against influenza-associated hospitalization in children, with a pooled seasonal IVE against hospitalization of 53.3% for any influenza. These findings support the recommendation for annual influenza vaccination in children to prevent influenza illness and its complications. These papers do not mention ITSA (Interrupted Time Series Analysis) or ML (Machine Learning) in assessing the effectiveness of influenza vaccination in children. Therefore, there is a need to explore the vaccine effect on several cases of medically attended respiratory illness using these methodologies.

2.1.3 Analyzing the effects of vaccination programs

A 2016 study to assess the impact of rotavirus vaccination on emergency department (ED) visits and hospital admissions for acute diarrhea in children under five years in Brazil highlighted the changes following RV1 vaccination. The researchers examined and compared all reported cases of diarrhea at the emergency department and hospital admissions recorded during consultations. The odds ratio of the frequencies of consultations and hospital admissions in the pre-and post-vaccine periods were calculated for statistical analysis. The odds ratio for rotavirus positive was also computed for the pre- and post-vaccine periods, and a chi-squared test was performed to determine whether there was a difference between the periods under investigation. In the end, a 40% decrease in hospital admissions and a 6% decrease in ED visits was attributed to the introduction of the RV1 vaccination (Paulo et al., 2016). Similarly, this study compares the rates of medically attended respiratory illness among children under five years pre and post-vaccination periods.

A 2017 study in Tuscan (Italy) to examine hospitalizations for pneumonia over 12 years for pre-vaccination and vaccination periods for pneumococcal pediatric vaccination reported a 29.1% decline in hospitalization cases for the targeted age group, 0 to 9 years. Between 2002 and 2014, the average annual hospitalization rate for 0-9 years was 989.2 hospitalizations/100,000 residents. This rate indicated a consistent downward

trend, with the highest value recorded in 2002 being 1,415.3/100,000 and the lowest in 2014 being 628/100,000. However, the research uncovered the expected indirect effect on older people was not reported, justifying the Tuscan recommendation to extend the vaccination to subjects for the group older than 64 years (Boccalini et al., 2017). This researcher used the incidence of the disease to monitor changes in the level of disease among the target population.

Interrupted time series analysis was used to assess the impact of pneumococcal conjugate vaccination on hospitalized childhood pneumonia in Taiwan. The researchers compared pneumonia trends throughout the study periods in the individual age groups using segmented regression based on an autoregressive error model. The quarter-time unit was selected to better fit the model and account for the notable seasonal variation in pneumonia outcomes. Stepwise autoregression was used to determine the autoregressive error model's order using quarterly data and an initial order of 5. Following the implementation of the nationwide PCV13 vaccination program, there was a decrease in the prevalence of lobar/pneumococcal pneumonia, lower respiratory tract infections, and pneumococcal parapneumonic illnesses (Lee et al., 2022).

2.2.0 Predictive Modeling for Respiratory Illness

Predictive modeling for respiratory illness has been explored in several studies. One study proposed an automated computational framework for patient-specific deposition modeling, which could optimize treatment plans for respiratory diseases based on patient-specific features such as breathing patterns and lung morphology. The researchers utilized the image processing approach to produce 3D patient respiratory geometries, evaluate airway and lung morphology, and assess deposition compared to in vivo data. However, they experienced difficulty in capturing upper airway anatomy and glottis in CT images. Also, there were challenges in representing instantaneous patient-specific flow patterns observed in vivo. Ultimately, they developed image processing and mathematical modeling pipelines for patient-specific drug deposition predictions and agreed with experimental regional deposition measurements. This paper does not discuss the accuracy of the predictive model (Williams et al., 2023).

Another study used machine learning algorithms, specifically logistic regression, to predict respiratory diseases based on respiratory data, providing insight for better decision-making. The researchers employed a human respiratory model and responses and a machine-learning algorithm for disease prediction. The model matched up to 95.6% with the natural respiratory system, and the logistic regression model gave a distinct separation between diseases (Naskar, n.d.). Additionally, a forecasting model using ARIMA and LSTM methods was developed to predict the occurrence of influenza-like illness and respiratory diseases based on air pollutant data, demonstrating the potential for disease prediction and preventive measures. Considering

everything, the ARIMA displayed better accuracy in the five-year dataset, and LSTM generally outperformed ARIMA by three to seven times (Tsan et al., 2022). The study compared time series and recurrent neural network models for prediction. This research intends to apply similar methodologies to assess the effect of vaccination and compare their accuracies.

Furthermore, the effect of meteorological factors and air pollution on hospital visits for respiratory diseases was evaluated using support vector regression, highlighting the correlation between air pollutants and respiratory diseases. Support Vector Regression (SVR) was used to build regression models, while machine learning was combined with meteorological and air pollution data. To summarize, meteorological factors and air pollution were correlated with respiratory diseases. In addition, machine learning could predict hospital visits for respiratory diseases. (Yang et al., 2023).

Finally, a study analyzed the correlation between air pollutants and respiratory diseases. It used the ARIMA and ridge regression models to predict air pollutant concentrations and the number of respiratory diseases in urban populations. The researchers performed an analysis of the correlation between air pollutants and respiratory diseases. They carried out prediction using the ARIMA model and ridge regression model. Ultimately, air pollutants increase the prevalence of respiratory diseases. The ARIMA and ridge regression models proved useful for prediction (Zhu, 2022).

The limitations of the machine learning models used in these papers include the need for continuous retraining as new data arrives and the potential for variability in model performance. Future research should explore methods for dynamic retraining and improving model stability. The papers also mention the need for further research to validate and refine the predictive models using more extensive and more diverse datasets. The generalizability of the models to different populations and settings is also a limitation that needs to be addressed. Another limitation is the reliance on retrospective data analysis, which may introduce biases and limitations regarding data quality and availability. Prospective studies with larger sample sizes and standardized data collection methods are needed to overcome this limitation.

2.2.1 Interrupted Time Series Analysis in Public Health

Interrupted Time Series Analysis (ITSA) is used in public health research to evaluate the impact of interventions or events on outcomes over time. ITSA involves analyzing data collected before and after an intervention or event to determine if a significant change in the outcome of interest exists. Several studies have utilized ITSA to assess the effects of COVID-19-related public health restrictions on various health

outcomes, such as sexually transmitted and blood-borne infections (STBBI) testing utilization. The researchers employed interrupted time series analysis using GetCheckedOnline program data and segmented generalized least-squared regression modeling. They struggled with the fact that there was limited evidence of the long-term impacts of COVID-19 restrictions on digital STBBI testing utilization. In summary, COVID-19 restrictions initially reduced digital STBBI testing but increased later. Additionally, the researchers concluded that sustained increases in digital testing highlighted the need for accessible and appropriate testing. (Iyamu et al., 2023).

Another application of ITSA involved research on the incidence rate of communicable diseases. Using autoregressive integrated moving average (ARIMA) models, the researchers utilized Interrupted time series analysis. They investigated the analysis of changes in the incidence rate of infectious diseases before and after the COVID-19 epidemic. The research faced challenges with the relatively short collection period after the COVID-19 epidemic and the quick change in policies, requiring further study in the future. As a whole, prevention and control measures during the COVID-19 epidemic significantly impacted notifiable infectious diseases in China. Additionally, the measures most affected respiratory and intestinal contagious diseases (Schaffer et al., 2021). And in the detection of gastroenteritis outbreaks. This was possible through a time series study, which applied segmented linear regression analysis. Notably, presentations increased after water contamination events. Segmented regressions could detect low-present outbreaks (Yuen et al., 2023).

ITSA has also been used to examine the impact of public health policies, such as implementing a public health product tax, on purchasing habits. The researchers used retrospective, descriptive analysis of tax bases and income. They were followed by interrupted time series analysis using the generalized least squares method. In general, it can be said that the public health product tax (PHPT) did not decrease households' unhealthy food purchasing trend. On the other hand, the PHPT generated revenue for health care and health-promoting programs (Csákvári et al., 2023).

Furthermore, research was conducted on an empirical comparison of statistical methods used for meta-analysis of Interrupted Time Series (ITS) studies in public health. The study did not specifically discuss the application of ITS analysis in public health. They applied two ITS analysis estimation methods: fixed-effect and four random-effects meta-analyses. They could only obtain a small number of datasets for analysis; therefore, assumptions made during analysis may not hold. Ultimately, the choice of statistical method had minimal impact on effect estimates and variances, and confidence intervals and p-values could vary depending on the statistical method (Korevaar et al., 2024). Overall, ITSA is a valuable tool in public health research for understanding the effects of interventions and events on health outcomes and informing evidence-based decision-making.

2.2.2 Machine Learning Techniques in Healthcare Forecasting

Machine learning techniques have been increasingly used in healthcare forecasting, including applications in predicting disease outbreaks and vaccination effectiveness. These techniques utilize patient information such as demographics, medical history, and diagnosis to create predictive models for clinical diagnosis, therapeutic efficacy, and prognosis. A recent study explored machine learning techniques, specifically extreme gradient boosting (XGBoost) and logistic regression, to predict measles outbreaks at the county level in the United States. Extreme gradient boosting (XGBoost) and logistic regression were used as supervised learning approaches; in contrast, hierarchical density-based spatial clustering of applications with noise (HDBSCAN) and unsupervised random forest (uRF) were used as clustering algorithms. In summary, XGBoost provided more accurate predictions of measles cases at the county level than logistic regression. Additionally, integrating unsupervised machine learning approaches with supervised models requires further investigation. The limitations highlighted in this research include the optimal strategy for integrating unsupervised machine learning with supervised models, which required further investigation; logistic regression models had higher sensitivity but lower positive predictive value and specificity compared to XGBoost models (Ru et al., 2023).

Furthermore, machine learning can be used to forecast seasonal epidemic peaks. Another 2023 study used this for the forecast specifically for the respiratory syncytial virus (RSV). The researchers explored the machine learning process for generating short-term forecasts and the comparison of different machine learning regression models for forecasting peaks. The limitations were that including seasonality in forecast models could result in overfitting and best-fitting models for typical seasons containing different variables than atypical seasons. Ultimately, machine learning models for forecasting epidemic peaks should not include seasonality, and the best-fitting models for typical seasons differ from atypical seasons (Morbey et al., 2023). A 2022 study investigated the use of machine learning techniques in forecasting the trend of COVID-19 spread with respect to vaccination using long short-term memory (LSTM) and gated recurrent unit (GRU). In conclusion, the research proposed a model to forecast Covid-19 spread based on vaccination data. Additionally, the model allowed researchers to approach COVID-19 spread prediction from a new perspective (Sunil et al., 2022).

In the case of influenza, machine learning methodologies have been used to carry out prediction and forecasting. A 2021 study evaluated the predictive performance of a neural network approach for estimating influenza activity in the United States, using a gated recurrent unit neural network approach and baseline machine learning methods. Overall, gated recurrent unit neural networks improved influenza prediction at long-term time horizons. The highlighted shortcomings included traditional healthcare-based surveillance

systems with inherent reporting delays and the neural network approach lacking improvement by real-time internet search data (Aiken et al., 2021).

Another recent study explored machine learning techniques, specifically Extreme Gradient Boosting (XGBoost), in forecasting influenza incidence as an ordinal variable. They deployed extreme gradient boosting (XGBoost) machine learning methods and statistical and mechanistic models. Several accuracy indicators were used to assess the machine learning framework's prediction performance, and it was compared to baseline models. The most accurate prediction method was the XGBoost model, whose accuracy increased as prediction time horizons increased. According to their findings, there was a chance that the machine learning framework for influenza-like illness (ILI) forecasting could be widely used as an effective public health tool in the future. The deficiencies highlighted by the research; the results were not directly comparable to other forecasting studies, and the XGBoost model had not been validated against other data sources (Wang et al., 2023). The use of machine learning, along with technologies like IoT, cloud computing, and fog computing, has the potential to revolutionize healthcare solutions and improve patient outcomes.

The above researchers discuss using machine learning and deep learning for disease prediction and forecasting. However, they do not specifically mention applications assessing vaccination effectiveness. The following research will explore the use of deep learning to evaluate the effect of an influenza vaccination program on the number of cases of medically attended respiratory illness.

2.2.3 Evaluating Model Precision in Interrupted Time Series Analysis and Deep Learning Models

Interrupted time series (ITS) models were compared regarding predictive accuracy. The performance of different models, including ARIMA, K-Nearest Neighbors (KNN), Support Vector Regression (SVR), and Long-Short Term Memory (LSTM), was evaluated using metrics such as Mean Squared Error (MSE), Mean Absolute Error (MAE), Median Absolute Error (Median AE), and Root Mean Squared Error (RMSE) in the context of classical and machine learning techniques for time series forecasting. The results showed that the KNN algorithm had better accuracy than the other models, particularly in the middle and long terms. Additionally, machine learning outperformed ARIMA for shorter-term predictions (Jaousse et al., 2023).

Another 2021 study examined ordinary least squares (OLS), generalized least squares (GLS), Newey-West (NW), ARIMA, and restricted maximum likelihood (REML) statistical methods for estimating ITS model parameters. OLS was preferred for series with fewer than 12 points, while REML was preferred for longer series. All methods produced unbiased estimates of the level and slope change, although all methods

underestimated the extent of autocorrelation. Therefore, the standard errors (SEs) were generally too small, and the confidence intervals did not achieve the nominal level of coverage. The Durbin-Watson test for the presence of autocorrelation performed poorly, except in cases of long series and high underlying autocorrelation. Ultimately, it was necessary to exercise caution when interpreting the results from all methods, as the confidence intervals were likely to be excessively narrow (Turner et al., 2021).

AIC and BIC are commonly used methods for evaluating time series models. AIC (Akaike Information Criterion) and BIC (Bayesian Information Criterion) are used to balance the complexity and accuracy of the model. AIC and BIC are often used to select the appropriate lag order for autoregressive (AR) models. In the case of two recent studies, differential operation was used to obtain a stable sequence and selection of appropriate models using AIC and BIC. The research proposed an adaptive model based on AIC and BIC ARIMA predictions. However, despite AIC and BIC being commonly used methods for choosing lag order. Analytic approximate cross-validation (APCV) was a more consistent method for choosing lag order in autoregressive models. The results of simulations also demonstrated that APCV performed comparably and, in certain instances, even better than AIC, AICc, and BIC (Y. Zhang & Meng, 2023); Han, 2022). Overall, AIC and BIC are widely used criteria for evaluating time series models, but alternative methods like APCV and BIC may provide more consistent and efficient results.

Machine learning models can be assessed using metrics such as mean absolute error (MAE), mean squared error (MSE), root mean squared error (RMSE), and R-squared. When quantifying average error, MAE is more interpretable than MSE or RMSE (Kumar, 2023). RMSE and MAE are widely used metrics for evaluating models, but there is confusion over their use. RMSE is optimal for normal (Gaussian) errors, while MAE is optimal for Laplacian errors. However, MAE can be decomposed into bias, proportionality, and unsystematic error. The three-part decomposition of MAE provides more straightforward information on the model-error distribution (Robeson & Willmott, 2023; Hodson, 2022).

R-squared is another metric used to evaluate the goodness of fit of a regression model. It measures the proportion of the variance in the dependent variable that the independent variables can explain. In their investigation, they demonstrated that R-squared was more accurate and illustrative than symmetric mean absolute percentage error (SMAPE) using several use cases and examples: In actuality, a high R-squared score was only produced if the regression accurately predicted the majority of the ground truth elements for each ground truth group, taking into account their distribution. Rather than considering their distribution, SMAPE concentrated on the relative distance between each predicted value and its matching ground truth element. SMAPE proved ineffective in their investigation for identifying subpar regression models (Chicco et al., 2021). All the above-discussed metrics provide valuable information about the performance and

accuracy of machine learning models in predicting and analyzing various phenomena.

