# PREDICTING FUTURE CRIMES USING TIME SERIES FORECASTING

*Tejas Ramesh*
*G01445777*

*Saksham Nayyar*
*G01462522*

*Sahil Shrikrishna Zele*
*G01465963*

*Swapneel Suhas Vaidya*
*G01459609*

*Faculty:*
*Prof. Keren Zhou*
*(kzhou6@gmail.com)*

*Fall,2023*

## 1.INTRODUCTION

Our project, predicting future crimes using time series analysis, focuses on predicting crime incidents in Baltimore by analyzing crime data for the city of Baltimore from 2011-2016. The research question driving this endeavor is "Can time series forecasting models provide accurate predictions of future crime incidents in Baltimore, aiding law enforcement in proactively addressing public safety concerns?"

Crime prediction models, utilizing historical data, have diverse applications. They optimize law enforcement resources, enhance preventive measures through community engagement, inform policy development, aid crisis response, and contribute to smart city initiatives. However, ethical considerations, transparency, and community involvement are vital for responsible implementation and ongoing evaluation of these predictive policing applications.

In our report, we commence by providing a clear definition of the task at hand. Subsequently, we detail the dataset and highlight crucial pre-processing steps. The following sections elaborate on our models, the conducted experiments, and the outcomes obtained. Lastly, we delineate prospective avenues for further research and engage in a discussion on the broader implications of our findings within the context of driver-behavior analysis.

## 2. DATASET

The dataset utilized for analysis originates from the City of Baltimore. Baltimore is an important seaport in Maryland state of United States of America. The city is significantly known for its high crime rate which ranks higher than the national average. The city government along with the office of Mayor provide public access to crime data in the BaltimoreOpenData portal ( https://data.baltimorecity.gov/). This data is updated every week. The open data portal maintains organized primary and secondary data published by the city council, local authorities, police department and public bodies. The data set used in the project contains detailed information of crimes from 2011 to 2016 that have occurred in the city of Baltimore, USA. This crime data set has been downloaded from the public safety domain of the Open Baltimore portal (https://data.baltimorecity.gov/Public-Safety/BPD-Part-1-Victim-Based-Crime-Data/wsfq-mvij), Although the link is currently not accessible.

A total of 286,000 records of crimes are reported. Each crime record comes with both spatial (longitude and latitude) and temporal (date and time of occurrence) information along with the type of crime.
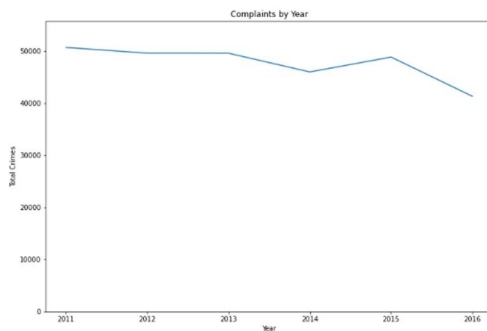
The dataset includes following features: date and time of occurrence, crime code, address where it happened, description of the crime along with (if the crime was committed inside or outside), weapon used, post where the crime is reported, district,

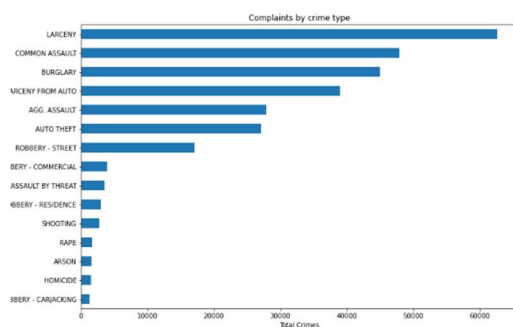neighborhood, location (longitude and latitude) and total incidents.

Here is the snapshot of the dataset:





As illustrated in Figure 1, the graph depicts the total number of crimes from 2011 to 2016. The city of Baltimore experienced an average of 47,000 crimes, reaching its peak in 2011 at approximately 50,000 and gradually decreasing to 42,000 by 2016. This figure represents around 0.244% of the national average, as per the FBI's national crime data.



