# Time Series Forecasting of NYC Crime

Nickolaus White
Michigan State University
1535 Woodbrook St
East Lansing, MI 48823
1-517-983-2237

whiten29@msu.edu

Stephen Lee
Michigan State University
136 Reniger Ct
East Lansing, MI, 48823
1-248-912-8263

leestep9@msu.edu

Mark Carravallah
Michigan State University
1850 Abbot Rd
East Lansing, MI 48823
1-248-761-9336

carrava8@msu.edu

## ABSTRACT

This project aimed to develop an interactive dashboard facilitating the selection of Airbnb accommodations within New York City. The dashboard integrates crime rate reports from 2014-2016 and forecasts future crime rates based on user-input parameters such as borough, crime type, time of day, and date. Additionally, the application allows users to predict crime rates 30 days in advance, offering a comparative analysis against actual reported data within the specified timeframe. The dashboard further incorporates an interactive map enabling users to filter and select Airbnbs based on criteria like price, borough, review rating, host identity verification, and instant bookability.

The utility and significance of this project lie in providing users with valuable insights into crime forecasts, enhancing their confidence in selecting safe vacation locations while optimizing their budget. The challenge of working with older data containing frequent NaN values was successfully addressed, and the transition to Streamlit posed another challenge, particularly in implementing multiple filters and an interactive map.

Key accomplishments include the creation of a comprehensive Jupyter Notebook documenting data collection, preprocessing, time series model creation, and results. The successful conversion of this documentation into a Streamlit application empowers users to explore the datasets independently.

In terms of results, the project compared ARIMA and SARIMA models with standardization, ultimately identifying SARIMA as the most effective model. The success of SARIMA in outperforming ARIMA can be attributed to its ability to account for seasonality and incorporate exogenous variables, providing a more nuanced and accurate prediction of crime rates in the specified context.

## Keywords
ARIMA, SARIMA, RMSE, Time Series, Crime Forecasting, NYC, Streamlit, Data Mining, Airbnb, Web Application

## 1. INTRODUCTION

We aim to create a useful tool for vacationers looking to stay in New York City but might have reservations about safety — Our solution is an interactive dashboard with the ability to specifically tailor one's specific wants for an Airbnb listing. Using reported crime data spanning 2014-2016 this app allows users a comprehensive tool to forecast and compare crime rates across specific boroughs in the city.

The fusion of crime rate data with Airbnb selection criteria responds to the intrinsic need of travelers to make informed decisions about their stay. Our goal is to create a single interface with an interactive map coupled with customizable filters so that they may confidently choose their upcoming vacation stay. By allowing users to predict and compare crime rates 30 days in advance, we provide a unique feature that enhances the decision-making process. Furthermore, the inclusion of an interactive map allows users a "one-stop-shop" website to choose Airbnbs based on filters such as price, borough, review rating, host identity verification, and instant bookability.

This project was not without its challenges. The original database did not have the necessary size to perform accurate forecasting and after the merging of two other databases the size and scope of the database required extensive data preprocessing. Additionally, transitioning our code to Streamlit, especially when incorporating multiple filters and an interactive map, proved to be a complex task.

In our exploration of predictive models, a notable finding emerged — SARIMA outshone ARIMA due to the seasonality of the data and therefore became our primary model for success.

The majority of our project was compiled in a Jupyter Notebook which covers the data collection, preprocessing, model creation, and results. Additionally we created a user-friendly Streamlit application in Python, allowing for users to explore datasets independently and with their own specifications.

## 2. RELATED WORK
This section describes the related work and the pros and cons found in similar projects.

## 2.1 OTHER RESEARCH
The development of interactive dashboards for decision-making in urban accommodation selection, especially in the context of integrating crime rate forecasts, has garnered attention from both the application domain and data mining communities. Our literature review reveals significant contributions that can be organized into several categories based on the nature of the related work.

1. Crime Rate Forecasting:

   - Gorr, W., Harries, R. (2003). "Introduction to Crime Forecasting." Provides foundational insights into crime forecasting techniques, emphasizing the importance of accurate predictions in enhancing public safety.

- Chainey, S., Tompson, L., & Uhlig, S. (2008). "The utility of hotspot mapping for predicting spatial patterns of crime." Explores spatial crime analysis and hotspot mapping, laying the groundwork for forecasting methodologies based on historical crime data.