2.3 Conclusion

In conclusion, the reviewed studies have highlighted the use of different approaches to assess the impact of vaccination programs. Research methodologies such as observational studies, Randomized Controlled Trials (RCTs), and the utilization of routinely collected medical databases to assess the rates of infectious diseases in the population, especially in cases where vaccination status is unknown, have been employed. Various statistical approaches have been investigated to assess the impact of vaccination programs across different population groups; examples include estimation and comparison of disease rates before and after vaccination, calculation of odds ratios, odds of disease occurrence before and after vaccination, monitoring of the average hospitalization rates and the use Interrupted Time Series Analysis (ITSA). The effectiveness of models such as ARIMA, SVR, GRU, and Long Short-Term Memory (LSTM) in predicting respiratory illnesses has been demonstrated, with previous research underscoring their ability to analyze time series data.

OLS, GLS, Segmented linear regression, ARIMA, ARIMAX, and REML are some of the commonly used ITSA implementations. XGBoost, GRU, LSTM, SVR, and ARIMA ML models have been recently utilized in public healthcare, especially for time series data. The metrics used to assess the accuracy of the above methods include but are not limited to using R-squared, MSE, MAE, RMSE, SMAPE, AIC, BIC, and APVC.

In summary, combining traditional statistical-based techniques and machine learning approaches can provide valuable insights and determine the effects of intervention in public health. However, each method has advantages and disadvantages, and the best strategy will rely on every aspect of the data and the desired outcomes of the research.

2.5 Research Gap

There is a need to explore ITS and ML techniques in a contrasting manner in assessing the effect of public health interventions. The two methods will aid in developing a predictive model, which will be applied to limited disease data collected routinely in medical databases. Investigating these methods is helpful for scenarios where vaccine administration in a population is required using limited routinely collected data. From the literature review, this has not been widely explored.

Chapter 3: Methodology

This chapter describes the steps undertaken to realize the objectives outlined in Chapter 1.

3.0 Study Design

This is a retrospective case-control multisite study. The primary sites of focus for this study are Mombasa and Nakuru counties. Njoro and Nakuru North are in Nakuru, and Jomvu and Likoni are in Mombasa. Where retrospective data collection has been conducted within each sub-counties that participated in the influenza vaccine demonstration project: classified as cases. Controls are the Molo sub-county in Nakuru and the Changamwe sub-county in Mombasa, which were not involved in the vaccine demonstration project. The process for identifying the participating sites in each of these sub-counties mirrored the implementation of the vaccine demonstration project. The two sites were selected as part of influenza surveillance sentinel sites since 2006; at the county referral hospitals in these counties. Notably, Nakuru has had frequent Severe Acute Respiratory Infection (SARI) outbreaks (Kna, 2022;Kahenda, n.d.). Data from the surveillance platform was useful in selection the above counties for the vaccine demonstration study.

3.1. Data Sources and Collection Methods

The data was sourced from the District Health Information System (DHIS), also called the Kenya Health Information System (KHIS). Data extraction was carried out by downloading the aggregated datasets following the indicators used to identify cases of respiratory illness. The number of medically attended respiratory illnesses was quantified by the number of lower, upper, and chronic respiratory illnesses (Asthma) in children under five years. The data was entered monthly for the respective sub countries. The different study periods were defined as the baseline phase spanning from January 2019 to October 2019, the 2-year intervention phase from November 2019 to November 2021, and the subsequent 1-year post-intervention phase from December 2021 to December 2022. The figure below shows the indicators used for data extraction.

Table 1 Medically attended respiratory illness indicators in the DHIS.

County	Indicators	Frequency of collection	Case/Control Site
Mombasa	Likoni Sub County Upper Respiratory Tract Infections <5 yrs	Monthly	Case
	Likoni Sub County MOH 705A Rev 2020_ Lower Respiratory	Monthly	Case

	Tract Infections <5 yrs		
	Likoni Sub County Asthma <5 yrs	Monthly	Case
	Jomvu Sub County Upper Respiratory Tract Infections <5 yrs	Monthly	Case
	Jomvu Sub County MOH 705A Rev 2020_ Lower Respiratory Tract Infections <5 yrs	Monthly	Case
	Jomvu Sub County Asthma <5 yrs	Monthly	Case
	Changamwe Sub County Upper Respiratory Tract Infections <5 yrs	Monthly	Control
	Changamwe Sub County MOH 705A Rev 2020_ Lower Respiratory Tract Infections <5 yrs	Monthly	Control
	Changamwe Sub County Asthma <5 yrs	Monthly	Control
Nakuru	Nakuru North Sub County Upper Respiratory Tract Infections <5 yrs	Monthly	Case
	Nakuru North Sub County MOH 705A Rev 2020_ Lower Respiratory Tract Infections <5 yrs	Monthly	Case
	Nakuru North Sub County Asthma <5 yrs	Monthly	Case
	Njoro Sub County Upper Respiratory Tract Infections <5 yrs	Monthly	Case
	Njoro Sub County MOH 705A Rev 2020_ Lower Respiratory Tract Infections <5 yrs	Monthly	Case
	Njoro Sub County Asthma <5 yrs	Monthly	Case
	Molo Sub County Upper Respiratory Tract Infections <5 yrs	Monthly	Control
	Molo Sub County MOH 705A Rev 2020_ Lower Respiratory Tract Infections <5 yrs	Monthly	Control
	Molo Sub County Asthma <5 yrs	Monthly	Control

3.2 Data Preprocessing

This is a crucial step in research, and what follows extraction is pre-processing and transformation. This encompasses checking for completeness, consistency, accuracy, cleaning data cleaning, and checking for completion in readiness for modeling.

3.2.1 Data Cleaning

This involves using all the variables in the correct format, updating data types, generating proper labels, and generating other columns from existing columns. The objective is to ensure that the final analytic

dataset is reliable and complete.

3.2.2 Data Exploration

This involves understanding the structure of the data, attributes, and underlying patterns. Approaches used here include checking for descriptive summaries (frequencies, percent, mean, median, standard deviation, and variance); this allows for the identification of potential errors or anomalies.

3.2.3 Handling Missing Mata

Handling missing data is crucial to data analysis and machine learning tasks. Missing data can occur due to various reasons, such as equipment failure, human error, or simply because the data was not collected. Dealing with missing data appropriately ensures that analysis results are accurate and unbiased. This research examined summary statistics and visualizations using programming libraries that provide functions to detect missing values. The number of upper respiratory illness cases was not entered between January 2019 and September 2020. Imputation was selected as a method of handling the missing data using the means.

Table 2 Handling missing values.

Means of handling missing data	Steps
Deleting Rows or Columns	One straightforward approach is to remove rows or columns with missing values entirely. While this approach is simple, it may lead to the loss of valuable information, especially if the missing values are non-random.
Imputation:	Imputation involves replacing missing values with substituted values. Standard imputation methods include:
	Mean/Median/Mode Imputation: Replace missing values with the mean, median, or mode of the observed values in the column.
	Forward Fill/Backward Fill: Propagate the last known value forward or backward to fill in missing values in time series data.
	Interpolation: Estimate missing values based on the values of adjacent data points.

	Regression Imputation: Predict missing values using regression models based on other variables in the dataset.
	K-Nearest Neighbors (KNN) Imputation: Replace missing values with the average nearest neighbors' values in multi-dimensional space.

3.2.5 Outlier Detection and Handling

The research used box plots and cluster plots to detect outliers. This is an essential step in data analysis and modeling to ensure the accuracy and reliability of results. However, the data had no outliers.

3.2.6 Data Type Conversion

Data type conversion, also known as data type casting or type conversion, refers to converting a value from one data type to another. This prepares the data for different operations or functions requiring specific data types' inputs.

3.3 Data Transformation

3.3.1 Creation of Categorical Variables

An intervention variable was created following the definition of the study period, coded as 1 during intervention and 0 for pre- or post-intervention.

3.4. Preprocessing Steps for ITS Analysis

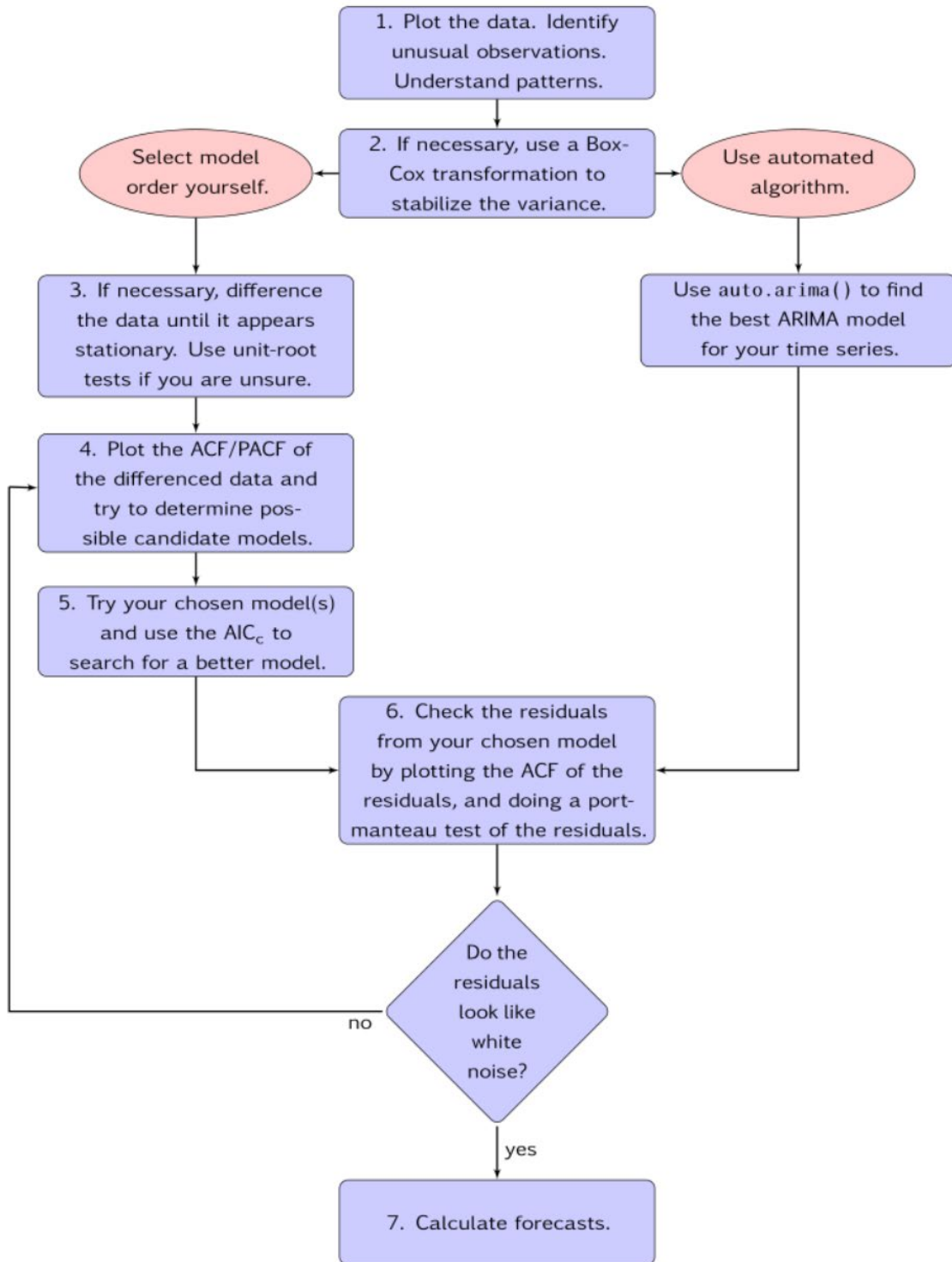


Figure 4 Flow chart for ARIMA model selection. Source: (8.7 ARIMA Modelling in R | Forecasting, n.d.)

3.4.1 ARIMA Model Specification

Three factors delineate the ARIMA model, which is specified as $ARIMA(p,d,q)$ where p , d , and q represent the quantity of lagged (or historical) observations to take into account for autoregression, the number of

iterations the original observations undergo differencing, and the magnitude of the moving average window, respectively. The subsequent equation illustrates a typical autoregressive model. As implied by the nomenclature, the updated values of this model are contingent solely upon a weighted linear amalgamation of its preceding values. The p antecedent values are symbolized as AR(p) or an autoregressive model of order p.

Equation 1

$$Y_t = \alpha + \beta_1 Y_{t-1} + \beta_2 Y_{t-2} + \dots + \beta_p Y_{t-p} + \epsilon_t$$

The moving average:

Equation 2

$$y_t = c + \epsilon_t + \theta_1 \epsilon_{t-1} + \theta_2 \epsilon_{t-2} + \dots + \theta_q \epsilon_{t-q},$$

The computation of the future value y_t involves the consideration of errors ϵ_t originating from the preceding model. Consequently, every subsequent term delves deeper into the historical errors to adjust for the inaccuracies of the present model. The determination of the value of q depends on the extent of the historical window under examination. Thus, the model described above can be specifically identified as a moving average of order q or succinctly denoted as MA(q).

3.4.2 Determining Stationarity

Time series data should be rendered stationary in order to eliminate any apparent correlation and collinearity with previous data. Within stationary time-series data, the characteristics or magnitude of a specific observation are independent of the timestamp at which it is recorded. To eliminate this correlation, ARIMA employs differencing techniques to achieve stationarity in the data. Differencing, in its basic form, consists of calculating the disparity between two consecutive data points. The accurate representation of this is modeled by

Equation 3

$$\begin{aligned}y'_t &= y_t - y_{t-1}. \\y'_t &= y_t - y_{t-1} = y_t - By_t = (1 - B)y_t.\end{aligned}$$

Where the backshift operator denoted by B is defines as

Equation 4

$$By_t = y_{t-1}.$$

Differencing for the purpose of achieving stationarity in data may not always be a straightforward process. It is possible that employing multiple rounds of differencing could offer further improvements, if deemed necessary. Specifically, applying differencing to the data d times results in a dth-order differenced dataset. If $d=2$,

Equation 5

$$\begin{aligned}y''_t &= y'_t - y'_{t-1} \\&= (y_t - y_{t-1}) - (y_{t-1} - y_{t-2}) \\&= y_t - 2y_{t-1} + y_{t-2}.\end{aligned}$$

Or,

Equation 6

$$y''_t = y_t - 2y_{t-1} + y_{t-2} = (1 - 2B + B^2)y_t = (1 - B)^2y_t$$

A generality is observed to be emerging in this context. Consequently, a d-order differenced series can be expressed as:

Equation 7

$$(1 - B)^d y_t$$

The final ARIMA model equation:

Equation 8

$$y'_t = c + \phi_1 y'_{t-1} + \dots + \phi_p y'_{t-p} + \theta_1 \varepsilon_{t-1} + \dots + \theta_q \varepsilon_{t-q} + \varepsilon_t$$

the symbol y'_t represents a differenced series of the first order. It is possible to apply differencing \mathbf{d} times. Regarding the parameters c , ϕ_1 and θ_i , they undergo updating through maximum likelihood estimation (MLE), similar to the process in linear regression. Nevertheless, in contrast to linear regression, the estimation and prediction of a time series using ARIMA allow for the generation of a wide array of future forecasts based on a single trained model, eliminating the need for external data in the inference stage.

3.4.3 Autocorrelation

Time series data frequently exhibit correlations with past observations, leading to a lack of independence in their distribution. This phenomenon is known as autocorrelation or serial correlation. When autocorrelation is present, the standard assumptions of regression analysis are not met. To address this issue, it is common practice to difference the non-stationary data in order to eliminate autocorrelation. Therefore, any required data transformations should be carried out prior to autocorrelation testing (Schaffer et al., 2021).