Delving deeper into the nuances revealed by Fig. 2, it is intriguing to note the stark variations in the frequency of different crime types. Notably, categories like 'homicide' and 'Robbery-Carjacking' stand out with remarkably low occurrence, painting a

distinctive picture of the city's safety landscape. In contrast, crimes such as 'larceny' and 'common assault' dominate the dataset, hinting at the prevalence of certain criminal behaviors.
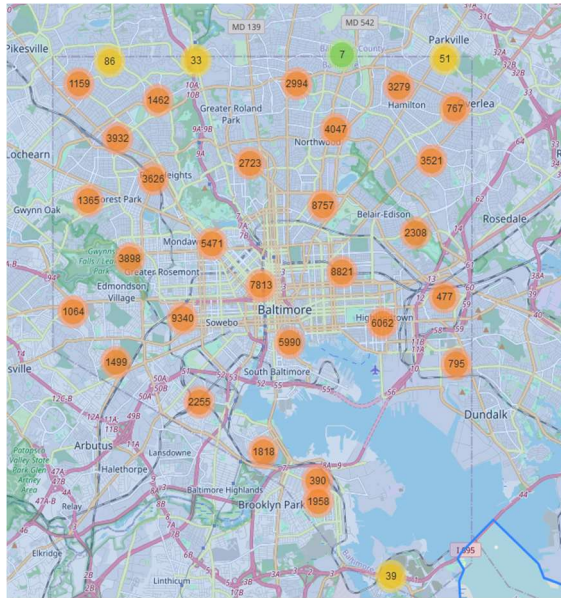
## 3. METHODOLOGY

### 3.1 PREPROCESSING

In the preliminary stage of our research endeavor, emphasis was placed on the preprocessing and exploratory analysis of the dataset, laying the foundation for subsequent time series analyses. A primary focus was directed toward rectifying discrepancies in the time and date columns, characterized by irregularities and missing entries. A systematic approach was employed to standardize the time format, followed by the amalgamation of date and time columns to establish uniform timestamps adhering to the ISO 8061 standard. This standardization was imperative for ensuring consistency and precision in temporal data representation. Concurrently, an examination of the 'Inside/Outside' feature revealed diverse categorizations ('Inside', 'Outside', 'I', 'O', 'NAN'), prompting a normalization to the binary distinctions 'I' or 'O.' This harmonization not only refined the dataset but also facilitated a more cohesive and dependable exploration of crime trends associated with spatial attributes.

### 3.2 EXPLORATORY DATA ANALYSIS

In the exploratory analysis of crime data pertaining to Baltimore from 2011 to 2016, the dataset encompasses comprehensive details, including crime date, time, location, crime description, crime code, and weapon information. These attributes serve as crucial elements for discerning underlying patterns and trends within the dataset. To conduct a thorough Exploratory Data Analysis (EDA), the spatial aspects of the dataset were accentuated through the utilization of Longitude and Latitude attributes. Employing JSON mapping and the FOLIUM library, the geographical distribution of crimes was visually represented. In furtherance of spatial analysis, the location column underwent segmentation into distinct Latitude and

Longitude columns, enhancing the granularity of spatial investigations.



The representation of criminal incidents in Baltimore involves mapping clustered points on the city map. This visual approach, utilizing Folium for an overhead view, helps illustrate the spatial distribution of crimes. Observing these patterns raises questions about the impact of urban features on crime concentrations. Factors like street layouts, lighting, and surveillance may influence where crimes occur. Data preprocessing involved standardizing column formats and handling missing values. Notably, columns like weapon descriptions and locations were adjusted to address missing information, marked as 'NA.' This ensured a consistent dataset for meaningful insights.