2. Interactive Dashboards:

- Lu, X., Wylie, B., & Jiao, J. (2018). "Developing a GIS-Based Interactive Dashboard for Crime Analysis and Decision Support." Presents a GIS-based interactive dashboard for crime analysis, showcasing the potential of visual tools in aiding decision-making.

3. Data Mining Techniques:

- Fotheringham, A. S., Brunsdon, C., & Charlton, M. (2000). "Quantitative Geography: Perspectives on Spatial Data Analysis." Offers a comprehensive overview of spatial data analysis techniques, laying the groundwork for spatial modeling in urban applications.

4. Time Series Forecasting:

- Box, G. E. P., Jenkins, G. M., & Reinsel, G. C. (1994). "Time Series Analysis: Forecasting and Control." A seminal work on time series analysis, introducing concepts that underlie traditional forecasting methods like ARIMA.

- Hyndman, R. J., & Athanasopoulos, G. (2018). "Forecasting: principles and practice." A comprehensive resource on forecasting principles, providing insights into the application of SARIMA models and the importance of standardization.

5. Challenges in Data Mining:

- Krumm, J. (2008). "Inference attacks on location tracks." Explores privacy concerns in location-based data, shedding light on the challenges associated with handling sensitive information in urban decision support systems.

## 2.2 PROS AND CONS

- Existing crime rate forecasting models provide valuable insights but may lack adaptability to specific urban contexts.

- GIS-based dashboards offer visualizations but might not sufficiently address the integration of diverse datasets for decision support.

- Traditional time series forecasting methods like ARIMA may overlook the complexities of urban crime dynamics, while more advanced techniques like SARIMA offer improved predictive accuracy. Thoroughly studied and well received models for Time Series Forecasting are Recurrent Neural Networks (RNN) and Long Short-Term Memory (LSTM) RNNs. These models consistently provide the best accuracy for Time Series Forecasting, however we could not get these models to run properly.

This survey of related work provides a foundation for understanding the landscape of interactive urban decision support systems, integrating crime rate forecasts and Airbnb selection criteria. The identified gaps and insights from these studies guide our efforts in contributing to this evolving field.

## 3. PROBLEM STATEMENT

In the realm of selecting Airbnb accommodations within New York City, a pertinent challenge arises in providing users with a comprehensive and informed decision-making tool. The problem domain encompasses the integration of crime rate data spanning 2014 to 2016, enabling the development of an interactive dashboard. This dashboard not only reports historical crime rates but also forecasts future crime occurrences based on specified boroughs, crime types, times of day, and dates. The predictive capability extends to a 30-day horizon, allowing users to anticipate and assess safety conditions in advance, with the option to compare predictions against actual occurrences within the last 30 days of the dataset.

The data mining task at the core of this project involves time series forecasting, a subset of predictive analytics. Specifically, the goal is to utilize historical crime rate data to predict future crime rates for different boroughs in New York City. The predictive analysis aligns with the broader application-oriented objective of aiding users in selecting Airbnb accommodations. This entails leveraging insights from crime rate forecasts to enhance users' confidence in their choice of stay, merging safety considerations with economic factors.

Comparative analysis of time series forecasting models, specifically ARIMA and SARIMA with standardization, revealed SARIMA as the superior choice. SARIMA's effectiveness in capturing both seasonal and non-seasonal components in crime rate data contributed to more accurate predictions. The selection of SARIMA as the best model reflects its aptitude in handling the temporal complexities inherent in crime rate forecasting, thereby enhancing the precision of the interactive dashboard.

## 4. METHODOLOGY

Data Collection: The study commenced with the acquisition of pertinent datasets to form the foundation for analysis. Crime rate data for New York City spanning the years 2014 to 2016 was systematically gathered, encompassing essential details such as crime types, boroughs, time of day, and specific dates. Additionally, comprehensive information regarding Airbnb listings in New York City, including price, borough, review ratings, host identity verification, and instant bookability, was systematically collected.

Data Preprocessing: A crucial phase involved the meticulous preprocessing of acquired datasets to ensure their reliability and uniformity. The challenge of prevalent NaN (Not a Number) values within the crime dataset was effectively addressed through judicious imputation techniques, preserving the integrity of the data. Rigorous data cleaning procedures were implemented across both crime and Airbnb datasets, rectifying discrepancies and errors to ensure consistency and accuracy.