Autocorrelation functions (ACFs) serve as useful tools for examining both stationarity and autocorrelation within a time series. Specifically, an ACF illustrates the degree of correlation between individual observations and preceding values at different time lags, with a lag representing the temporal gap between a given observation and its antecedent values. Complementary to the ACF, the partial autocorrelation function (PACF) captures the correlation between a current observation and past values not accounted for by correlations at lower-order lags. To illustrate, the PACF value at lag 4 denotes the correlation between an observation Y_t and the previous observation at lag 4 Y_{t-4} , adjusted for correlations with Y_{t-3} , Y_{t-2} , and Y_{t-1} . In the context of a stationary series, the autocorrelation portrayed in the ACF plot should exhibit rapid decay; conversely, in a non-stationary series, the ACF will demonstrate a slower rate of decay (Schaffer et al., 2021).

3.4.3 Seasonality

Seasonality pertains to fluctuations with a predetermined or recognized frequency, manifesting at consistent temporal spans, for instance, time of year or day of the week. The presence of seasonality in temporal sequences of health-related data is widespread and may stem from inherent factors like meteorological cycles, or organizational procedures like the impact of weekends or holidays. The degree of seasonality is contingent upon the time unit of the series; for example, seasonality is infrequent in time series measured annually. When dealing with seasonal monthly data, it is probable to observe notable autocorrelation at a lag of 12 in the ACF graph. Within the realm of ARIMA modeling, addressing seasonality typically involves implementing the seasonal difference. Specifically, in the case of monthly data, one would compute the

disparity between each data point and its preceding value at a lag of 12 ($Y_t - Y_{t-12}$). For quarterly data, a lag of 4 would be employed. It is important to note that in the process of calculating the seasonal difference for monthly data, the initial 12 data points are omitted due to the inability to compute the seasonal difference for these observations (Schaffer et al., 2021).

3.4.4 Transfer Functions

Another benefit of ARIMA models is their capacity to extend beyond the fundamental intervention impact shapes and represent more intricate impacts through the use of "transfer functions." These functions elucidate the correlation between the intervention and the outcome series Y_t , altering the connection between the inputs (step change, pulse, ramp) and the time series to depict more intricate relationships, like gradual level shifts or a pulse that diminishes gradually over time while also accounting for lagged effects.

Equation 9

$$Y_t = \mu + \frac{\omega_0 + \omega_1 B + \omega_2 B^2 + \dots + \omega_h B^h}{1 - \delta_1 B - \delta_2 B^2 - \dots - \delta_r B^r} X_t + \varepsilon_t$$

the backshift operator B is denoted as $B^p Y_t = Y_{t-p}$. Within the transfer function framework, ω_0 signifies the initial impact value during the intervention period (T), δ denotes the rate of decay, and X_t symbolizes the intervention variable (such as step change, pulse, or ramp). Researchers are required to provide values for h and r ; h indicates the timing of the effect, whereas r characterizes the decay behavior (Schaffer et al., 2021).

Function	Values for h and r	Transfer function	Response i at times 0 through k post-intervention	Form of response	Interpretation
Step function $S_t = \begin{cases} 0, & \text{if } t < T \\ 1, & \text{if } t \geq T \end{cases}$	$h = 0,$ $r = 0$	ω_0	$i_0 = \omega_0$ $i_1 = \omega_0$ $i_2 = \omega_0$ \dots $i_k = \omega_0$		The time series increases by ω_0 immediately following the intervention, and remains at this new level for the duration of the study period.
	$h = 0,$ $r = 1$	$\frac{\omega_0}{(1 - \delta_1 \beta)}$ $(\delta_1 < 1)$	$i_0 = \omega_0$ $i_1 = \omega_0(1 + \delta_1)$ $i_2 = \omega_0(1 + \delta_1 + \delta_1^2)$ \dots $i_k = \omega_0(1 + \delta_1 + \delta_1^2 + \dots + \delta_1^k)$		The time series increases by ω_0 immediately following the intervention, and increases by $\omega_0 \delta_1^k$ each subsequent time point until it reaches a new level, calculated by $\frac{\omega_0}{(1 - \delta_1)}$.
Pulse function $P_t = \begin{cases} 0, & \text{if } t \neq T \\ 1, & \text{if } t = T \end{cases}$	$h = 0,$ $r = 0$	ω_0	$i_0 = \omega_0$ $i_1 = 0$ $i_2 = 0$ \dots $i_k = 0$		The time series increases by ω_0 immediately following the intervention and returns to baseline immediately afterwards.
	$h = 0,$ $r = 1$	$\frac{\omega_0}{(1 - \delta_1 \beta)}$ $(\delta_1 < 1)$	$i_0 = \omega_0$ $i_1 = \omega_0 \delta_1$ $i_2 = \omega_0 \delta_1^2$ \dots $i_k = \omega_0 \delta_1^k$		The time series increases by ω_0 the time of the intervention, and decays by $(1 - \delta_1)$ each subsequent time point.
Ramp function $R_t = \begin{cases} 0, & \text{if } t < T \\ t - T + 1, & \text{if } t \geq T \end{cases}$	$h = 0,$ $r = 0$	ω_0	$i_0 = \omega_0$ $i_1 = 2\omega_0$ $i_2 = 3\omega_0$ \dots $i_k = (k + 1)\omega_0$		The time series increases by ω_0 at each time point.

Figure 5 Description of transfer functions within interrupted time series analysis in the context of ARIMA.
Source: (Schaffer et al., 2021)

3.5 Building the ARIMA Model

3.5.1 Plotting Data to Understand Underlying Patterns

Plotting data is essential for comprehending patterns: Prior to engaging in model fitting, it is imperative to visually represent the time series data to identify patterns, including existing trends, seasonal variations, and outliers. Plots were generated for the time series data for cases and controls.

3.5.2 Choosing Model Parameters (p, d, and q)

Multiple methodologies exist for selecting ARIMA parameters, which can be arrived at through either an examination of the data characteristics or through empirical methods involving model fitting and performance assessment. An exploration of ACF and PACF plots is essential for detecting the overall

correlation present in the dataset. Additionally, the calculation of rolling mean and standard deviation, as well as the application of statistical tests like the Augmented Dickey Fuller (ADF) test, are valuable techniques for determining the characteristics of the time series data. This research used ACF and PACF plots alongside ADT test to determine the overall correlation and to test for stationarity of the dataset.

3.5.2 Determining the Seasonal Decomposition of the Time Series.

This study generated a decomposed plot of the time series data. The different components are discussed in the subsequent chapter.

3.5.3 Plotting the Time Series Data with Intervention

A time series plot with intervention components was plotted.

3.6.4 Fitting the ARIMA Model

An automated ARIMA model was fitted on the time series data. This resulted in auto selection of the final model parameters. The model summary and the metric for assessing model accuracy were presented in table format.

3.7.5 Forecasting the ARIMA Model

A 12-month forecast was generated from the ARIMA model. The confidence intervals of the predictions were presented in a table format and a plot of the forecast was also generated.

3.6 Data Preparation for LSTM-RNN Modeling

3.6.1 Data preprocessing

The cleaned dataset was converted into data samples and split into training and testing sites.

3.6.2 Feature scaling

Entails the adjustment of the range of features to a common scale, usually within a predefined range such as $[0,1]$ or $[-1,1]$. The rationale behind scaling the dataset is due to the presence of features with different

scales, which could impact the performance of machine learning algorithms. The research utilized the standard scaling method by converting them to have a zero mean and unit variance.

3.7 LSTM-RNN Specification

3.7.1. Architecture and Model Setup

The core component of an LSTM network resides in its cell, also known as the cell state, serving as a mechanism to retain a certain amount of memory within the LSTM model, thereby enabling it to recollect previous information. It has 3 gates; input, forget and output. The gates utilized in Long Short-Term Memory (LSTM) models are represented by the sigmoid activation functions, which result in an output value ranging between 0 and 1. Typically, this output value tends to be either 0 or 1 in the majority of instances.

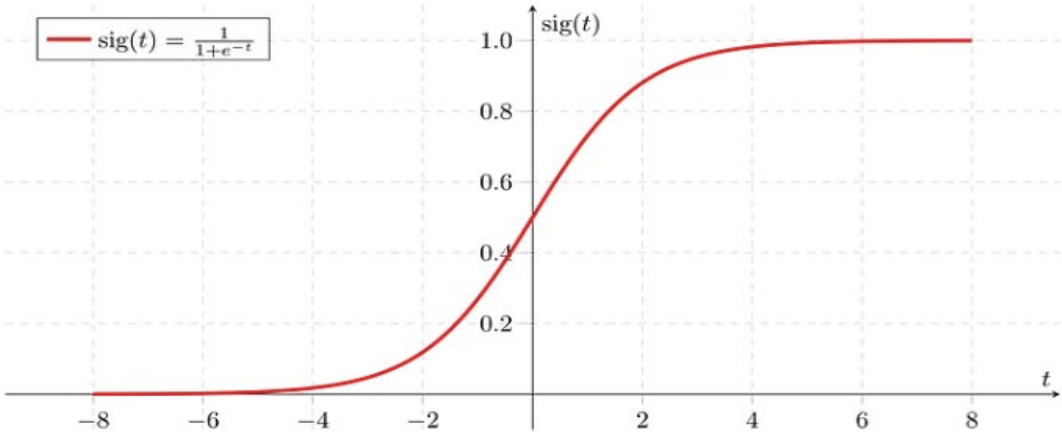


Figure 6 Sigmoid functions for LSTM gates.²

A sigmoid function is employed in gate mechanisms in order to ensure that the output values remain strictly positive and facilitate a definitive decision-making process regarding the inclusion or exclusion of specific features. A value of "0" signifies complete blockage by the gates, while a value of "1" indicates unrestricted passage through the gates. The equations for the gates in LSTM are:

² Source: [LSTM and its equations. LSTM stands for Long Short Term Memory... | by Divyanshu Thakur | Medium](#)

Equation 10

$$\begin{aligned} i_t &= \sigma(w_i[h_{t-1}, x_t] + b_i) \\ f_t &= \sigma(w_f[h_{t-1}, x_t] + b_f) \\ o_t &= \sigma(w_o[h_{t-1}, x_t] + b_o) \end{aligned}$$

i_t → represents input gate. f_t → represents forget gate. o_t → represents output gate. σ → represents sigmoid function. w_x → weight for the respective gate (x) neurons. h_{t-1} → output of the previous LSTM block (at timestamp $t - 1$). x_t → input at current timestamp. b_x → biases for the respective gates (x).

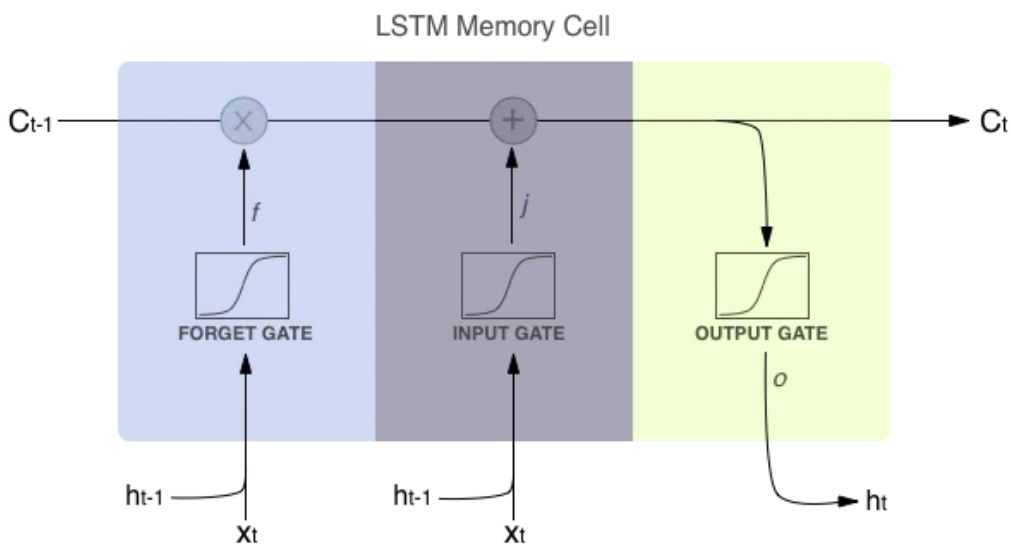


Figure 7 LSTM memory cell. ³

The initial equation pertains to the Input Gate, elucidating the nature of the incoming data to be encoded within the cell state, as elaborated further below. Subsequently, the following equation corresponds to the forget gate, dictating which pieces of information are to be discarded from the cell state. Lastly, the third equation governs the output gate, responsible for bestowing activation upon the ultimate output of the LSTM block at timestamp 't'.

Equation 11

$$\begin{aligned} \tilde{c}_t &= \tanh(w_c[h_{t-1}, x_t] + b_c) \\ c_t &= f_t * c_{t-1} + i_t * \tilde{c}_t \\ h_t &= o_t * \tanh(c^t) \end{aligned}$$

³ Source: [LSTM and its equations. LSTM stands for Long Short Term Memory... | by Divyanshu Thakur | Medium](#)

$c_t \rightarrow$ cell state (memory) at timestamp (t). $\tilde{c}_t \rightarrow$ represents candidate for cell state at timestamp (t). note* others are same as above.

To obtain the memory vector for the present timestamp c_t , the candidate is computed. Subsequently, it can be observed from the equation that at any given timestamp, the cell state is aware of the information it should discard from the preceding state (i.e., $f_t * c_{t-1} + i_t * \tilde{c}_t$) and the information it should prioritize from the current timestamp (i.e., $i_t * \tilde{c}_t$).

Note: * this operation signifies the multiplication of vectors on an element-by-element basis.

3.7.2 Model Compilation:

This step involves the following:

3.7.2.1 Training the NN

This fine tunes weights and biases within neurons aimed at aligning the networks outputs closely to the actual values during testing. This happened through the employment of random weights to generate an initial forecast based on the output values (El-Amir & Hamdy, 2020).

3.7.2.2 Layers, Neurons and Hyperparameters

These were manual specified through trial and error and also following previous work done in the same field and in some cases setting a rule of thumb. For instance, the quantity of neurons that constitute the network, the quantity of layers it possesses, and the methodologies employed for training the network are not estimated by the model (*Applied Deep Learning with TensorFlow 2*, n.d.).

3.7.2.3 Loss Function

It is mandatory to specify the loss function, optimizer, and metrics to monitor during training. Common loss functions for sequence prediction tasks include categorical cross-entropy or mean squared error, depending on the nature of the problem. Mean squared error was chosen for the model.

3.7.2.4 Training Procedure

The NN utilizes an optimizer algorithm which is designed to minimize the loss function during training. It operates by identifying the local minimum of a specific function. The computation is carried out by a central processing unit (CPU) that possesses knowledge of the exact values within the training dataset and continuously strives to reach the point where the loss function is minimized. The whole process is 3-step. Step1; involves forward and backward propagation to acquire both the function value and gradient. Subsequently, a new step is proposed with an increment determined by the current step. The final step

involves integrating this increment into the original function, followed by the repetition of the entire process for a specified number of iterations (Lv et al., 2017).

3.7.2.5 Batch Sizes and Number of Epochs

When training a LSTM model, a batch refers to an iteration consisting of either one sample or multiple samples employed for the purpose of data prediction (*Applied Deep Learning with TensorFlow 2*, n.d.). We selected batch size of 32 which has been approved to be of good choice (Reimers & Gurevych, 2017). An LSTM neural network undergoes multiple iterations of training to reduce the model's error. The completion of one full pass of the entire training dataset through the model signifies the conclusion of an epoch (*Applied Deep Learning with TensorFlow 2*, n.d.). We specified 250 for this model.