### 3.2.1 DECOMPOSITION
We performed seasonal decomposition with additive model based on the incident count for the crimes. Seasonal decomposition refers to decomposing the time series into three components; namely Trend component, Seasonal component, Residual component. Together, we can state that,
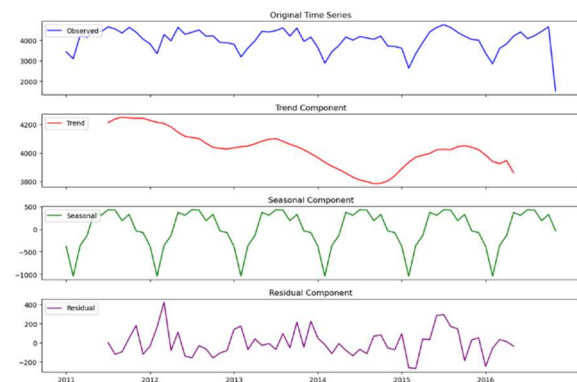
Observed Time series = Trend + Seasonality + Residual

Trend represents the overall long-term movement of data, meaning it suggests the total incidents of crime that have occurred throughout the duration of the data.

Seasonal components capture the regular and repeating patterns at fixed intervals of time. It represents systematic fluctuations in the data. In this dataset it represents how the crime incidents are affected for a short duration, and how the pattern repeated over a duration of time.

Residual components are the error remainder, it represents the random variation in the series that is not accounted in trends and seasonal components.



To get an idea of seasonality and trends in the data, we had to first get the data aligned with those features. We first split all the data representing crimes for years. To achieve this, we used the CrimeDate feature, which was initially preprocessed to get the time stamp values and created a series containing Year and total complaints made in that year.

### 3.2.2 TREND ANALYSIS
Using the time stamp values, we computed a series suggesting the complaints made during the months. This datapoint helped us in deriving the seasonality in terms of whether higher number of crimes were reported in a particular month.

Another useful insight we gathered was regarding the trends was that crimes reported on Fridays were visibly higher than on other days of the week. It also gives insights that crimes were lower on Saturdays and Sundays.



The graph below shows the monthly crime data for all the years in consideration, we can see the seasonality, in the crimes, for all the years, a similar trend is followed for 6 years, this is a useful insight in data.
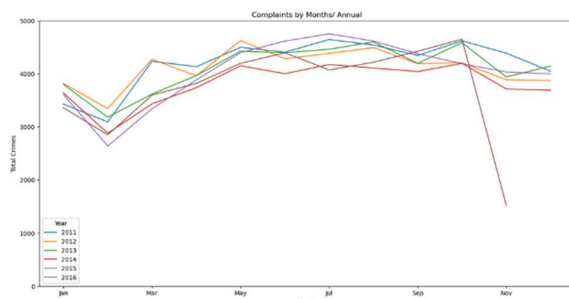


From the graph, we can see that October saw the highest number of crimes being reported. The trend also suggests that the complaints made throughout a given year were in a certain range, and not much varying than the number of complaints in other years.

Other trends we observed were in terms of location, Downtown neighborhood saw most reported crimes

## 3.3 MODELS

Time series analysis involves employing methods to scrutinize time-dependent data, extracting significant statistics and other data characteristics. Time series forecasting, on the other hand, entails utilizing a model to anticipate future values based on previously observed data.

Time series methods are commonly applied to non-stationary data, such as those related to economics, weather, stock prices, and retail sales. In this project, we employed different models to analyze historical crime data and subsequently make forecasts about future criminal activities.

### 3.3.1 LINEAR REGRESSION

Linear regression is a statistical method commonly used for predicting the value of a dependent variable based on one or more independent variables. While linear regression is not specifically designed for time series forecasting, it can be applied to time series data with some considerations.

Linear regression was employed as a baseline model to capture linear relationships between time and crime counts. The model was trained using historical crime data, and predictions were made based on the linear relationship identified during training.