Model Creation: Employing Jupyter Notebook facilitated an in-depth exploration of crime data, revealing discernible patterns, trends, and potential predictors crucial for forecasting. Various time-series forecasting models, including AutoRegressive Integrated Moving Average (ARIMA) and Seasonal ARIMA (SARIMA), were rigorously examined. Following a comparative analysis, SARIMA emerged as the optimal model for crime rate prediction. Standardization techniques were applied to ensure uniform contribution of all features in optimizing the performance of the SARIMA model.

Prototype Development: The translation of methodology and models from the Jupyter Notebook to an interactive Streamlit application marked a pivotal step. This facilitated dynamic user engagement with the datasets. The integration of multiple filters into the Streamlit interface allowed users to personalize crime rate forecasts by selecting specific boroughs, crime types, time of day, and dates. Furthermore, an interactive map feature was incorporated to empower users to explore and select Airbnb listings based on various criteria, enhancing the user experience.

Evaluation and Comparison: The performance of the SARIMA model was critically evaluated by comparing its predictions with actual crime rates, with a specific focus on a 30-day forecast period. A comprehensive Jupyter Notebook documenting the entire process, including data collection, preprocessing, model creation, and obtained results, was generated for transparency and reproducibility.

Challenges and Solutions: Challenges encountered during the study were systematically addressed to ensure the robustness of the methodology. Effective strategies were implemented to handle NaN values in the crime dataset, and the conversion of code to Streamlit was navigated successfully, resulting in a seamless and user-friendly interface with interactive features.

Importance and Usefulness: The significance of forecasting crime rates for potential Airbnb users was underscored, providing valuable insights to facilitate informed decisions about their stay. Emphasis was placed on the importance of conducting a comparative analysis between crime rates and Airbnb listing features, enabling users to assess the correlation between safety considerations and the value of their selected accommodation.

In summary, the methodology employed in this study collectively integrated data collection, preprocessing, model development, and interactive prototype creation. The selection of SARIMA as the forecasting model was justified based on its superior performance, particularly after standardization. The Streamlit application not only offers users a comprehensive tool to explore and evaluate Airbnb listings in the context of crime rates but also enhances their overall vacation planning experience.

## 5 EXPERIMENTAL EVALUATION

This section describes the experimental setup and results we obtained.

## 5.1 EXPERIMENTAL SETUP

The initial dataset comprised crime reports with multiple attributes, reflecting the various aspects of each incident. Preprocessing involved the following steps:

1. Date-Time Conversion: The date and time columns were combined into a single datetime object to facilitate time-series analysis.

2. Indexing: The datetime object was set as the DataFrame index to allow for resampling and time-based filtering.

3. Cleaning: Rows with invalid datetime values were removed to ensure data quality.

4. Aggregation: Data was resampled to a daily frequency, aggregating counts to reflect daily crime incidents.

5. Attribute Selection: Only relevant attributes contributing to the time series analysis were retained, such as 'CMPLNT_FR_DT', 'CMPLNT_FR_TM', and 'BORO_NM'. Before preprocessing, the dataset contained 1410315 instances (number of rows before cleaning) and 24 attributes (24 number of columns before dropping irrelevant ones). After preprocessing, the dataset was restructured into separate time series, with each series representing the daily crime count for a borough from 2014 onwards.

6. Evaluation Measures: To evaluate the forecasting models, the Mean Squared Error (MSE) was calculated to measure the average squared difference between the estimated values and the actual value. This provided a clear indicator of forecasting accuracy, with lower values indicating better performance.

7. Software: The results were generated using Python 3.8, with the pandas library for data manipulation and preprocessing, and the statsmodels library for time series modeling. Specifically, the SARIMAX class from statsmodels was used to fit Seasonal ARIMA models. Python's Matplotlib library was employed for visualizing the data and the forecasting results.
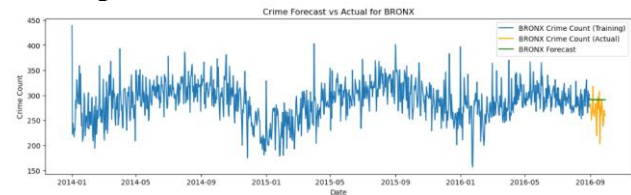
## 5.2 Experimental Results