3.8.4 Evaluation Metrics

In order to determine the effectiveness of various models in forecasting, it is important to conduct an evaluation of the forecast. RMSE and MAE were selected owing to the fact that they are common measures used when determining the accuracy and rate of error for different models as discussed in chapter 2.

Equation 12

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{n}}$$

Equation 13

$$MAE = \sum_{i=1}^n \frac{|\hat{y}_i - y_i|}{n},$$

Where $\hat{y}_i \rightarrow$ represents the forecasted values and $y_i \rightarrow$ represents the observed values and n is the number of forecasts.

3.9 Comparison of Traditional Time Series Analysis and LSTM -RNN algorithm

The accuracy of predictions was assessed by comparing the respective metrics, while the visualizations were evaluated in the context of the final forecasts.

3.10 Ethical considerations

This study was approved by the ethics committee of the National Commission for Science, Technology & Innovation (NACOSTI), Strathmore University ethical review committee, Ministry of Health Kenya granted access to the use of data from DHIS.

3.11 Limitations

The study may be subject to certain limitations, including potential confounding factors not accounted for in the analysis, the possibility of incomplete or inaccurate data, and the potential for unmeasured changes in healthcare delivery or other factors that may affect the outcome variable.

Chapter 4: Results

4.0 Introduction

The findings and outcomes of the data analysis are presented in this chapter. Tables and graphs have been used to illustrate the findings. The project objectives were considered in the analysis.

4.1 Exploratory Data Analysis

Table 3. Comparison of Intervention Effects on Cases and Control Groups Over Four Years (2019-2022).

Characteristic	2019, N = 12	2020, N = 12	2021, N = 12	2022, N = 12	Overall, N = 48
Intervention					
0	10 (83%)	0 (0%)	1 (8.3%)	12 (100%)	23 (48%)
1	2 (17%)	12 (100%)	11 (92%)	0 (0%)	25 (52%)
Cases	6,956 (6,460, 7,557)	5,683 (4,485, 8,248)	13,426 (11,407, 15,184)	10,968 (9,912, 12,468)	8,864 (6,677, 12,027)
Control	1,879 (1,533, 2,103)	1,657 (1,394, 1,900)	3,385 (3,039, 4,044)	3,547 (2,972, 3,884)	2,419 (1,776, 3,452)
*n (%); Median (IQR), Intervention, 1: Period of intervention, 0 = period pre or post intervention					

The dataset covered the years 2019 to 2022, comprising 48 observations in total. The intervention spanned from 2019 to 2021, lasting for a period of 25 months. The dataset recorded 451,674 cases and 130,170 controls. Median case numbers varied annually, ranging from 6,956 to 10,968, with corresponding interquartile ranges (IQR) between 6,460 and 12,468. Similarly, median control numbers fluctuated, ranging from 1,879 to 3,547, with IQRs from 1,533 to 3,884. These findings highlight the fluctuations observed in case and control numbers throughout the observed period.

4.2 Time Series Plot of The Cases and Controls

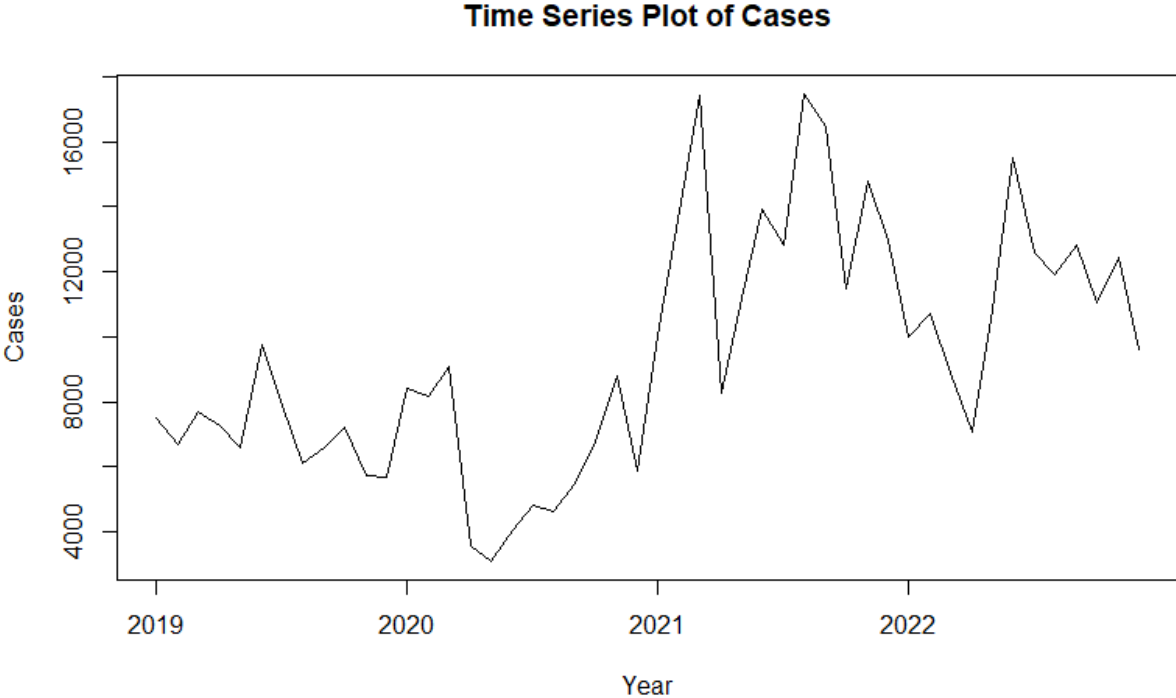


Figure 8 Retrospective Time Series Plots of Cases Data from 2019 to 2022.

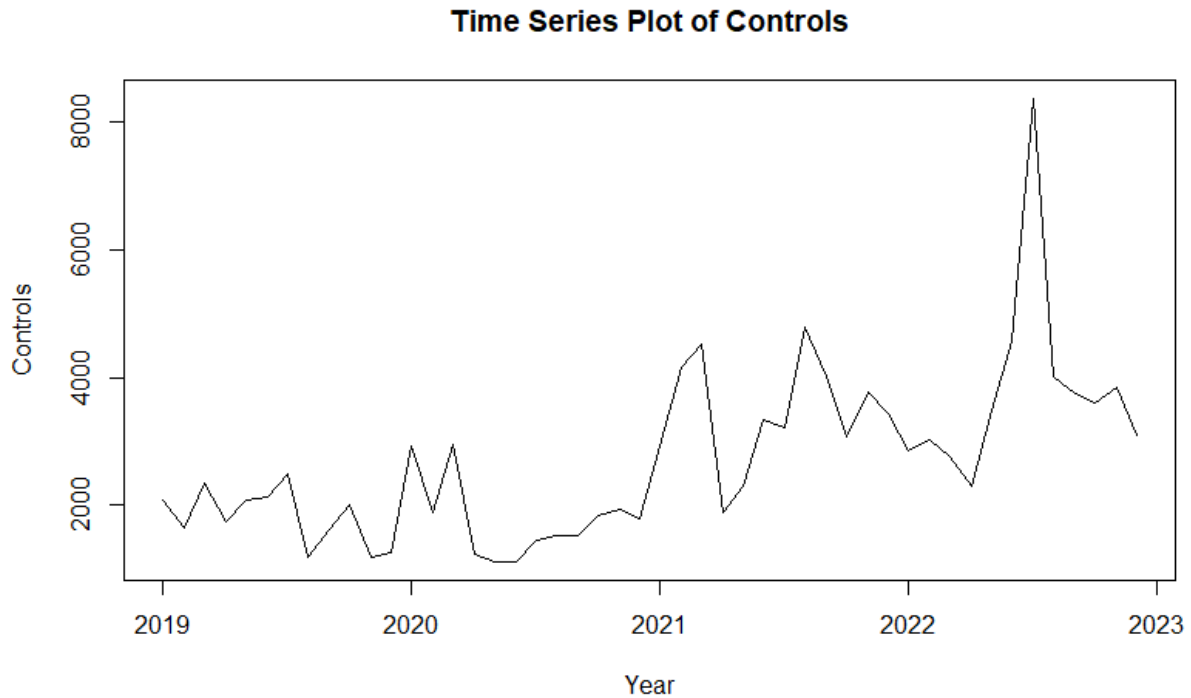


Figure 9 Retrospective Time Series Plots of Control Data from 2019 to 2022.

Figures 8 and 9, highlight the time series plots of cases and control over the 4-year period. They depict a similar rising trend in the number of cases of respiratory illness. Prior to 2021, there were no significant fluctuations in the number cases or controls. However, the implementation of the intervention commenced in November 2019 to November 2021. Figure 10 below is a representation of the pattern observed in the number of cases of medically attended respiratory illness during this period. A sudden decline was observed in 2020, which was transient, succeeded by a gradual but steady rise in case numbers. In contrast, the year 2021 exhibited a sudden surge in cases, accompanied by considerable variations. Overall, 2021 accounted for the highest number of reported cases of medically attended respiratory illness throughout the study period.

Time Series Plot of Cases with Intervention



Figure 10 Time series plot illustrating the intervention period spanning from November 2019 to November 2021.

4.3 Building an ARIMA Model

4.3.1 Determining Stationarity

The Dickey-Fuller test was conducted to determine the stationarity of the time series.

Table 4 Augmented Dickey-Fuller Test

Augmented Dickey-Fuller Test		
Metric	Cases	Controls
Dickey-Fuller	-2.1458	-3.3328
Lag order	3	3
P-value	0.5159	0.07759

The stationarity results for both cases and control with p-values greater than 0.05 means we do not have sufficient evidence to conclude that the data is stationary. As such, to properly use an ARIMA model, a differenced time series data was utilized.

4.3.2 Autocorrelation Plots

ACF and PACF plots were used to identify correlation, lag orders, and aid in model selection for the time series data of cases and controls. In the ACF (Figure 11a), there is significant autocorrelation that gradually dies off at lag 7, similar to Figure 11c where it dies off at lag 3. Nevertheless, when examining the PACF plot (as depicted in Figure 11b), it is evident that the autocorrelation observed at larger lags can be fully accounted for by the autocorrelation observed at smaller lags.

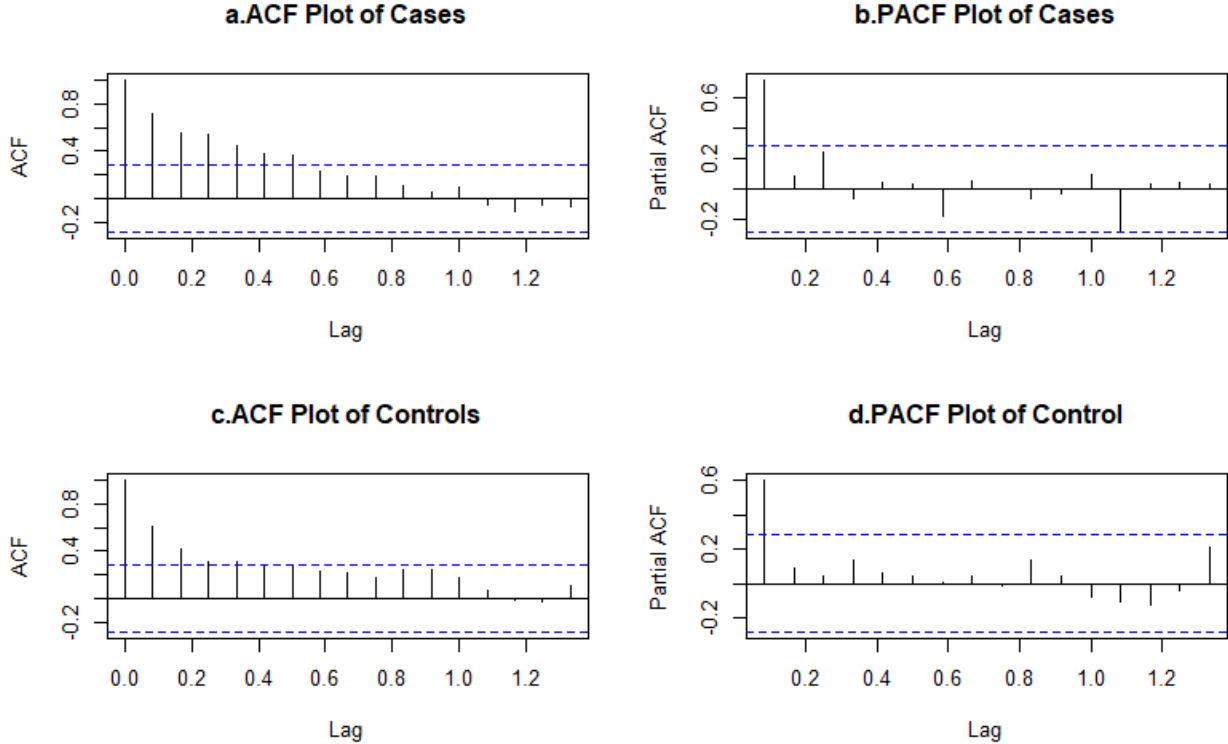


Figure 11 AFC and PACF plots for the cases and control time series before differencing.

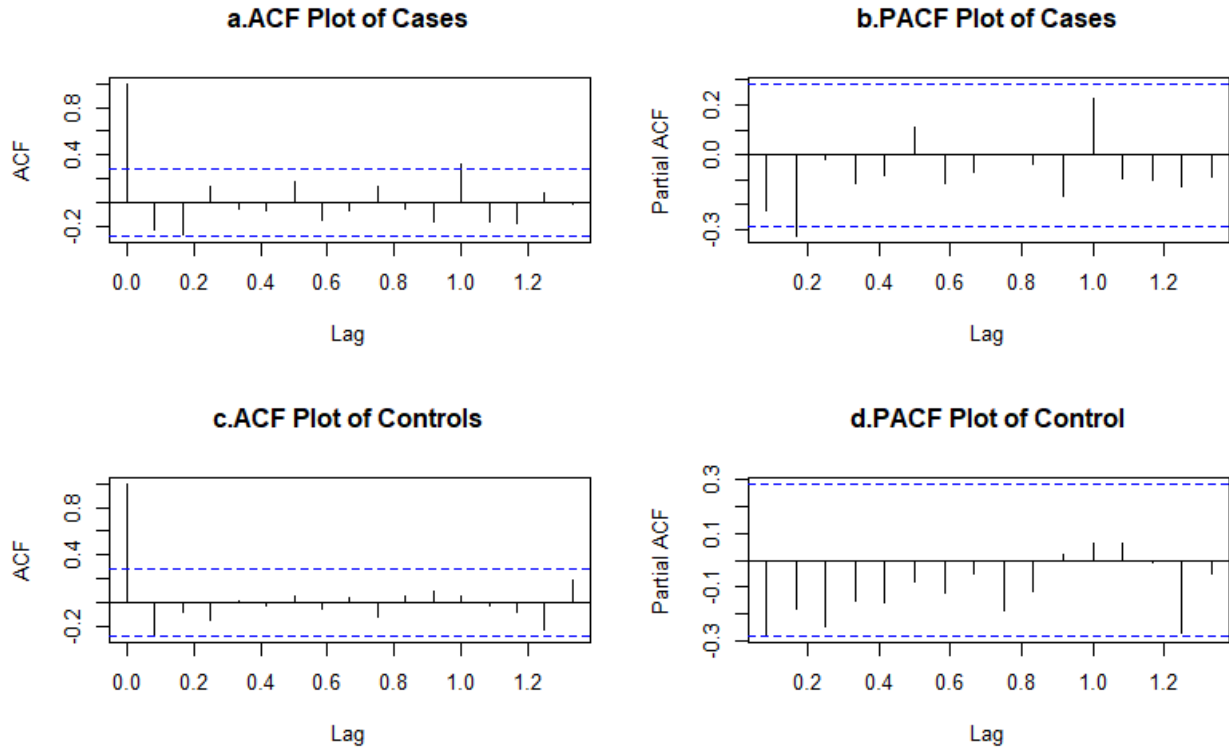


Figure 12 ACF and PACF plots for the cases and control time series post differencing.