### 3.3.2 EXPONENTIAL SMOOTHING

Exponential smoothing is a popular time series forecasting method. It operates on the principle of assigning weights to past observations, with the weights decreasing exponentially. This ensures that recent observations have a greater influence on the forecast than older ones. The fundamental parameter governing this weighting is the smoothing parameter (alpha).

As time series data often exhibits trends and seasonality, exponential smoothing can be extended

to handle these components. Double Exponential Smoothing incorporates trend information, while Triple Exponential Smoothing (Holt-Winters method) considers both trend and seasonality, providing a more robust forecasting tool.

### 3.3.2.1 SINGLE EXPONENTIAL SMOOTHING

Single Exponential Smoothing is a foundational method in this approach. The formula for SES involves updating the forecast based on the exponentially weighted average of past observations. The simplicity of SES makes it an effective tool for quick and efficient forecasting.

The single exponential smoothing formula is given by:

$s_t = \alpha x_t + (1 - \alpha)\, s_{t-1} = s_{t-1} + \alpha\,(x_t - s_{t-1})$

### 3.3.2.2 DOUBLE EXPONENTIAL SMOOTHING

This method is also called Holt's trend corrected or second-order exponential smoothing. This method is used for forecasting the time series when the data has a linear trend and no seasonal pattern. The primary idea behind double exponential smoothing is to introduce a term to consider the possibility of a series showing some form of trend. This slope component is itself updated through exponential smoothing.

The double exponential smoothing formulas are given by:

$S_1 = x_1$

$B_1 = x_1 - x_0$

For t>1,

$s_t = \alpha x_t + (1 - \alpha)(s_{t-1} + b_{t-1})$

$\beta_t = \beta(s_t - s_{t-1}) + (1 - \beta)b_{t-1}$

Here,

$s_t$ = smoothed statistic, it is the simple weighted average of current observation $x_t$

$s_{t-1}$ = previous smoothed statistic

$\alpha$ = smoothing factor of data; $0 < \alpha < 1$

t = time period

$b_t$ = best estimate of trend at time t

$\beta$ = trend smoothing factor; $0 < \beta < 1$

### 3.3.2.3 TRIPLE EXPONENTIAL SMOOTHING

In this method, exponential smoothing applied three times. This method is used for forecasting the time series when the data has both linear trend and seasonal pattern. This method is also called Holt-Winters exponential smoothing.

The triple exponential smoothing formulas are given by:

$$
\begin{aligned}
s_0 &= x_0 \\
s_t &= \alpha \frac{x_t}{c_{t-L}} + (1 - \alpha)(s_{t-1} + b_{t-1}) \\
b_t &= \beta(s_t - s_{t-1}) + (1 - \beta)b_{t-1} \\
c_t &= \gamma \frac{x_t}{s_t} + (1 - \gamma)c_{t-L}
\end{aligned}
$$

Here,

$s_t$ = smoothed statistic, it is the simple weighted average of current observation $x_t$

$s_{t-1}$ = previous smoothed statistic

$\alpha$ = smoothing factor of data; $0 < \alpha < 1$

t = time period

$b_t$ = best estimate of a trend at time t

$\beta$ = trend smoothing factor; $0 < \beta < 1$

$c_t$ = sequence of seasonal correction factor at time t

$\gamma$ = seasonal change smoothing factor; $0 < \gamma < 1$

### 3.3.3 NAÏVE MODEL

The Naive Model operates on the assumption that the future values of a time series are solely influenced by the most recent observation. It simplifies forecasting by projecting the next time step's value as equal to the latest observed data point.

Formula

Mathematically, the Naive Model can be expressed as:

$s(t+1)=s(t)$ where:

- $s(t+1)$ is the forecasted value for the next time step.

- $s(t)$ is the most recently observed value.

The Naive model operates under the assumption that trends remain constant and rely solely on the latest observation. However, it falls short in addressing the complexities introduced by seasonality, cycles, and other intricate structures embedded within the time series.