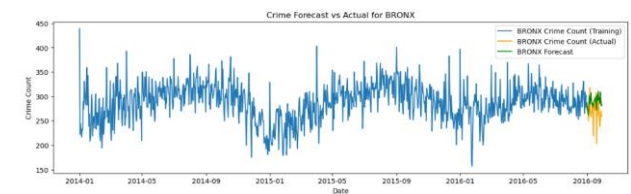**Figure 1:** Arima model forecast of Bronx



**Figure 2:** Sarima model forecast of Bronx

1. https://github.com/nicktony/CSE881-Final-Project

2. The ARIMA model served as the baseline, providing insights into non-seasonal trends. However, its performance, with an RMSE of 1059.46, indicated the need for a more nuanced approach. This led to the exploration of the SARIMA model, which incorporated seasonal components to better capture the data's inherent patterns. The SARIMA model outperformed the baseline ARIMA model, achieving a lower RMSE of 998.07. This improvement highlighted the model's ability to capture complex, multi-year patterns in crime rates. Figures 1 and 2 in the document visually contrast the forecast accuracy of both models against actual crime counts, with the SARIMA model demonstrating closer alignment with real data.

## 5.3 DISCUSSION

This study was aimed at forecasting daily crime rates in urban boroughs using time series analysis. A critical part of the research involved comparing the performance of ARIMA and SARIMA models to identify the most effective approach.

The initial baseline model, ARIMA, provided a foundational understanding of the non-seasonal trends in the data. However, its performance, as gauged by the root mean square error (RMSE), indicated a need for a more sophisticated model that could encapsulate both the trend and seasonal fluctuations. The ARIMA model's RMSE stood at 1059.46, which, while reasonable,

suggested room for significant improvement. In contrast, the SARIMA model, with its incorporation of seasonal components, offered a marked improvement. Uniquely, the seasonal period was set to 32, a non-standard choice guided by the observed data patterns starting from 2014. This approach captured a specific recurring cycle in the crime data that was not aligned with the typical annual or weekly seasonal patterns. The chosen parameters for the SARIMA model were (1,1,2)x(0,2,2,32), reflecting the identified 32-period cycle.

The SARIMA model achieved a notably lower RMSE of 998.07, underscoring its superior predictive accuracy over the ARIMA model. This reduction in RMSE highlights the model's enhanced capability in capturing the complex, multi-year patterns in the dataset. Figures 1 and 2 illustrate the comparative forecast accuracy of the ARIMA and SARIMA models, respectively, against the actual crime counts. The SARIMA model's forecasts are observed to closely align with the actual data, demonstrating its effectiveness.

The results of this experiment underscore the importance of model selection in time series analysis and the value of customizing the model to the specific characteristics of the dataset. The SARIMA model's success, particularly with an unconventional seasonal period of 32, opens avenues for further exploration into similar non-standard patterns in other datasets.

# 6. STREAMLIT APPLICATION

Streamlit is a powerful and user-friendly Python library designed for creating interactive web applications with minimal effort. It enables data scientists and developers to transform data scripts into shareable web applications through a straightforward and intuitive approach. With Streamlit, complex data analyses, visualizations, and machine learning models can be seamlessly converted into interactive dashboards, requiring minimal coding and without the need for a separate web development framework. Its simplicity lies in its ability to automatically refresh the web application in response to changes in the underlying code, making it an ideal choice for projects like the Time Series Forecasting of NYC Crime. Streamlit not only accelerates the development of data-driven applications but also provides an accessible means for users to interact with and derive insights from complex datasets.

In the realm of urban safety and tourism planning, the Time Series Forecasting of NYC Crime project emerges as a valuable tool. Leveraging advanced statistical models, specifically ARIMA (AutoRegressive Integrated Moving Average) and SARIMA (Seasonal ARIMA), this project aims to predict crime rates in different boroughs of New York City. The data used for this analysis spans the years 2014-2016, providing a comprehensive temporal view of crime incidents.