Figure 12a and 12c show the eliminated autocorrelation after differencing, compared to figure 11a and 11c. In this instance, the plots of the autocorrelation function (ACF) and partial autocorrelation function (PACF) of the stationary (i.e., differenced) series do not provide significant assistance in determining the p and q parameters, as they do not align with any of the choices outlined in Figure 4. Consequently, an automated procedure, specifically the `auto.arima()` function within the `forecast` package for R, was utilized to ascertain the terms of the ARIMA model.

4.3.4 Fitting an Automated ARIMA Model

For our model, a specific value of $d = 1$ was predetermined (to induce stationarity), along with $D = 1$ (attributed to the existence of seasonality); however, the algorithm selected the most suitable values for p , d , P , and Q . The model identified by the algorithm as the most optimal based on information criteria was $(0,1,2) \times (0,0,1)_{12}$. Stated differently, the autocorrelation order of the model (p) was 0, the moving average order (q) was 2, the autocorrelation order of the seasonal component (P) was 0, and the moving average order of the seasonal component (Q) was 1. Introducing a first-order difference ($d = 1$) and a first-order seasonal difference ($D = 1$) into the model served the purpose of removing trends and enabling stationarity.

4.3.5 Residual check

The residual plots can be observed in Figure 4. No discernible pattern or notable autocorrelation is present in the residuals, which exhibit a distribution that is consistent with normality.

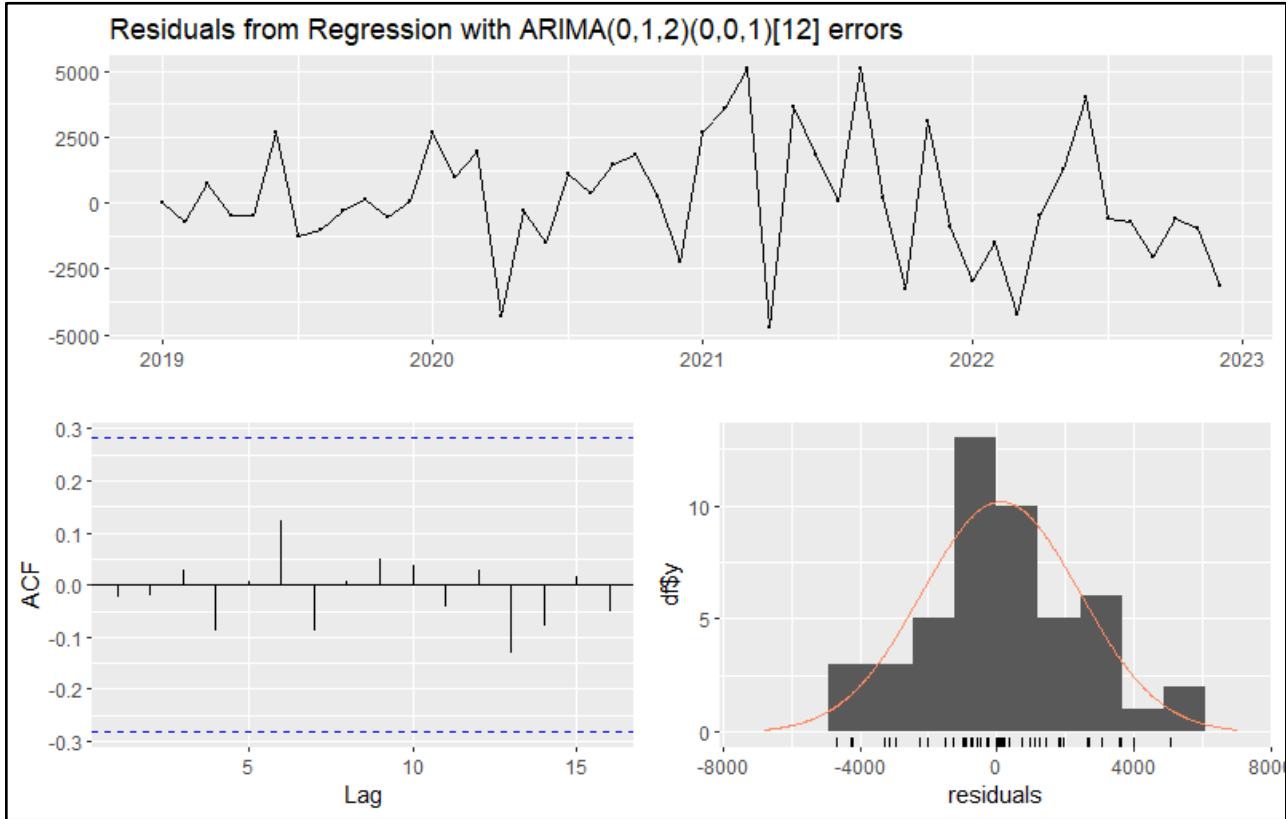


Figure 13 Residual check for the final model, ARIMA (0,1,2) (0,0,1)₁₂.

As observed in Figure 13 above, there is no discernible pattern or notable autocorrelation within the residuals, whose distribution is that from the normal distribution. The p-value obtained from the Ljung-Box test for white noise is 0.9972 when considering 24 lags, we did not reject the null hypothesis since our final model had a good fit.

4.3.6 Final ARIMA Model

The estimated step change was – 698 cases of medically attended respiratory illness (95% CI – 4,477 to – 3,481) while the estimated change in slope was – 248 cases of medically attended respiratory illness per

month (95% CI – 490 to – 5. Figure 14 shows the predicted values in the absence of the intervention plot alongside the observed values. There was an indication that the administration of influenza vaccine among children below 5 years was associated with an immediate and sustained decline of 698 of the number of cases of medically attended respiratory illness, with a further decline of 248 every month. Put differently, there was a reduction 946 (698 +248) of cases of medically attended respiratory illness among children below 5 years in January 2022 compared to what was expected in the absence of influenza vaccine administration.

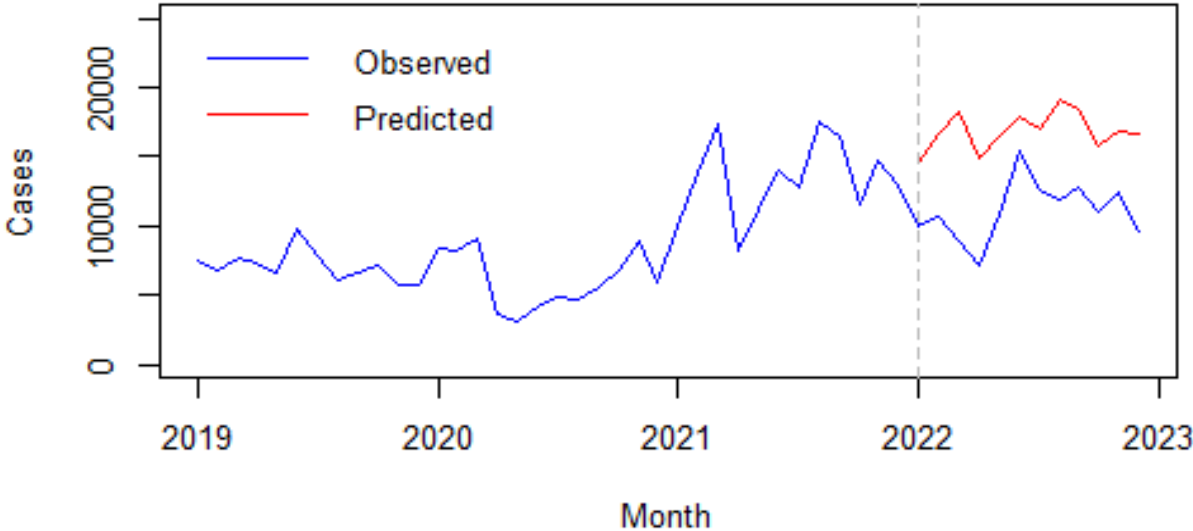


Figure 14 Observed Values and Predicted Values in Absence of Intervention based on ARIMA Model.

Table 5 Error Measures of the ARIMA Model.

Forecast Accuracy							
N	ME	RMSE	MAE	MPE	MAPE	MASE	ACF1
48	119.1282	2281.2	1745.033	-2.64518	19.69101	0.435732	-0.02265

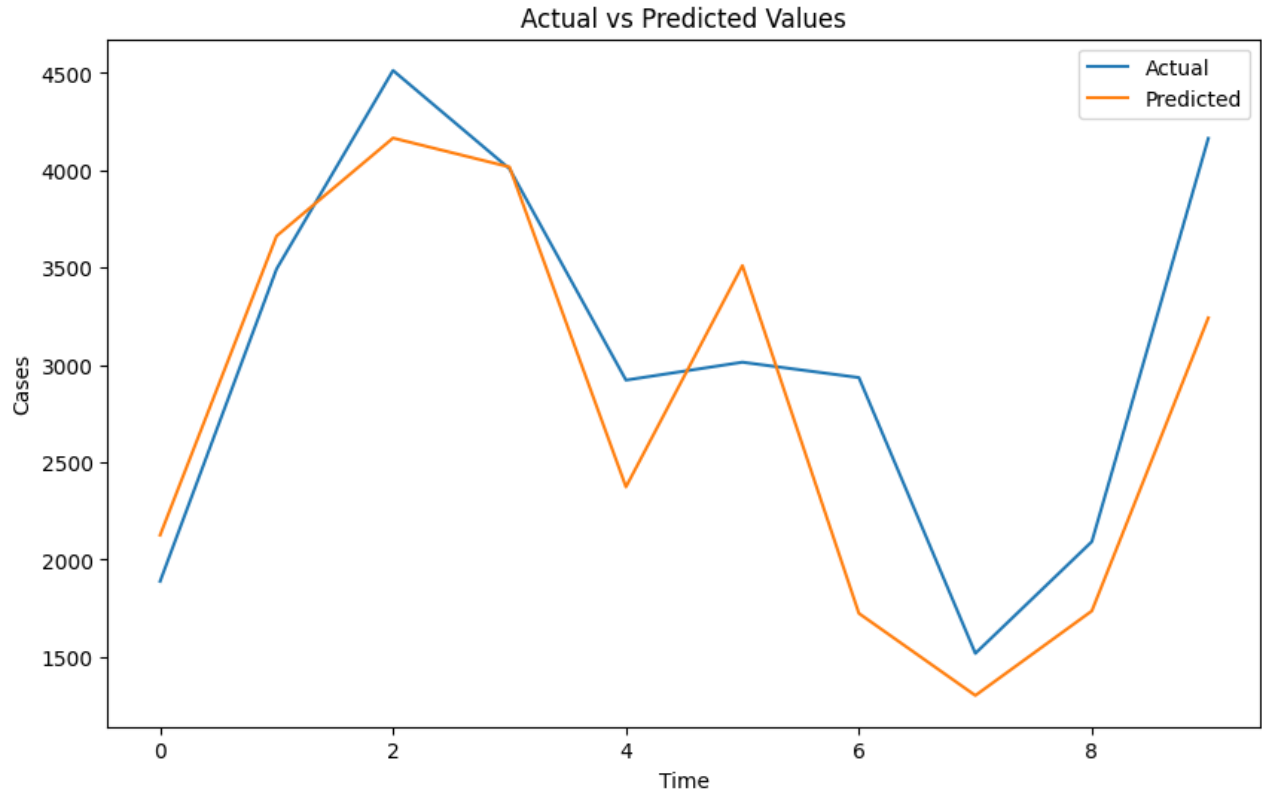


Figure 15 The observed number of medically attended respiratory illness and values predicted by LSTM from 2019 to 2022

Table 6 The forecasting performance of the two models

Model	RMSE	MAE
ARIMA (01,2) (0,0,1) ₁₂	2281.2	1745.033
LSTM	535.0	415.8

Chapter 5: Discussion

Looking at the error scores, RMSE and MAE results in Table 6, the LSTM model outperforms the ARIMA model. LSTM outperforms the ARIMA by a margin that is very significant. Overfitting is a problem in many neural networks and dropout is used as the chosen method in order to mitigate this, naturally the hope is that this would work. Looking at Figure 15, it shows a good fit according to the present standards for a neural network, likely confirming that the risk of overfitting is successfully combated. The ARIMA model, Figure 14 above shows the forecasted values of the number of medically attended respiratory illness in the absence of an intervention, the level and of change and slope allows for the estimation of decline in the number of cases post the intervention period. The would be an increase in cases without the intervention.

ITS model with transfer functions incorporated provides estimation using the step and ramp as showed from the ARIMA model. There was evidence that giving the influenza vaccine to children under the age of five was linked to a rapid and persistent decrease in the number of respiratory illnesses that required medical attention; a sudden decline of 698 cases with an additional monthly fall of 248 cases. Stated differently, compared to what was anticipated in the absence of influenza vaccine delivery, there were 946 ($698 + 248$) fewer instances of medically attended respiratory illness among children under the age of five in January 2022. On the other hand, the LSTM model provided estimates in the cases throughout the study periods. Both models proved useful in answering objective one, which was to build predictive models.

Comparing Figure 14 and 15, the number so cases had a slightly decreasing trend at the beginning and post intervention, showing a possible as association with the administration of the influenza vaccine throughout the study period. However, the decline was not very significant it can be argued the ability of both models to accurately predict the effect of the vaccine was limited by the fact the models did not incorporate other possible confounders and solely relied of the routinely disease data collected for the study period.

As mentioned previously, the estimation of rates and incidences of infectious diseases are quite complex. This does not mean that depending solely on the retrospective cases of respiratory illnesses that were treated medically is a very limited approach for the above prediction is insignificant. This disregards all basic model elements, including the question of whether the influenza vaccine and the type of the virus that is circulating are appropriately matched. It is also impossible to predict increases in the number of cases brought on by unforeseen conditions. The model would probably be far more accurate if indicators could be appropriately incorporated on this. The scope of this dissertation is limited because it solely models respiratory disorders as an index.

Even well-meaning health initiatives might have unexpected results because there is often no data to justify their implementation. In order to discover both intended and unintentional outcomes, provide input to regulators and legislators, enhance the delivery of healthcare, and influence future public health policy, it is imperative that health interventions be evaluated (Lu et al., 2018; Shaw et al., 2019). Like with other analyses, researchers who are interested in assessing interventions should make use of instruments that are appropriate for the specific study topic at hand. Relying on too straightforward methods can result in results that are skewed or misleading (Soumerai et al., 2015). We have emphasized how crucial it is to account for seasonality, autocorrelation, and trends. Segmented regression can also help with these problems to some extent. Usually, this involves adding time and season as covariates to the model, which is sufficient to get rid of simple autocorrelation. And that is why ARIMA is preferred in such scenarios.

The process of choosing the best ARIMA model can occasionally be difficult, subjective, and time-consuming because, as our example illustrates, conventional methods that depend on ACF/PACF plots to determine model orders are frequently uninformative. However, over time, there have been initiatives to streamline and automate the model selection procedure. We used the prediction package for R to apply `auto.arima()`, one such technique that we selected because it is user-friendly and convenient. Although ARIMA modeling is now more widely available thanks to these advancements, it still requires an expert user to apply and interpret the data correctly, just like any other automated statistical technique.