### 3.3.4 MOVING AVERAGE

Moving averages are a fundamental tool in time series analysis, employed to smooth out fluctuations and identify underlying trends in data. They are particularly useful in revealing patterns and making data more interpretable for forecasting purposes. One popular method is the Moving Average (MA) model.

A moving average is calculated by taking the average of a subset of data points within a specified window or period. The Simple Moving Average (SMA) involves equally weighting all data points in the window, while the Weighted Moving Average (WMA) assigns different weights to different points.

The mathematical formula for the Moving Average model is:

$$MA(t) = \frac{1}{n}\sum_{i=1}^{n} Y(t-1)$$

Where:

- $MA(t)$ is the moving average at time $t$,

- $n$ is the order of the moving average model,

- $Y(t-i)$ represents the observed values in the time series.

Moving averages play a vital role in smoothing time series data, reducing noise, and revealing underlying trends. They are particularly effective in highlighting long-term patterns while minimizing the impact of short-term fluctuations.

In our project, we employed four variations of moving averages, distinguished solely by the increment in the window size over which data points are averaged.

- **3-Point Moving Average:**

$$MA(t) = \frac{1}{3}\sum_{i=1}^{3} Y(t-i)$$

- **6-Point Moving Average:**

$$MA(t) = \frac{1}{6}\sum_{i=1}^{6} Y(t-i)$$

- **9-Point Moving Average:**

$$MA(t) = \frac{1}{9}\sum_{i=1}^{9} Y(t-i)$$

- **12-Point Moving Average:**

$$MA(t) = \frac{1}{12}\sum_{i=1}^{12} Y(t-i)$$

### 3.3.5 ARIMA MODEL

One powerful tool in this realm of time series forecasting is the Autoregressive Integrated Moving Average (ARIMA) model. ARIMA is a versatile forecasting method that combines autoregression (AR), differencing (I), and moving averages (MA) to model a wide range of time series patterns.

ARIMA is well-suited for handling non-stationary time series data. The autoregressive component (AR) captures the relationship between an observation and several lagged observations. The differencing component (I) transforms a non-stationary time series into a stationary one by computing differences

between consecutive observations. The moving average component (MA) deals with the influence of past white noise on the current value.

The ARIMA model is denoted as ARIMA (p, d, q), where:

- **p**: Order of the autoregressive component.

- **d**: Degree of difference.

- **q**: Order of the moving average component.

$$Y_t = \mu + \phi_1 Y_{t-1} + \phi_2 Y_{t-2} + \ldots + \phi_p Y_{t-p} + \epsilon_t - \theta_1 \epsilon_{t-}$$

Where:

- $Y_t$ is the observed value at time $t$.
- $\mu$ is the mean of the series.
- $\phi_1, \phi_2, \ldots, \phi_p$ are autoregressive coefficients.
- $\epsilon_t$ is white noise at time $t$.
- $\theta_1, \theta_2, \ldots, \theta_q$ are moving average coefficients.

Extensions of ARIMA include Seasonal ARIMA (SARIMA), which incorporates seasonal components. SARIMA is denoted as SARIMA (p, d, q) (P, D, Q) s, where *P, D, Q* are similar to *p*,*d*,*q* but for the seasonal component, and *s* is the seasonality parameter.

In our project we have used SARIMAX (Seasonal Autoregressive Integrated Moving Average with exogenous factors). SARIMAX includes exogenous variables in addition to the components of SARIMA. Exogenous variables are external factors that can affect the time series but are not influenced by it. The model allows for incorporating the influence of these external variables in the forecasting process.

To select the best parameter, we are performing a grid search over different combinations of parameters for a Seasonal Autoregressive Integrated Moving Average with Exogenous Factors (SARIMAX) model. And then selecting the parameters which gives us minimum RMSE.

