

ANALYSIS OF 911 CALLS FOR DETROIT AND NEW YORK CITY

Swadha Bhatt
Computer Science
Towson University
Towson, USA
sbhatt8@students.towson.edu

Krishna Sruthi Velaga
Computer Science
Towson University
Towson, USA
kvelaga1@students.towson.edu

Meghana Desai Jitendrakumar
Computer Science
Towson University
Towson, USA
mdesai2@students.towson.edu

Abhishek Rangi
Computer Science
Towson University
Towson, USA
arangi1@students.towson.edu

Abstract – The study provides a comprehensive analysis of 911 emergency call data from Detroit and New York City, aiming to enhance public safety and optimize emergency responses. Utilizing advanced data mining techniques, we explored patterns and trends within the call data to assist various city departments, including police, hospitals, and emergency services, in improving operational efficiencies and resource allocation. Detroit, known for its high crime rates, and New York City, noted for its dense population, presented unique challenges and opportunities for understanding urban emergency response dynamics. Our methodology included predictive modeling, anomaly detection, and spatial analysis, which facilitated insights into the frequency, distribution, and nature of emergency calls. The findings are expected to contribute to targeted interventions and strategic planning, ultimately fostering safer urban environments. This project underscores the potential of data analytics in public safety operations, providing a blueprint for other cities to leverage technology in crisis management and response optimization.

Keywords – 911 Calls, Public Safety, Response Optimization, Predictive Modeling, Anomaly Detection, Resource allocation

I. INTRODUCTION

In this country, whenever there is any sort of emergency, we dial 911 with the hope that we can receive immediate help to resolve the situation that we are currently in. When an individual calls 911 they are immediately connected to a 911 operator who asks us that individual series of questions to understand what sort of emergency response would be best suited for that victim. 911 call services can often mean the difference between life and death. Whether it's administering first aid, extinguishing fires, or apprehending suspects, timely intervention can save lives and mitigate the severity of injuries. The existence of 911 services contributes to a sense of security and confidence within communities, knowing that help is just a phone call away and times of crisis.

For our significant project in this class, we have initially considered analyzing the 911 emergency calls for Baltimore and New York City. Baltimore City has long been noted for its violence, intertwined deeply with complex societal issues, making it challenging to derive clear insights from the data available. However, upon further review, we found that the Baltimore dataset did not provide sufficient information suitable for the various data mining models we plan to utilize.

Instead, we have shifted our focus to analyzing the 911 calls from Detroit. According to research from the Rochester Institute of Technology, Detroit had the highest homicide rate among cities of similar population sizes, with a rate of 50 per 100,000 residents in recent years. In 2022, Detroit was again ranked with the third-highest homicide rate across our entire sample. This stark statistic underscores the urgent need for deeper analysis and understanding, which could significantly impact emergency response strategies and public safety measures. New York City, on the other hand, is not known to have as much violence as Detroit however, because it is the most populated city in the nation, it will be very interesting to compare the analysis between the 911 calls of Detroit vs New York City.

Overall, our problem statement for our project is to analyze 911 emergency calls for Detroit, and New York City with the hope of helping both cities become safer places for their citizens. By analyzing the 911 data, we will help five different departments of each city. Those five different departments include the police department, hospitals, ambulance, district leaders, and neighborhoods. Listed below are the five departments and how those departments will use the analyzed 911 call data.

Police Department:

- Utilize the data to identify high-crime areas and hotspot locations for targeted patrolling and crime prevention efforts.
- Analyze trends in emergency calls related to criminal activities to inform law enforcement strategies and resource allocation.

Hospitals:

- Use the data to anticipate surges in demand for medical services based on patterns of emergency medical calls.
- Coordinate with emergency medical services (EMS) to optimize hospital resources and staff allocation for timely response to emergencies.

District Leaders:

- Use the data to inform community policing strategies and prioritize public safety initiatives based on local emergency call patterns.
- Collaborate with law enforcement agencies and community organizations to address neighborhood-specific safety concerns identified through the analysis.

Neighborhoods:

- Empower community members with information about prevalent safety concerns and emergency response protocols based on analyzed 911 call data.
- Establish neighborhood watch programs or community patrols in response to identified safety issues and trends.