The application of the idea of a network that trains itself on historical data to provide predictions in healthcare setting to forecast the number of cases of medically attended seems naïve. It is unlikely that this model, or any one similar to it, which makes predictions based only on historical data, will be able to predict medically attended respiratory illness cases in the future with any degree of accuracy and consistency. A similar setting might not result in a disastrous outcome if an intervention were used, but it would not be the revolutionary success that it would have been if it had been able to accurately make forecasts. This is not to argue that artificial intelligence and machine learning have no significance; on the contrary, the ability to identify early warning signs is incredibly useful in domains where predictions based on historical data can be made. Regarding their application in healthcare, given that epidemiologists are using machine learning more and more every year, they clearly have some utility, even though they are not a near-magical tool that can predict values in the future based just on historical data.

Chapter 6: Conclusion and Recommendation

Conclusion ITS analysis, especially when combined with a control series, is a powerful study design for assessing population-level health intervention impacts, and its use is increasing. ITS analysis, using ARIMA modelling is a useful tool, as it can account for underlying trends, autocorrelation and seasonality and allows for flexible modelling of different types of impacts.

This was done by fitting ARIMA models and a LSTM neural network to the same dataset and evaluating their performance. Possible improvements to the model that would make it slightly more sophisticated would be to try and accurately incorporate indicators of disease ecology. In the study, the LSTM models performed better in forecasting than the ARIMA. Furthermore, it is possible that the results would be different if different datasets were used in the forecasts. Forecasts for one dataset is presented in this study, but if these same methods were applied on other disease data or other time series data, other conclusions might be drawn. A conclusion can be drawn on the effectiveness of the two presented models though, where the LSTM neural network consistently seems to outperform the ARIMA. However, considering that none of the models make highly accurate forecasts, using either one to try and turn a assess changes after an intervention seems to be a futile endeavor.

References

- Arabi Belaghi, R., Beyene, J., & McDonald, S. D. (2021). Prediction of pre-term birth in nulliparous women using logistic regression and machine learning. *PLoS One*, *16*(6), e0252025. <https://doi.org/10.1371/journal.pone.0252025>
- Belongia, E. A., Simpson, M. D., King, J. P., Sundaram, M. E., Kelley, N. S., Osterholm, M. T., & McLean, H. Q. (2016). Variable influenza vaccine effectiveness by subtype: A systematic review and meta-analysis of test-negative design studies. *The Lancet. Infectious Diseases*, *16*(8), 942–951. [https://doi.org/10.1016/S1473-3099\(16\)00129-8](https://doi.org/10.1016/S1473-3099(16)00129-8)
- Bengio, Y., Simard, P., & Frasconi, P. (1994). Learning long-term dependencies with gradient descent is difficult. *IEEE Transactions on Neural Networks*, *5*(2), 157–166. <https://doi.org/10.1109/72.279181>
- Bernal, J. L., Cummins, S., & Gasparrini, A. (2017). Interrupted time series regression for the evaluation of public health interventions: A tutorial. *International Journal of Epidemiology*, *46*(1), 348–355. <https://doi.org/10.1093/ije/dyw098>
- Brousseau, N., Green, H., Andrews, N., Pryse, R., Baguelin, M., Postema, A., Ellis, J., & Pebody, R. (2015). Impact of influenza vaccination on respiratory illness rates in children attending private boarding schools in England, 2013–2014: A cohort study. *Epidemiology and Infection*, *143*, 1–11. <https://doi.org/10.1017/S0950268815000667>
- Campbell, A. P., Ogokeh, C., Weinberg, G. A., Boom, J. A., Englund, J. A., Williams, J. V., Halasa, N. B., Selvarangan, R., Staat, M. A., Klein, E. J., McNeal, M., Michaels, M. G., Sahni, L. C., Stewart, L. S., Szilagyi, P. G., Harrison, C. J., Lively, J. Y., Rha, B., Patel, M., & New Vaccine Surveillance Network (NVSN). (2021). Effect of Vaccination on Preventing Influenza-Associated Hospitalizations Among Children During a Severe Season Associated With B/Victoria Viruses, 2019–2020. *Clinical Infectious Diseases*, *73*(4), e947–e954. <https://doi.org/10.1093/cid/ciab060>
- CDC. (2020, October 6). *Benefits of Flu Vaccination During 2018-2019 Flu Season*. Centers for Disease Control and Prevention. <https://www.cdc.gov/flu/about/burden-averted/2019-2020.htm>
- CDC. (2022, October 3). *Flu Symptoms & Complications*. Centers for Disease Control and Prevention. <https://www.cdc.gov/flu/symptoms/symptoms.htm>
- Cho, K., van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., & Bengio, Y.

- (2014). *Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation* (arXiv:1406.1078). arXiv. <https://doi.org/10.48550/arXiv.1406.1078>
- Clemens, J., Brenner, R., Rao, M., Tafari, N., & Lowe, C. (1996). Evaluating New Vaccines for Developing Countries: Efficacy or Effectiveness? *JAMA*, 275(5), 390–397. <https://doi.org/10.1001/jama.1996.03530290060038>
- Collen, M. F. (1986). Origins of medical informatics. *The Western Journal of Medicine*, 145(6), 778–785.
- Cruz, J. A., & Wishart, D. S. (2007). Applications of machine learning in cancer prediction and prognosis. *Cancer Informatics*, 2, 59–77.
- Dawa, J. A., Chaves, S. S., Nyawanda, B., Njuguna, H. N., Makokha, C., Otieno, N. A., Anzala, O., Widdowson, M.-A., & Emukule, G. O. (2018). National burden of hospitalized and non-hospitalized influenza-associated severe acute respiratory illness in Kenya, 2012-2014. *Influenza and Other Respiratory Viruses*, 12(1), 30–37. <https://doi.org/10.1111/irv.12488>
- de Hoog, M. L. A., Venekamp, R. P., Damoiseaux, R. A. M. J., Schilder, A. G. M., Sanders, E. A. M., Smit, H. A., & Bruijning-Verhagen, P. C. J. L. (2019). Impact of Repeated Influenza Immunization on Respiratory Illness in Children With Pre-existing Medical Conditions. *Annals of Family Medicine*, 17(1), 7–13. <https://doi.org/10.1370/afm.2340>
- Dennis, J., Ramsay, T., Turgeon, A. F., & Zarychanski, R. (2013). Helmet legislation and admissions to hospital for cycling related head injuries in Canadian provinces and territories: Interrupted time series analysis. *BMJ (Clinical Research Ed.)*, 346, f2674. <https://doi.org/10.1136/bmj.f2674>
- Ellis, K., Godbole, S., Marshall, S., Lanckriet, G., Staudenmayer, J., & Kerr, J. (2014). Identifying Active Travel Behaviors in Challenging Environments Using GPS, Accelerometers, and Machine Learning Algorithms. *Frontiers in Public Health*, 2. <https://www.frontiersin.org/articles/10.3389/fpubh.2014.00036>
- Elman, J. L. (1991). Distributed representations, simple recurrent networks, and grammatical structure. *Machine Learning*, 7(2), 195–225. <https://doi.org/10.1007/BF00114844>
- Emukule, G. O. (n.d.). *The epidemiological and economic burden of influenza in Kenya: Implications for public health action*. 182.
- Emukule, G. O., Paget, J., van der Velden, K., & Mott, J. A. (2015). Influenza-Associated Disease Burden in Kenya: A Systematic Review of Literature. *PLoS ONE*, 10(9), e0138708. <https://doi.org/10.1371/journal.pone.0138708>
- Esteva, A., Kuprel, B., Novoa, R. A., Ko, J., Swetter, S. M., Blau, H. M., & Thrun, S. (2017). Dermatologist-level classification of skin cancer with deep neural networks. *Nature*, 542(7639), 48

115–118. <https://doi.org/10.1038/nature21056>

Ferdinands, J. M., Olsho, L. E. W., Agan, A. A., Bhat, N., Sullivan, R. M., Hall, M., Mourani, P. M., Thompson, M., Randolph, A. G., & Pediatric Acute Lung Injury and Sepsis Investigators (PALISI) Network. (2014). Effectiveness of influenza vaccine against life-threatening RT-PCR-confirmed influenza illness in US children, 2010-2012. *The Journal of Infectious Diseases*, *210*(5), 674–683. <https://doi.org/10.1093/infdis/jiu185>

Fretheim, A., & Tomic, O. (2015). Statistical process control and interrupted time series: A golden opportunity for impact evaluation in quality improvement. *BMJ Quality & Safety*, *24*(12), 748–752. <https://doi.org/10.1136/bmjqs-2014-003756>

Gers, F. A., Schmidhuber, J., & Cummins, F. (2000). Learning to Forget: Continual Prediction with LSTM.

Neural Computation, *12*(10), 2451–2471. <https://doi.org/10.1162/089976600300015015>

Grijalva, C. G., Nuorti, J. P., Arbogast, P. G., Martin, S. W., Edwards, K. M., & Griffin, M. R. (2007).

Decline in pneumonia admissions after routine childhood immunisation with pneumococcal conjugate vaccine in the USA: A time-series analysis. *Lancet (London, England)*, *369*(9568), 1179–1186. [https://doi.org/10.1016/S0140-6736\(07\)61339-4](https://doi.org/10.1016/S0140-6736(07)61339-4)

Retrieved March 30, 2024, from <https://otexts.com/fpp2/arma-r.html>

Aiken, E. L., Nguyen, A. T., Viboud, C., & Santillana, M. (2021). Toward the use of neural networks for influenza prediction at multiple spatial resolutions. *Science Advances*, *7*(25), eabb1237.

<https://doi.org/10.1126/sciadv.abb1237>

Applied Deep Learning with TensorFlow 2. (n.d.). Retrieved April 5, 2024, from

<https://link.springer.com/book/10.1007/978-1-4842-8020-1>

Barone-Adesi, F., Gasparrini, A., Vizzini, L., Merletti, F., & Richiardi, L. (2011). Effects of Italian smoking regulation on rates of hospital admission for acute coronary events: A country-wide study. *PloS One*, *6*(3), e17419. <https://doi.org/10.1371/journal.pone.0017419>

Bassat, Q., Blau, D. M., Ogbuanu, I. U., Samura, S., Kaluma, E., Bassey, I.-A., Sow, S., Keita, A. M., Tapia, M. D., Mehta, A., Kotloff, K. L., Rahman, A., Islam, K. M., Alam, M., El Arifeen, S., Gurley, E. S., Baillie, V., Mutevedzi, P., Mahtab, S., ... Child Health and Mortality Prevention Surveillance (CHAMPS) Network. (2023). Causes of Death Among Infants and Children in the

- Child Health and Mortality Prevention Surveillance (CHAMPS) Network. *JAMA Network Open*, 6(7), e2322494. <https://doi.org/10.1001/jamanetworkopen.2023.22494>
- Boccalini, S., Varone, O., Chellini, M., Pieri, L., Sala, A., Berardi, C., Bonanni, P., & Bechini, A. (2017). Hospitalizations for pneumonia, invasive diseases and otitis in Tuscany (Italy), 2002-2014: Which was the impact of universal pneumococcal pediatric vaccination? *Human Vaccines & Immunotherapeutics*, 13(2), 428–434. <https://doi.org/10.1080/21645515.2017.1264796>
- Boddington, N. L., Pearson, I., Whitaker, H., Mangtani, P., & Pebody, R. G. (2021). Effectiveness of Influenza Vaccination in Preventing Hospitalization Due to Influenza in Children: A Systematic Review and Meta-analysis. *Clinical Infectious Diseases*, 73(9), 1722–1732. <https://doi.org/10.1093/cid/ciab270>
- Chicco, D., Warrens, M. J., & Jurman, G. (2021). The coefficient of determination R-squared is more informative than SMAPE, MAE, MAPE, MSE and RMSE in regression analysis evaluation. *PeerJ Computer Science*, 7, e623. <https://doi.org/10.7717/peerj-cs.623>
- Chimmula, V. K. R., & Zhang, L. (2020). Time series forecasting of COVID-19 transmission in Canada using LSTM networks. *Chaos, Solitons, and Fractals*, 135, 109864. <https://doi.org/10.1016/j.chaos.2020.109864>
- Coronavirus*. (n.d.). Retrieved February 13, 2024, from <https://www.who.int/health-topics/coronavirus>
- Csákvári, T., Elmer, D., Németh, N., Komáromy, M., Mihály-Vajda, R., & Boncz, I. (2023). Assessing the impact of Hungary's public health product tax: An interrupted time series analysis. *Central European Journal of Public Health*, 31(1), 43–49. <https://doi.org/10.21101/cejph.a7284>
- Dawa, J., Emukule, G., Barasa, E., Widdowson, M.-A., Anzala, O., Leeuwen, E., Baguelin, M., Chaves, S., & Eggo, R. (2020). Seasonal influenza vaccination in Kenya: An economic evaluation using dynamic transmission modelling. *BMC Medicine*, 18, 223. <https://doi.org/10.1186/s12916-020-01687-7>
- El-Amir, H., & Hamdy, M. (2020). A Tour Through the Deep Learning Pipeline. In H. El-Amir & M. Hamdy (Eds.), *Deep Learning Pipeline: Building a Deep Learning Model with TensorFlow* (pp. 57–

- 84). Apress. https://doi.org/10.1007/978-1-4842-5349-6_3
- Han, C. (2022). *Cross-validation for autoregressive models*. [University of Louisville].
<https://doi.org/10.18297/etd/3958>
- Hariton, E., & Locascio, J. J. (2018). Randomised controlled trials - the gold standard for effectiveness research: Study design: randomised controlled trials. *BJOG: An International Journal of Obstetrics and Gynaecology*, *125*(13), 1716. <https://doi.org/10.1111/1471-0528.15199>
- HealthLeaders. (n.d.). *COVID-19: Oregon Hospitals Share Data, Create Real-Time Bed Capacity System*. Retrieved February 16, 2024, from <https://www.healthleadersmedia.com/innovation/covid-19-oregon-hospitals-share-data-create-real-time-bed-capacity-system>
- Hodson, T. O. (2022). Root-mean-square error (RMSE) or mean absolute error (MAE): When to use them or not. *Geoscientific Model Development*, *15*(14), 5481–5487. <https://doi.org/10.5194/gmd-15-5481-2022>
- Hood, N., Flannery, B., Gaglani, M., Beeram, M., Wernli, K., Jackson, M. L., Martin, E. T., Monto, A. S., Zimmerman, R., Raviotta, J., Belongia, E. A., McLean, H. Q., Kim, S., Patel, M. M., & Chung, J. R. (2023). Influenza Vaccine Effectiveness Among Children: 2011–2020. *Pediatrics*, *151*(4), e2022059922. <https://doi.org/10.1542/peds.2022-059922>
- Iaousse, M., Jouilil, Y., Bouincha, M., & Mentagui, D. (2023). A Comparative Simulation Study of Classical and Machine Learning Techniques for Forecasting Time Series Data. *International Journal of Online and Biomedical Engineering (iJOE)*, *19*(08), Article 08.
<https://doi.org/10.3991/ijoe.v19i08.39853>
- Iyamu, I., Pedersen, H., Ablona, A., Chang, H.-J., Worthington, C., Grace, D., Grennan, T., Wong, J., Salmon, A., Koehoorn, M., & Gilbert, M. (2023). Evaluating the Impact of the COVID-19–Related Public Health Restrictions on Access to Digital Sexually Transmitted and Blood-Borne Infection Testing in British Columbia, Canada: An Interrupted Time Series Analysis. *Sexually Transmitted Diseases*, *50*(9), 595. <https://doi.org/10.1097/OLQ.0000000000001833>
- Kahenda, M. (n.d.). *Helpless mothers cry for answers as deadly influenza strikes hard*. Health. Retrieved