# 4. EXPERIMENTS

## 4.1 SETUP

In the initial stages of setting up the forecasting framework and establishing evaluation parameters, a crucial step involved partitioning the dataset into training and testing subsets. The temporal data spanning from 2011 to 2016 was segmented into weekly intervals, with corresponding total incident counts. To address missing data for several weeks in November and December 2016, mean imputation was employed, computed based on analogous weeks in the preceding years (2011 to 2015). The training dataset encompassed the period from the first week of 2011 to the week ending on 08/30/2015, constituting approximately 70% of the entire dataset. Subsequently, the testing dataset incorporated the remaining weeks, and the designated evaluation metrics, including Root Mean Squared Error (RMSE), Normalized Mean Squared Error (NMSE), and U-statistics, were applied for assessment.

# 5. RESULTS

## 5.1 EVALUATION METRIC

RMSE (ROOT MEAN SQUARED ERROR):

The Root Mean Squared Error (RMSE) is a commonly used metric to measure the accuracy of a predictive model, including time series forecasting models. It provides a way to quantify the average magnitude of the errors between predicted and actual values. The RMSE is expressed in the same units as the predicted and actual values, making it interpretable and easy to understand.

The RMSE is calculated as follows:

RMSE = $(\sum_{i=1}^{n}(y_i - \hat{y}_i)^2/n)^{0.5}$

Where:

n is the number of observations.

$y_i$ is the actual value at time i.

$\hat{y}_i$ is the predicted value at time i.

## NMSE (Normalized Mean Squared Error):

The Normalized Mean Squared Error (NMSE) is a metric used to evaluate the performance of a forecasting model. It is a normalized version of the Mean Squared Error (MSE) and provides a relative measure of how well the model is performing compared to a simple baseline model. NMSE is particularly useful when you want to understand the proportion of the variance in the predicted values relative to the variance in the actual values.

The NMSE is calculated as follows:

NMSE=MSE/Var(y)

Where:

MSE is the Mean Squared Error, calculated as the average of the squared differences between the predicted and actual values.

Var(y) is the variance of the actual values.

## U-Statistic (Theil's U-Statistics):

Theil's U-Statistic is a statistical measure used to evaluate the forecast performance of a model. It is used in the context of time series forecasting. Theil's U assesses the relative accuracy of a forecasting model by comparing the forecasted values to the actual values. It was introduced by Henri Theil, a Dutch economist, and is part of a family of statistical measures designed to evaluate forecasting accuracy.

Theil's U is calculated as follows:

$$U=\left(\sum_{t=1}^{T}((y_t-\hat{y}_t)/y_t)^2 / \sum_{t=1}^{T}(y_t/\bar{y})^2\right)^{0.5}$$

Where:

T is the total number of time periods.

$y_t$ is the actual value at time t.

$\hat{y}_t$ is the forecasted value at time t.

$\bar{y}$ is the mean of the actual values.
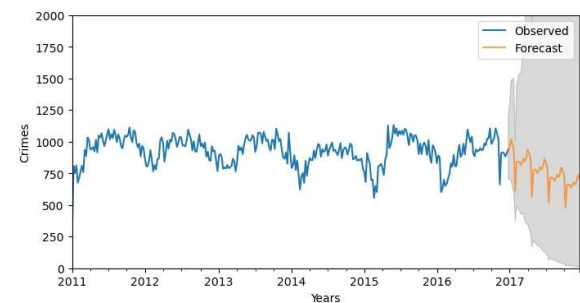
## 5.2 Model Performance

Below is an image, how all the models compare to each other using the evaluation metrics defined above. We can see, 3-point moving average is the best metrics according to all evaluation metrics.