This project is very much connected to the field of data mining because while analyzing the 911 call data, we will find patterns that include but are limited to the common types of emergencies, the common peak time calls, and geographic distribution. While discovering patterns in the data we will then be able to efficiently help the five main departments - listed above- in each of the cities use the 911 data efficiently. The second reason the analysis of 911 is connected to data mining is through predictive modeling. Our goal while analyzing the data is to create predictive models with the hope to forecast future trends in the 911 call volume. Predicting these trends will have many of the five departments that were listed previously. The third way that this project is connected to data mining is anomaly detection. Through anomaly detection, we will detect unusual patterns including spikes in call volume, unexpected changes in padder, etc.

II. LITERATURE REVIEW

Mining 911 calls in New York City: Temporal Patterns, Detection, and Forecasting

Researchers in the paper 'Mining 911 calls in New York City: Temporal Patterns, Detection, and Forecasting' had three main objectives while analyzing the 911 call of New York City.

Objectives:

- Identifying Patterns: Analyzing when and where emergency calls happen to understand local patterns over various times and places. This gives a clearer picture than just using broad crime statistics.
- Predicting Future Needs: Developing models that predict where and when police resources will be needed, based on several factors that might

influence crime. This reduces the need for commanders to rely solely on their judgment.

- Spotting Unusual Events: Detecting emergency calls that are out of the ordinary, allowing commanders to respond quickly to unexpected or extreme situations.

Methodology

The first objective was fulfilled with the help of temporal behavior patterns. They created a visual analysis in which three distinct behavioral patterns for emergency call times throughout the day were created based on K means clustering.

- Cluster 0: This group shows a higher volume of non-crime related calls during workday hours, which typically peak in the middle of the day.
- Cluster 1: This pattern has a peak of crime-related calls in the evening hours on weekdays.
- Cluster 2: Here, there is a peak in crime calls late at night on weekends, which aligns with what NYPD officers have noted anecdotally.

These clusters were formed by analyzing emergency call data, both crime and non-crime, across different days of the week and times, and then standardizing the data (z-normalizing). The goal was to minimize the difference between the calls in a cluster and the average pattern of calls within that cluster.

The second objective was fulfilled by creating a predictive model which followed the rolling forecast method. Rolling Forecast Method: The model is updated using the latest data. After predicting call demand for a given day, that day's data is added to the training set, and the model is retrained for the next day's prediction.

Initial Training Data: The model starts with a base of historical data, specifically 90 eight-hour periods (referred to as boundary), which provides the foundation for initial predictions.

Continuous Learning: As each prediction is made, the corresponding actual data is used to retrain the model, so it continually improves and adapts over time.

Model Types: Two main types of models were tested:

Random Forest Regression: Chosen for its robustness in handling complex, high-dimensional data and large sample sizes. It works by creating multiple decision trees and combining their predictions. The model used 100 trees, with more trees not significantly improving results.

Poisson Regression: A statistical model that predicts counts, like the number of calls, as a log-linear function of the features. It is suitable for modeling events that occur at a constant rate within a fixed period.

Performance:

Initial Random Forest Model:

- Predicted daily call counts.
- $R^2 = .7$ $\rho = 0.83$ – Good at detecting patterns.

Shift-Specific Predictions:

- Adapted for 8-hour shifts.
- Lower accuracy: $R^2 \approx 0.5$, $\rho = 0.7$.
- Overestimates low and underestimates high call counts

The third objective was accomplished by spatial clustering. To be more specific, the model is an adaptation of Kull Dorff's population-based Poisson scan statistic, widely used to detect disease outbreaks. It operates on a uniform two-dimensional grid, evaluating rectangular regions for unusually high call counts. The counts in these regions are assumed to follow a Poisson distribution, where the call rate is unknown, but the expected historical call count is known.

An Integrated Model For Crime Prediction Using Temporal And Spatial Factors

Crime prediction has garnered significant attention in both governmental and academic spheres due to its potential to enhance public safety. Various predictive policing tools, such as PredPol, have been implemented, using historical data to predict potential crime spots and timings.

Spatial and Temporal Dynamics: Understanding the interplay between spatial and temporal factors is crucial for enhancing crime prediction accuracy. Spatial data encompasses geographical and urban metrics, while temporal data involves the analysis of the timing and sequence of criminal activities.

Methodological Approaches: Existing methodologies include Kernel Density Estimation (KDE) for identifying crime hotspots and advanced statistical models like Random Walk and Self-Exciting Point Processes, which consider the dependency of current events on historical data. Gaussian Process Regression (GPR) has also been noted for its effectiveness in modeling spatial relationships.