- April 5, 2024, from <https://www.standardmedia.co.ke/health/health-science/article/2000198402/helpless-mothers-cry-for-answers-as-deadly-influenza-strikes-hard>
- Kna. (2022, August 24). Kenya: Govt Contains Swine-Flu Outbreak in Nakuru. *Capital FM*.
<https://allafrica.com/stories/202208240075.html>
- Korevaar, E., Turner, S. L., Forbes, A. B., Karahalios, A., Taljaard, M., & McKenzie, J. E. (2024). Comparison of statistical methods used to meta-analyse results from interrupted time series studies: An empirical study. *BMC Medical Research Methodology*, 24, 31. <https://doi.org/10.1186/s12874-024-02147-z>
- Kumar, V. (2023). Crime Data Analysis Using Machine Learning Models. In M. Botto-Tobar, M. Zambrano Vizueté, S. Montes León, P. Torres-Carrión, & B. Durakovic (Eds.), *Applied Technologies* (pp. 296–309). Springer Nature Switzerland. https://doi.org/10.1007/978-3-031-24985-3_22
- Lee, J.-T., Lin, J.-W., Chen, H.-M., Wang, C.-Y., Lu, C.-Y., Chang, L.-Y., & Huang, L.-M. (2022). Impact of pneumococcal conjugate vaccination on hospitalized childhood pneumonia in Taiwan. *Pediatric Research*, 92(4), 1161–1167. <https://doi.org/10.1038/s41390-021-01772-4>
- Lu, C. Y., Simon, G., & Soumerai, S. B. (2018). Counter-Point: Staying Honest When Policy Changes Backfire. *Medical Care*, 56(5), 384–390. <https://doi.org/10.1097/MLR.0000000000000897>
- Lv, K., Jiang, S., & Li, J. (2017). *Learning Gradient Descent: Better Generalization and Longer Horizons* (arXiv:1703.03633). arXiv. <https://doi.org/10.48550/arXiv.1703.03633>
- Malone, B., Simovski, B., Moliné, C., Cheng, J., Gheorghe, M., Fontenelle, H., Vardaxis, I., Tennøe, S., Malmberg, J.-A., Stratford, R., & Clancy, T. (2020). Artificial intelligence predicts the immunogenic landscape of SARS-CoV-2 leading to universal blueprints for vaccine designs. *Scientific Reports*, 10(1), 22375. <https://doi.org/10.1038/s41598-020-78758-5>
- Matranga, D., Bono, F., & Maniscalco, L. (2021). Statistical Advances in Epidemiology and Public Health. *International Journal of Environmental Research and Public Health*, 18(7), Article 7.
<https://doi.org/10.3390/ijerph18073549>
- Mealing, N., Hayen, A., & Newall, A. T. (2016). Assessing the impact of vaccination programmes on

- burden of disease: Underlying complexities and statistical methods. *Vaccine*, 34(27), 3022–3029.
<https://doi.org/10.1016/j.vaccine.2016.04.014>
- Morbey, R. A., Todkill, D., Watson, C., & Elliot, A. J. (2023). *Machine learning forecasts for seasonal epidemic peaks: Lessons learnt from an atypical respiratory syncytial virus season* (p. 2023.06.29.23291793). medRxiv. <https://doi.org/10.1101/2023.06.29.23291793>
- Mpox (monkeypox)*. (n.d.). Retrieved February 13, 2024, from <https://www.who.int/news-room/factsheets/detail/monkeypox>
- Naskar, I. (n.d.). System Modeling and Prediction of Respiratory Diseases using Machine Learning. *YMER Digital*. Retrieved March 1, 2024, from https://www.academia.edu/102716575/System_Modeling_and_Prediction_of_Respiratory_Diseases_using_Machine_Learning
- Newall, A. T., Reyes, J. F., Wood, J. G., McIntyre, P., Menzies, R., & Beutels, P. (2014). Economic evaluations of implemented vaccination programmes: Key methodological challenges in retrospective analyses. *Vaccine*, 32(7), 759–765. <https://doi.org/10.1016/j.vaccine.2013.11.067>
- Nyawanda, B. O., Murunga, N., Otieno, N. A., Bigogo, G., Nyiro, J. U., Vodicka, E., Bulterys, M., Nokes, D. J., Munywoki, P. K., & Emukule, G. O. (2023). Estimates of the national burden of respiratory syncytial virus in Kenyan children aged under 5 years, 2010–2018. *BMC Medicine*, 21(1), 122. <https://doi.org/10.1186/s12916-023-02787-w>
- Paulo, R. L. P., Rodrigues, A. B. D., Machado, B. M., & Gilio, A. E. (2016). The impact of rotavirus vaccination on emergency department visits and hospital admissions for acute diarrhea in children under 5 years. *Revista Da Associação Médica Brasileira*, 62(6), 506–512. <https://doi.org/10.1590/1806-9282.62.06.506>
- Reduction in Acute Gastroenteritis Hospitalizations among US Children After Introduction of Rotavirus Vaccine: Analysis of Hospital Discharge Data from 18 US States* | *The Journal of Infectious Diseases* | *Oxford Academic*. (n.d.). Retrieved February 28, 2024, from <https://academic.oup.com/jid/article/201/11/1617/850599>

- Reimers, N., & Gurevych, I. (2017). *Optimal Hyperparameters for Deep LSTM-Networks for Sequence Labeling Tasks* (arXiv:1707.06799). arXiv. <https://doi.org/10.48550/arXiv.1707.06799>
- Robeson, S. M., & Willmott, C. J. (2023). Decomposition of the mean absolute error (MAE) into systematic and unsystematic components. *PLOS ONE*, *18*(2), e0279774. <https://doi.org/10.1371/journal.pone.0279774>
- Ru, B., Kujawski, S., Afanador, N. L., Baumgartner, R., Pawaskar, M., & Das, A. (2023). Predicting Measles Outbreaks in the United States: Evaluation of Machine Learning Approaches. *JMIR Formative Research*, *7*(1), e42832. <https://doi.org/10.2196/42832>
- S, V., S, S., H, V., & S, S. (2018). *Disease Prediction Using Machine Learning Over Big Data* (SSRN Scholarly Paper 3458775). <https://doi.org/10.2139/ssrn.3458775>
- Schaffer, A. L., Dobbins, T. A., & Pearson, S.-A. (2021). Interrupted time series analysis using autoregressive integrated moving average (ARIMA) models: A guide for evaluating large-scale health interventions. *BMC Medical Research Methodology*, *21*(1), 58. <https://doi.org/10.1186/s12874-021-01235-8>
- Scopus preview - Scopus - Document details - Resolving the pneumococcal vaccine controversy: Are there alternatives to randomized clinical trials?* (n.d.). <https://doi.org/10.1093/clinids/6.5.589>
- Shakeel, S. M., Kumar, N. S., Madalli, P. P., Srinivasaiah, R., & Swamy, D. R. (2021). COVID-19 prediction models: A systematic literature review. *Osong Public Health and Research Perspectives*, *12*(4), 215–229. <https://doi.org/10.24171/j.phrp.2021.0100>
- Sharmila, L., Dharuman, C., & Venkatesan, P. (2017). Disease Classification Using Machine Learning Algorithms-A Comparative Study. *International Journal of Pure and Applied Mathematics*, *114*, 1–10.
- Shaw, J., Murphy, A. L., Turner, J. P., Gardner, D. M., Silvius, J. L., Bouck, Z., Gordon, D., & Tannenbaum, C. (2019). Policies for Deprescribing: An International Scan of Intended and Unintended Outcomes of Limiting Sedative-Hypnotic Use in Community-Dwelling Older Adults. *Healthcare Policy = Politiques De Sante*, *14*(4), 39–51. <https://doi.org/10.12927/hcpol.2019.25857>

- Soumerai, S. B., Starr, D., & Majumdar, S. R. (2015). How Do You Know Which Health Care Effectiveness Research You Can Trust? A Guide to Study Design for the Perplexed. *Preventing Chronic Disease, 12*, E101. <https://doi.org/10.5888/pcd12.150187>
- Subrahmanya, S. V. G., Shetty, D. K., Patil, V., Hameed, B. M. Z., Paul, R., Smriti, K., Naik, N., & Somani, B. K. (2022). The role of data science in healthcare advancements: Applications, benefits, and future prospects. *Irish Journal of Medical Science, 191*(4), 1473–1483. <https://doi.org/10.1007/s11845-021-02730-z>
- Sunil, S., Kumar, V. D., Babu, A., Thilak, G., & Udayan, D. (2022). Covid-19 spread Forecast with respect to vaccination based on LSTM and GRU. *Proceedings of the 2022 Fourteenth International Conference on Contemporary Computing, 32–36*. <https://doi.org/10.1145/3549206.3549213>
- Tsan, Y.-T., Chen, D.-Y., Liu, P.-Y., Kristiani, E., Nguyen, K. L. P., & Yang, C.-T. (2022). The Prediction of Influenza-like Illness and Respiratory Disease Using LSTM and ARIMA. *International Journal of Environmental Research and Public Health, 19*(3), Article 3. <https://doi.org/10.3390/ijerph19031858>
- Turner, S. L., Forbes, A. B., Karahalios, A., Taljaard, M., & McKenzie, J. E. (2021). Evaluation of statistical methods used in the analysis of interrupted time series studies: A simulation study. *BMC Medical Research Methodology, 21*(1), 181. <https://doi.org/10.1186/s12874-021-01364-0>
- Wang, H., Kwok, K. O., & Riley, S. (2023). *Forecasting influenza incidence as an ordinal variable using machine learning* (p. 2023.02.09.23285705). medRxiv. <https://doi.org/10.1101/2023.02.09.23285705>
- Williams, J., Ahlqvist, H., Cunningham, A., Kirby, A., Katz, I., Fleming, J., Conway, J., Cunningham, S., Ozel, A., & Wolfram, U. (2023). *Validated respiratory drug deposition predictions from 2D and 3D medical images with statistical shape models and convolutional neural networks* (arXiv:2303.01036). arXiv. <http://arxiv.org/abs/2303.01036>
- Yang, J., Xu, X., Ma, X., Wang, Z., You, Q., Shan, W., Yang, Y., Bo, X., & Yin, C. (2023). Application of machine learning to predict hospital visits for respiratory diseases using meteorological and air

- pollution factors in Linyi, China. *Environmental Science and Pollution Research International*, 30(38), 88431–88443. <https://doi.org/10.1007/s11356-023-28682-8>
- Yuen, A., Pourmarzi, D., Sarkis, S., Luisetto, C., Khatri, K., Bone, A., & Black, J. (2023). Interrupted time series segmented regression analysis for detecting waterborne disease outbreaks by syndromic surveillance. *Communicable Diseases Intelligence*, 47. <https://doi.org/10.33321/cdi.2023.47.5>
- Zhang, H., Yi, D., & Guan, Y. (2021). Timesias: A machine learning pipeline for predicting outcomes from time-series clinical records. *STAR Protocols*, 2(3), 100639. <https://doi.org/10.1016/j.xpro.2021.100639>
- Zhang, Y., & Meng, G. (2023). Simulation of an Adaptive Model Based on AIC and BIC ARIMA Predictions. *Journal of Physics: Conference Series*, 2449(1), 012027. <https://doi.org/10.1088/1742-6596/2449/1/012027>
- Zhu, Y. (2022). Prediction Modeling and Research on the Relationship Between Urban Air Pollutants and Respiratory Diseases. *Academic Journal of Science and Technology*, 2(3), Article 3. <https://doi.org/10.54097/ajst.v2i3.1531>
- [org/10.1016/S0140-6736\(07\)60564-9](https://doi.org/10.1016/S0140-6736(07)60564-9)
- Habebh, H., & Gohel, S. (2021). Machine Learning in Healthcare. *Current Genomics*, 22(4), 291–300. <https://doi.org/10.2174/1389202922666210705124359>
- Hasman, A., Mantas, J., & Zarubina, T. (2014). An abridged history of medical informatics education in europe. *Acta Informatica Medica: AIM: Journal of the Society for Medical Informatics of Bosnia & Herzegovina: Casopis Drustva Za Medicinsku Informatiku BiH*, 22(1), 25–36. <https://doi.org/10.5455/aim.2014.22.25-36>
- Influenza (Seasonal)*. (n.d.). Retrieved November 29, 2022, from [https://www.who.int/news-room/fact-sheets/detail/influenza-\(seasonal\)](https://www.who.int/news-room/fact-sheets/detail/influenza-(seasonal))
- Katz, M. A., Lebo, E., Emukule, G. O., Otieno, N., Caselton, D. L., Bigogo, G., Njuguna, H., Muthoka, P. M., Waiboci, L. W., Widdowson, M.-A., Xu, X., Njenga, M. K., Mott, J. A., & Breiman, R. F. (2016). Uptake and Effectiveness of a Trivalent Inactivated Influenza Vaccine in Children in Urban and Rural Kenya, 2010 to 2012. *The Pediatric Infectious Disease Journal*, 35(3), 322–329. <https://doi.org/10.1097/INF.0000000000001035>

- Li, L., Jiang, Y., & Huang, B. (2021). Long-term prediction for temporal propagation of seasonal influenza using Transformer-based model. *Journal of Biomedical Informatics*, *122*, 103894. <https://doi.org/10.1016/j.jbi.2021.103894>
- Linden, A. (2015). Conducting Interrupted Time-series Analysis for Single- and Multiple-group Comparisons. *The Stata Journal*, *15*(2), 480–500. <https://doi.org/10.1177/1536867X1501500208>
- Lopez Bernal, J. A., Andrews, N., & Amirthalingam, G. (2019). The Use of Quasi-experimental Designs for Vaccine Evaluation. *Clinical Infectious Diseases*, *68*(10), 1769–1776. <https://doi.org/10.1093/cid/ciy906>
- Mameli, C., Cocchi, I., Fumagalli, M., & Zuccotti, G. (2019). Influenza Vaccination: Effectiveness, Indications, and Limits in the Pediatric Population. *Frontiers in Pediatrics*, *7*, 317. <https://doi.org/10.3389/fped.2019.00317>
- Matheka, D. M., Mokaya, J., & Maritim, M. (2013). Overview of influenza virus infections in Kenya: Past, present and future. *The Pan African Medical Journal*, *14*, 138. <https://doi.org/10.11604/pamj.2013.14.138.2612>
- Mw, T., Hk, T., Ch, T., M, G., Tm, M., As, M., Et, M., Rk, Z., Fp, S., Db, M., Sm, O., Rj, G. K., Jr, B., Jm, F., Mm, P., & undefined. (2021). Influenza Vaccine Effectiveness Against Hospitalization in the United States, 2019-2020. *The Journal of Infectious Diseases*, *224*(5), 813–820. <https://doi.org/10.1093/infdis/jiaa800>
- Nunes, M. C., Cutland, C. L., Jones, S., Downs, S., Weinberg, A., Ortiz, J. R., Neuzil, K. M., Simões, E. A. F., Klugman, K. P., & Madhi, S. A. (2017). Efficacy of Maternal Influenza Vaccination Against All- Cause Lower Respiratory Tract Infection Hospitalizations in Young Infants: Results From a Randomized Controlled Trial. *Clinical Infectious Diseases*, *65*(7), 1066–1071. <https://doi.org/10.1093/cid/cix497>
- Omrani, H. (2015). Predicting Travel Mode of Individuals by Machine Learning. *Transportation Research Procedia*, *10*, 840–849. <https://doi.org/10.1016/j.trpro.2015.09.037>
- Otieno, N. A., Nyawanda, B. O., McMorrow, M., Oneko, M., Omollo, D., Lidechi, S., Widdowson, M., Flannery, B., Chaves, S. S., Azziz-Baumgartner, E., & Emukule, G. O. (2022). The burden of influenza among Kenyan pregnant and postpartum women and their infants, 2015–2020. *Influenza and Other Respiratory Viruses*, *16*(3), 452–461. <https://doi.org/10.1111/irv.12950>
- Penfold, R. B., & Zhang, F. (2013). Use of interrupted time series analysis in evaluating health care quality improvements. *Academic Pediatrics*, *13*(6 Suppl), S38-44. <https://doi.org/10.1016/j.acap.2013.08.002>