| | Test RMSE | Test NMSE | Test U-stat |
|---|---|---|---|
| RegressionOnTime | 114.539289 | 1.010685 | 0.123930 |
| Alpha=0.0976:SimpleExponentialSmoothing | 174.858279 | 2.355479 | 0.189194 |
| Alpha=1.490116e-08,Beta=3.694809e-16:DoubleExponentialSmoothing | 360.705570 | 10.023342 | 0.390279 |
| Alpha=0.07575,Beta=0.0541,Gamma=0.4107:TripleExponentialSmoothing | 162.265218 | 2.028419 | 0.175569 |
| Naive Model | 177.280822 | 2.421198 | 0.191815 |
| SimpleAverageModel | 114.419510 | 1.008572 | 0.123800 |
| 3pointTrailingMovingAverage | 57.398674 | 0.253811 | 0.062105 |
| 6pointTrailingMovingAverage | 75.150899 | 0.435086 | 0.081312 |
| 9pointTrailingMovingAverage | 88.331849 | 0.601093 | 0.095574 |
| 12pointTrailingMovingAverage | 98.690973 | 0.750346 | 0.106782 |
| Arima | 147.415625 | 1.674148 | 0.159502 |

In the image below, we have plotted a graph to show the comparison of test predictions between various forecast models used.



In the image below, we have plotted a graph for showing the lower and upper bound forecast received from SARIMAX model.



## 5.3 challenges (univariant data)

The execution of this final project was not without its challenges. Recognizing and addressing these challenges played a crucial role in the overall project

management and outcomes. The key challenges encountered include:

### 5.3.1 QUALITY OF DATA

One notable challenge pertained to the quality of the dataset. The data collected from the BaltimoreOpenData portal exhibited instances of incomplete or missing information. This required careful consideration and appropriate handling to ensure the reliability and accuracy of our analyses.

### 5.3.2 VISUALIZATION FEASIBILITY

Due to the granularity of the dataset and resource limitations, there were constraints on the feasibility of certain visualization methods. This limitation impacted the team's ability to explore certain aspects of the data in more detail, necessitating alternative approaches to convey meaningful insights.

### 5.3.3 DATA LOSS/SIMPLIFICATION DURING AGGREGATION

The aggregation of data to a higher level for weekly analysis posed challenges related to potential data loss and simplification. While this aggregation was necessary for certain aspects of the project, it introduced complexities in terms of maintaining the richness of features and nuanced patterns present in the original dataset.

Addressing these challenges required a collaborative and adaptive approach. The team adopted strategies such as data imputation for missing values and adjusting the granularity of analysis to mitigate the impact of challenges on the project's overall success. Despite these challenges, the team successfully navigated through them, contributing valuable insights to the field of crime prediction using time series forecasting.

## 6. RELATED WORK

The task of predicting future crimes using time series forecasting has garnered significant attention in the realm of predictive policing and urban analytics. Researchers and practitioners alike have explored various methodologies and models to address the complexities associated with crime prediction. In this section, we review existing work that aligns with our research focus, emphasizing key findings, methodologies, and limitations.

### 6.1 TIME SERIES FORECASTING IN CRIME PREDICTION

Several studies have delved into the application of time series forecasting techniques for predicting crime incidents. One notable work by [1] Makridakis et al., 1998 applied autoregressive integrated moving average (ARIMA) models to predict crime rates in a major urban center. Their findings suggested that ARIMA models effectively captured temporal patterns, providing valuable insights for law enforcement.

In a different approach, [2] Sun et al., 2019 explored the use of machine learning algorithms, including neural networks, to forecast crime occurrences. Their research highlighted the potential of advanced models in capturing non-linear relationships within crime data, offering enhanced predictive accuracy compared to traditional methods.

### 6.2 SPATIO-TEMPORAL ANALYSIS IN CRIME PREDICTION

Spatial considerations play a crucial role in crime prediction models. [3] Zhao et al., 2016 conducted a spatio-temporal prediction of crime events in Baltimore using Convolutional Long Short-Term Memory (CLSTM) neural networks. Their work demonstrated the significance of incorporating both spatial and temporal dimensions for accurate predictions, particularly in urban settings.