The centerpiece of the project is the "Crime Forecasting" dashboard, which offers users an interactive and insightful experience as shown in Figure 3. Users can specify their criteria by selecting a borough (or all boroughs), type of crime, time of day, and date. The application then employs the ARIMA and SARIMA models to predict crime rates for the selected parameters, projecting the trends 30 days into the future.
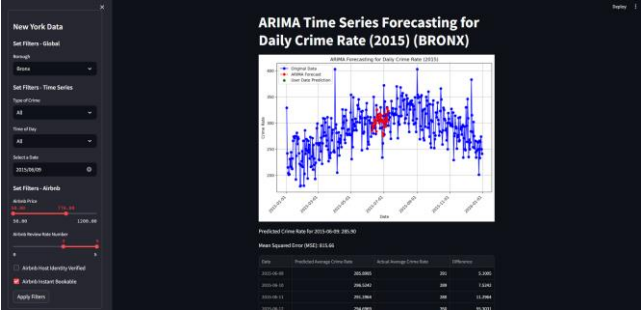


**Figure 3**: Preview of interactive crime forecasting dashboard

For the ARIMA model, predictions are made for the selected date, and the Root Mean Squared Error (RMSE) is calculated to assess the accuracy of the forecast. Additionally, a detailed table is provided, comparing the actual and predicted values along with their absolute differences.

Similarly, the SARIMA model estimates crime rates for the subsequent 30 days after August 30, 2016. A comparative table is presented, showcasing actual and predicted values, along with the absolute differences between them.

To enhance the practical utility of the project, the "Airbnb Interactive Map" offers users the ability to explore potential lodging options based on their preferences. Filters such as price, borough, review rating, host identity verification, and instant bookability allow users to select accommodations meeting their specific criteria. An interactive map interface enables users to click and explore Airbnb listings, with pricing and ratings conveniently displayed upon hover. This feature empowers users to not only consider crime rates but also factor in accommodation preferences for a comprehensive decision-making process. A preview of this map and the varying filters are shown in Figure 4 below.
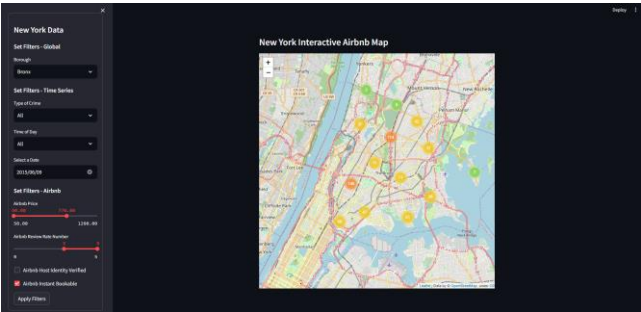


**Figure 4**: Preview of Airbnb interactive map

For users with a more analytical inclination, the "Dataframe Overviews" page provides a detailed exploration of the datasets used in the project, as shown in figure 5. This section offers insights into the Airbnb and NYC Crime datasets, including comprehensive descriptions of each column. This facilitates transparency and encourages data scientists or researchers to delve into the intricacies of the code and datasets, fostering a deeper understanding of the underlying methodologies.
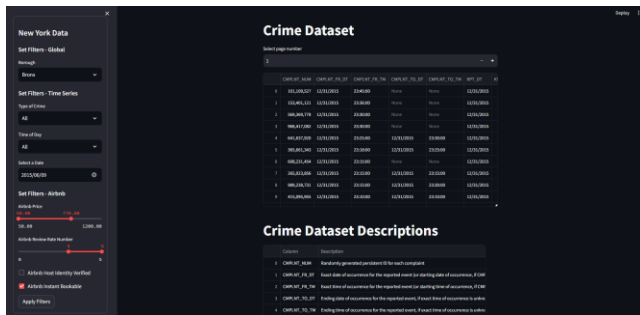
**Figure 5**: Preview of interactive data frame dashboard

The overarching goal of this Streamlit application is to assist individuals in researching NYC area crime trends and making informed decisions about their stay in the city. By seamlessly integrating crime forecasts, Airbnb options, and detailed dataset overviews, this project offers a multifaceted approach to urban safety and travel planning, promoting a safer and more enjoyable experience for residents and visitors alike.