- Perin, J., Mulick, A., Yeung, D., Villavicencio, F., Lopez, G., Strong, K. L., Prieto-Merino, D., Cousens, S., Black, R. E., & Liu, L. (2022). Global, regional, and national causes of under-5 mortality in 2000–19: An updated systematic analysis with implications for the Sustainable Development Goals. *The Lancet Child & Adolescent Health*, 6(2), 106–115. [https://doi.org/10.1016/S2352-4642\(21\)00311-4](https://doi.org/10.1016/S2352-4642(21)00311-4)
- Rajpurkar, P., Irvin, J., Zhu, K., Yang, B., Mehta, H., Duan, T., Ding, D., Bagul, A., Langlotz, C., Shpanskaya, K., Lungren, M. P., & Ng, A. Y. (2017). CheXNet: Radiologist-Level Pneumonia Detection on Chest X-Rays with Deep Learning. In *ArXiv e-prints*. <https://ui.adsabs.harvard.edu/abs/2017arXiv171105225R>
- Recommendations for influenza vaccine composition*. (n.d.). Retrieved November 1, 2022, from <https://www.who.int/teams/global-influenza-programme/vaccines/who-recommendations>
- Robot-assisted surgery for benign distal ureteral strictures: Step-by-step technique using the SP® surgical system—Kaouk—2019—BJU International—Wiley Online Library*. (n.d.). Retrieved December 7, 2022, from <https://bjui-journals.onlinelibrary.wiley.com/doi/abs/10.1111/bju.14635>
- Rose, A., Kissling, E., Emborg, H.-D., Larrauri, A., McMenamin, J., Pozo, F., Trebbien, R., Mazagatos, C., Whitaker, H., & Valenciano, M. (2020). Interim 2019/20 influenza vaccine effectiveness: Six European studies, September 2019 to January 2020. *Eurosurveillance*, 25(10), 2000153. <https://doi.org/10.2807/1560-7917.ES.2020.25.10.2000153>
- Sambala, E. Z., Cooper, S., Schmidt, B.-M., Walaza, S., & Wiysonge, C. S. (2021). Role of vaccines in preventing influenza in healthy children. *SAMJ: South African Medical Journal*, 111(3), 206–207. <https://doi.org/10.7196/samj.2021.v111i3.14478>
- Segaloff, H. E., Leventer-Roberts, M., Riesel, D., Malosh, R. E., Feldman, B. S., Shemer-Avni, Y., Key, C., Monto, A. S., Martin, E. T., & Katz, M. A. (2019). Influenza Vaccine Effectiveness Against Hospitalization in Fully and Partially Vaccinated Children in Israel: 2015-2016, 2016-2017, and 2017-2018. *Clinical Infectious Diseases: An Official Publication of the Infectious Diseases Society of America*, 69(12), 2153–2161. <https://doi.org/10.1093/cid/ciz125>
- Shinjoh, M., Sugaya, N., Furuichi, M., Araki, E., Maeda, N., Isshiki, K., Ohnishi, T., Nakamura, S., Yamada, G., Narabayashi, A., Nishida, M., Taguchi, N., Nakata, Y., Yoshida, M., Tsunematsu, K., Shibata, M., Munenaga, T., Hirano, Y., Ookawara, I., ... Keio Pediatric Influenza Research Group. (2019). Effectiveness of inactivated influenza vaccine in children by vaccine dose, 2013-18. *Vaccine*, 37(30), 4047–4054. <https://doi.org/10.1016/j.vaccine.2019.05.090>
- Siddiqui, M. K., Morales-Menendez, R., Huang, X., & Hussain, N. (2020). A review of epileptic seizure detection using machine learning classifiers. *Brain Informatics*, 7(1), 5. <https://doi.org/10.1186/s40708-020-00105-1>
- Sutskever, I., Martens, J., & Hinton, G. (2011). Generating text with recurrent neural networks.

Proceedings of the 28th International Conference on International Conference on Machine Learning, 1017–1024.

Tang, J., Liu, R., Zhang, Y.-L., Liu, M.-Z., Hu, Y.-F., Shao, M.-J., Zhu, L.-J., Xin, H.-W., Feng, G.-W., Shang,

W.-J., Meng, X.-G., Zhang, L.-R., Ming, Y.-Z., & Zhang, W. (2017). Application of Machine-Learning Models to Predict Tacrolimus Stable Dose in Renal Transplant Recipients. *Scientific Reports*, 7(1), Article 1. <https://doi.org/10.1038/srep42192>

Tenny, S., Kerndt, C. C., & Hoffman, M. R. (2022). Case Control Studies. In *StatPearls*. StatPearls Publishing. <http://www.ncbi.nlm.nih.gov/books/NBK448143/>

Thakkar, M., & Davis, D. C. (2006). Risks, barriers, and benefits of EHR systems: A comparative study based on size of hospital. *Perspectives in Health Information Management*, 3, 5.

Thompson, M. G., Kwong, J. C., Regan, A. K., Katz, M. A., Drews, S. J., Azziz-Baumgartner, E., Klein, N. P.,

Chung, H., Effler, P. V., Feldman, B. S., Simmonds, K., Wyant, B. E., Dawood, F. S., Jackson, M. L.,

Fell, D. B., Levy, A., Barda, N., Svenson, L. W., Fink, R. V., ... PREVENT Workgroup. (2019). Influenza Vaccine Effectiveness in Preventing Influenza-associated Hospitalizations During Pregnancy: A Multi-country Retrospective Test Negative Design Study, 2010-2016. *Clinical Infectious Diseases: An Official Publication of the Infectious Diseases Society of America*, 68(9), 1444–1453. <https://doi.org/10.1093/cid/ciy737>

Turing 1950 Article - A. M. Turing (1950) Computing Machinery and Intelligence. Mind 49: 433-460. - Studocu. (n.d.). Retrieved December 7, 2022, from <https://www.studocu.com/en-us/document/emory-university/cognition/turing-1950-article/13121870>

Vaccine Effectiveness: How Well Do Flu Vaccines Work? | CDC. (2022, August 25). <https://www.cdc.gov/flu/vaccines-work/vaccineeffect.htm>

Verani, J. R., Baqui, A. H., Broome, C. V., Cherian, T., Cohen, C., Farrar, J. L., Feikin, D. R., Groome, M. J.,

Hajjeh, R. A., Johnson, H. L., Madhi, S. A., Mulholland, K., O'Brien, K. L., Parashar, U. D., Patel, M. M., Rodrigues, L. C., Santosham, M., Scott, J. A., Smith, P. G., ... Zell, E. R. (2017). Case-control vaccine effectiveness studies: Data collection, analysis and reporting results. *Vaccine*, 35(25), 3303–3308. <https://doi.org/10.1016/j.vaccine.2017.04.035>

Wati, D. A. R., & Abadianto, D. (2017). Design of face detection and recognition system for smart home security application. *2017 2nd International Conferences on Information Technology, Information Systems and Electrical Engineering (ICITISEE)*, 342–347. <https://doi.org/10.1109/ICITISEE.2017.8285524>

What Is Machine Learning (ML)? (2020, June 26). *UCB-UMT*. <https://ischoolonline.berkeley.edu/blog/what-is-machine-learning/>

Woldaregay, A. Z., Årsand, E., Botsis, T., Albers, D., Mamykina, L., & Hartvigsen, G. (2019). Data-Driven Blood Glucose Pattern Classification and Anomalies Detection: Machine-Learning Applications in Type 1 Diabetes. *Journal of Medical Internet Research*, 21(5), e11030. <https://doi.org/10.2196/11030>



Interrupted Time Series and Machine Learning with Application to Effect of Influenza Vaccine

By
Cynthia Ombok Juma
144892





Final Decision

This is to certify that the application for ethics clearance submitted by:

Principal Investigator: Ms. Juma, Cynthia Ombok

Reference number: SU-ISERC1639/23

For Study: "Effect of Influenza Vaccine Administration: Comparative approach, Interrupted Times Series and Machine Learning"

Was reviewed and received the following status: "done"


Reviewer Comments

Final decision: **approved**

Comments sent:

The SU-ISERC wishes you all the best with this research undertaking.

19 April 2023 11:45:14



Cynthia Ombok Juma,
PO Box 72938-00200,
Nairobi, Kenya.

4th May 2023.

The Principal Secretary,
Ministry of Health,
State Department for Public Health and Professional Standards,
Afya House, Cathedral Road,
P.O. Box: 30016 - 00100,
Nairobi, Kenya.



D. Dr. Hellen Kiame
FYA
Amma
4/5/2023

Dear Madam,

RE: REQUEST FOR DHIS2 DATA FROM THE MINISTRY OF HEALTH, KENYA

I am writing to request for data from Ministry of Health Kenya District Health Information System (DHIS2) repository for my studies. I plan to use the data in my study titled "*Effect of Influenza Vaccine Administration: Comparative Approach, Interrupted Times Series and Machine Learning*" in partial fulfillment of my MSc Data Science and Analytics degree training requirement at Strathmore University.

Human influenza infections are a major cause of morbidity and mortality worldwide. In Kenya data on the burden of influenza disease are needed to inform influenza control policies. As part of generating data. In 2016 the Kenya National Immunization Technical Advisory Group (KENITAG) undertook a detailed review of influenza data in Kenya and recommended influenza vaccination for children aged 6-23 months as they have the highest number of cases with severe influenza disease. However, given the little experience of influenza vaccination in the public sector in Kenya, KENITAG recommended that a pilot influenza vaccination project is first implemented to generate additional data (including an assessment of the impact on cold chain management, waste management, optimal vaccine delivery strategies, and vaccine uptake) that could potentially inform a future national vaccination program in Kenya. Subsequently, the MOH with the support from partners conducted an influenza vaccination pilot demonstration project among children aged 6-23 months based on KENITAG's recommendations in selected sites within Mombasa and Nakuru Counties in 2018-2021.

There is a need to rigorously evaluate the impact of influenza vaccine administration in Kenya to document its effectiveness. I am therefore requesting for DHIS2 data on the number of respiratory illnesses (including Pneumonia, Asthma, URTI, and LRTI) among aged children >5 years for the period of 1st January 2019 to 31st December 2022 for Mombasa and Nakuru Counties. This planned study aims to compare the effectiveness of two methods for evaluating the impact of influenza vaccine administration in Kenya, that is, an Interrupted Time Series (ITS) design and a Recurrent Neural Network (RNN) design. The findings will inform the development of more effective methods for evaluating the impact of health interventions in similar settings, which will ultimately contribute to improving the health outcomes of vulnerable populations. The study has already been approved by Strathmore University and the National Commission for Science, Technology, and Innovation (NACOSTI).

I look forward to your kind consideration. Thank you.

Sincerely,

Cynthia Ombok Juma



REPUBLIC OF KENYA



NATIONAL COMMISSION FOR SCIENCE, TECHNOLOGY & INNOVATION

Ref No: 799926

Date of Issue: 29/April/2023

RESEARCH LICENSE



This is to Certify that Ms. Cynthia Ombok Juma of Strathmore University, has been licensed to conduct research as per the provision of the Science, Technology and Innovation Act, 2013 (Rev.2014) in Mombasa, Nakuru on the topic: Effect of Influenza Vaccine Administration: Comparative approach, Interrupted Times Series and Machine Learning for the period ending : 29/April/2024.

License No: NACOSTI/P/23/25489

799926

Applicant Identification Number

Walter Mwangi

Director General

NATIONAL COMMISSION FOR SCIENCE, TECHNOLOGY & INNOVATION

Verification QR Code



NOTE: This is a computer generated License. To verify the authenticity of this document, Scan the QR Code using QR scanner application.

See overleaf for conditions

THE SCIENCE, TECHNOLOGY AND INNOVATION ACT, 2013 (Rev. 2014)
Legal Notice No. 108: The Science, Technology and Innovation (Research Licensing) Regulations, 2014

The National Commission for Science, Technology and Innovation, hereafter referred to as the Commission, was established under the Science, Technology and Innovation Act 2013 (Revised 2014) herein after referred to as the Act. The objective of the Commission shall be to regulate and assure quality in the science, technology and innovation sector and advise the Government in matters related thereto.

CONDITIONS OF THE RESEARCH LICENSE

1. The License is granted subject to provisions of the Constitution of Kenya, the Science, Technology and Innovation Act, and other relevant laws, policies and regulations. Accordingly, the licensee shall adhere to such procedures, standards, code of ethics and guidelines as may be prescribed by regulations made under the Act, or prescribed by provisions of International treaties of which Kenya is a signatory to
2. The research and its related activities as well as outcomes shall be beneficial to the country and shall not in any way;
 - i. Endanger national security
 - ii. Adversely affect the lives of Kenyans
 - iii. Be in contravention of Kenya's international obligations including Biological Weapons Convention (BWC), Comprehensive Nuclear-Test-Ban Treaty Organization (CTBTO), Chemical, Biological, Radiological and Nuclear (CBRN).
 - iv. Result in exploitation of intellectual property rights of communities in Kenya
 - v. Adversely affect the environment
 - vi. Adversely affect the rights of communities
 - vii. Endanger public safety and national cohesion
 - viii. Plagiarize someone else's work
3. The License is valid for the proposed research, location and specified period.
4. The license any rights thereunder are non-transferable
5. The Commission reserves the right to cancel the research at any time during the research period if in the opinion of the Commission the research is not implemented in conformity with the provisions of the Act or any other written law.
6. The Licensee shall inform the relevant County Director of Education, County Commissioner and County Governor before commencement of the research.
7. Excavation, filming, movement, and collection of specimens are subject to further necessary clearance from relevant Government Agencies.
8. The License does not give authority to transfer research materials.
9. The Commission may monitor and evaluate the licensed research project for the purpose of assessing and evaluating compliance with the conditions of the License.
10. The Licensee shall submit one hard copy, and upload a soft copy of their final report (thesis) onto a platform designated by the Commission within one year of completion of the research.
11. The Commission reserves the right to modify the conditions of the License including cancellation without prior notice.
12. Research, findings and information regarding research systems shall be stored or disseminated, utilized or applied in such a manner as may be prescribed by the Commission from time to time.
13. The Licensee shall disclose to the Commission, the relevant Institutional Scientific and Ethical Review Committee, and the relevant national agencies any inventions and discoveries that are of National strategic importance.
14. The Commission shall have powers to acquire from any person the right in, or to, any scientific innovation, invention or patent of strategic importance to the country.
15. Relevant Institutional Scientific and Ethical Review Committee shall monitor and evaluate the research periodically, and make a report of its findings to the Commission for necessary action.

National Commission for Science, Technology and
Innovation(NACOSTI),
Off Waiyaki Way, Upper Kabete,
P. O. Box 30623 - 00100 Nairobi, KENYA
Telephone: 020 4007000, 0713788787, 073 5404245
E-mail: dg@nacosti.go.ke
Website: www.nacosti.go.ke