### 6.3 CHALLENGES AND ETHICAL CONSIDERATIONS

Predictive policing, including crime forecasting, raises ethical concerns regarding bias, transparency, and community engagement. [4] Lum and Sidelle, 2016

addressed these issues by proposing a framework for responsible implementation of predictive policing tools. Their work emphasized the importance of transparency in model development, community involvement, and ongoing evaluation to mitigate potential biases and enhance public trust.

## 6.4 LIMITATIONS OF EXISTING MODELS

Despite advancements, existing models exhibit limitations. [5] Huang et al., 2020 highlighted challenges related to data quality, resource constraints in visualization, and potential data loss/simplification during aggregation. Additionally, they identified a gap in online projects focusing on time series forecasting using the specific data source under consideration.

## 7. CONCLUSION

In conclusion, we have implemented and studied various forecasting models with respect to predicting the future crimes in the Baltimore city area of United States of America. With the challenges and limitations faced during our work here and using the various forecasting models we have implemented till now. We can see very good improvements in the Evaluation metrics and can confirm that there's still room for further improvement. Soon, we can perform more in the exploratory data analysis area to get a better grasp on the data and select a more precise and sophisticated model such as neural networks to get more accurate predictions. And publish a research paper, which can be helpful to other researchers studying in the same areas and aiding the law enforcement to study and act proactively on safety concerns in a better way.

## 8. DIVISION OF WORK

*TEJAS RAMESH (G01445777):*

**Role**: Literature review and implementation of the ARIMA model.

Achievements: Conducted a comprehensive literature review on time series forecasting. Implemented the ARIMA model, fine-tuned

parameters, and performed statistical evaluations.

*SAKSHAM NAYYAR (G01462522):*

**Role:** Data-focused tasks, including collection, pre-processing, and EDA and implementing Naïve and Linear regression models.

**Achievements:** Cleaned and standardized the dataset, implemented visualizations using Folium, provided key insights into crime patterns, and implemented Naïve and linear Regression models.

*SAHIL SHRIKRISHNA ZELE (G01465963):*

**Role**: Data preprocessing and Exponential Smoothing model implementation.

**Achievements**: Ensured dataset readiness by cleaning time and date columns. Implemented the Exponential Smoothing model and contributed to seasonal decomposition analysis.

*SWAPNEEL SUHAS VAIDYA (G01459609):*

**Role**: Exploration of related work and implementation of moving averages models

**Achievements**: Implemented all moving average models, conducted in-depth result analysis of the related work.

COLLABORATIVE TASKS:
- Model selection, hyperparameter tuning, and evaluation metrics implementation.
- Finalizing the report structure and content.
- Addressing challenges related to data quality and resource constraints.

## 9. REFERENCES

[1] Makridakis, S., Spengler, T., & Hibbert, M. (1998). Statistical forecasting methods. In Companion to Contemporary Economic Geography (pp. 110-136). Routledge.

[2] Sun, J., Wang, W., & Kamel, I. (2019). Crime forecasting using spatio-temporal recurrent neural networks. In Proceedings of the 27th ACM

International Conference on Information and Knowledge Management (pp. 2407-2416).

[3] Zhao, J., Esch, T., & Gómez-Hernández, J. A. (2016). Spatio-temporal prediction of crime hotspots using social media data. Crime Science, 1(1), 1-20.

[4] Lum, C., & Sidelle, A. (2016). Algorithmic fairness: When to trust your algorithms? In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (pp. 963-964).

[5] Huang, Y., Zhao, Q., & Zheng, Y. (2020). A systematic review of time series forecasting studies on crime prediction. Journal of Quantitative Criminology, 36(1), 1-32.

[Spatio-Temporal Prediction of Baltimore Crime Events Using CLSTM Neural Networks | IEEE Journals & Magazine | IEEE Xplore](#)

[Exponential Smoothing- Definition, Formula, Methods and Examples (byjus.com)](#)

[Microsoft Word - Timeseries_anls.doc (arxiv.org)](#)