## 7. FUTURE WORK

In the realm of future developments, several avenues present themselves for the advancement of the Time Series Forecasting of NYC Crime project. First and foremost, there is a potential exploration of more sophisticated forecasting models beyond SARIMA, such as Recurrent Neural Networks (RNN) and Long Short-Term Memory (LSTM) RNNs, which could offer a deeper understanding of intricate patterns within crime data. Real-time data integration emerges as another critical area for enhancement, wherein the incorporation of live crime data updates into forecasting models promises increased relevance and accuracy. Addressing the evolving nature of crime data, a focus on dynamic model parameters is recommended, aiming to automatically adjust parameters in response to changing crime patterns.

Implementing user feedback mechanisms stands as a valuable avenue for continuous refinement, allowing users to contribute to the improvement of crime rate forecasts. To broaden the application's utility, expanding its geographic scope beyond New York City is proposed, involving model adaptation to diverse urban contexts. Additionally, enriching Airbnb selection criteria by incorporating features like proximity to public transportation, amenities, and neighborhood characteristics could provide users with more comprehensive decision-making insights.

Privacy considerations become paramount, suggesting the implementation of mechanisms to anonymize or aggregate crime data while still delivering meaningful insights. Lastly, the optimization of the Streamlit application's performance, particularly as datasets grow in size, is crucial and necessitates exploration of efficient strategies for data loading, processing, and visualization.

In the context of Streamlit applications, optimizing performance is crucial for providing a seamless and responsive user experience. One effective strategy to enhance the speed of a Streamlit application is through the strategic use of caching functions. Caching involves storing the results of expensive function calls and returning the cached result when the same inputs occur again. This not only reduces redundant computations but also accelerates the responsiveness of the application.

In Streamlit, the `@st.cache` decorator plays a pivotal role in implementing caching. By applying this decorator to specific functions, users can explicitly instruct Streamlit to cache the results of those functions, avoiding their recomputation on subsequent runs with the same input parameters. This is particularly beneficial for functions involved in heavy computations, data loading, or any operation that remains constant across multiple user interactions.

For instance, when dealing with data retrieval or preprocessing functions that do not change frequently, caching can significantly improve response times. Consider a scenario where a function fetches and processes a large dataset. By caching this function, subsequent runs with the same dataset parameters will retrieve the precomputed result, mitigating the need for redundant processing.

Furthermore, caching can be particularly advantageous when dealing with computationally intensive machine learning models or complex calculations. By caching the results of these functions, users can experience faster interactions with the application, especially when exploring various parameters or configurations.

It's essential to strike a balance when implementing caching. While it optimizes performance, it's crucial to avoid over-caching functions that involve dynamic or frequently changing data, as this may lead to outdated results being presented to the user.

In summary, these future endeavors collectively contribute to the ongoing evolution of the project, ensuring its sustained relevance and efficacy in assisting users with informed decision-making regarding urban safety and travel planning.

## 8. REFERENCES

[1] Schaffer, A.L., Dobbins, T.A. & Pearson, SA. Interrupted time series analysis using autoregressive integrated moving average (ARIMA) models: a guide for evaluating large-scale health interventions. BMC Med Res Methodol 21, 58 (2021). https://doi.org/10.1186/s12874-021-01235-8

[2] Liu, J., Yu, F. & Song, H. Application of SARIMA model in forecasting and analyzing inpatient cases of acute mountain sickness. BMC Public Health 23, 56 (2023). https://doi.org/10.1186/s12889-023-14994-4

[3] Smith, J. (2023, January 15). LSTM Networks: A Detailed Explanation. Towards Data Science. https://towardsdatascience.com/lstm-networks-a-detailed-explanation-8fae6aefc7f9l

[4] Duncan A. Smith,Online interactive thematic mapping: Applications and techniques for socio-economic research,Computers, Environment and Urban Systems,Volume 57,2016,Pages 106-117,ISSN 0198-9715,https://doi.org/10.1016/j.compenvurbsys.2016.01.002.

[5] Jebb, A. T., Tay, L., Wang, W., & Huang, Q. (2015). Time series analysis for psychological research: Examining and forecasting change. Frontiers in Psychology, 6. https://doi.org/10.3389/fpsyg.2015.00727

[6] Mansion, G., Syndicate, K., & Amber Bldg, Hydri Chowk, Naya Nagar, Mira Road (E), Thane, Maharashtra, India. (2023). Stock Prediction Web-App Based on Python-Streamlit Using Data Analysis and Machine Learning. International Research Journal of Modern Education and Technical Science, 05, (08). https://www.doi.org/10.56726/IRJMETS43978