Integration of Data Sources: The effectiveness of crime prediction models is significantly enhanced by integrating diverse data sources. This integration allows for a comprehensive analysis that includes geographical data, Points of Interest (POI), urban dynamics, census data, and social media.

Advancements in Predictive Models: The Continuous Conditional Random Field (CCRF) model has been adapted into a Clustered-CCRF version to better handle spatial and temporal correlations by clustering similar areas based on their spatial characteristics. This method demonstrates improved prediction accuracy over traditional models like ARMA and Linear Regression, which are limited by their handling of multi-dimensional data and their inability to integrate multiple data sources effectively.

Challenges and Future Directions: Despite advancements, crime prediction remains complex due to the multifaceted nature of criminal activities and external influences. Ongoing research is necessary to refine existing models and explore new data sources and analytical techniques to improve prediction capabilities further.

Mining Patterns From 9-1-1 Calls Dataset: Identification, Prediction, And Resource Allocation [4]

The research in a study called "Mining Patterns from 9-1-1 Calls Dataset" tried to look at Montgomery County's 9-1-1 call data to improve emergency response systems. Their research focused on identifying patterns and showing future demands to optimize resource allocation and improve emergency response times.

Objectives:

Identifying Patterns: The study analyzed the time frame and location. patterns of the data to determine when and where emergency calls occurred. This study aids in identifying places and periods of high emergency happening, resulting in a more detailed understanding than broad statistics summaries.

Predicting Future Needs: The project aims to reduce dependency on human judgment by creating models to predict where and when emergencies may occur. This predictive technique helps to strategically plan the deployment of emergency response resources.

Resource Allocation: An important goal was to allocate resources effectively, identifying how to best spread emergency services such as police, fire, and medical help based on known trends and projected demand in the future.

Methodology:

The study used K-means clustering to divide 9-1-1 calls into various groups based on variables such as location, time, and kind of emergency. This strategy was important in identifying trends in the data. This technique was necessary for identifying patterns in the data:

Cluster 0: Typically, non-emergency calls peak around the afternoon.

Cluster 1: There was an increase in emergency calls in the evenings, indicating that services were mostly required at this time.

Cluster 2: High emergency call volumes occurred late at night on weekends, overlapping with periods of increased involvement in society.

Methods used:

K-means Clustering

In the study "Mining Patterns from 9-1-1 Calls Dataset," the K-means clustering approach is thoroughly employed to assess and group emergency call data based on many parameters such as location, call time, and emergency kind. This technique is critical for identifying trends that might help guide and improve emergency response strategies. Here is a more extensive explanation of how K-means clustering is used in this case. For

Implementation of k-means clustering, the first stage in K-means clustering is to randomly choose K starting centers from the dataset. These centers serve as the beginning points for each cluster, and their initial location has an important effect on the algorithm's output. The choice of K, or the number of important clusters, is critical and is frequently established using approaches such as the Elbow Method.

After setup, each data point in the dataset is assigned to the closest cluster. This assignment is based on the distance calculated by Euclid between each data point and the centroids. The objective is to reduce the distance between data points and the centroid of their respective clusters to create cohesive clusters. Each point is paired with the nearest center, resulting in clusters of similar points. The assignment and update methods are carried out iteratively until the method converges. Convergence occurs when the centroids no longer move significantly, indicating that more iterations will result in insignificant changes, or when the cluster assignments remain constant between iterations. This repeated procedure guarantees that the clusters are as highly accurate and indicative of the underlying patterns in the data as possible. The study paper's extensive use of K-means clustering tries to find specific patterns in 9-1-1 call data that are critical for strategic emergency response planning. By efficiently clustering the data, it is possible to identify places and periods with high emergency call frequencies and various sorts of emergencies, allowing for better informed resource allocation and emergency response planning.

Hotspot Analysis:

In the study "Mining Patterns from 9-1-1 Calls Dataset," hotspot analysis is used to find areas in Montgomery County with statistically significant concentrations of emergency calls. This method serves as crucial to determining which places may require more specialized emergency services and resources. The hotspot analysis in this research works as follows:

Hotspot analysis is used to identify locations in Montgomery County with high or low numbers of emergency calls. This geographical study employs the Getis-Ord G_i^* statistic, a widely developed tool for locating physical clusters or hotspots of activity. The method starts with identifying physical connection between the data points, which effectively calculates how calls are spread over geographical space. Each data point is then examined to see if the number of calls in its area is more or less than the average throughout the whole dataset.

The Getis-Ord G_i^* statistic assigns a Z-score to each data point, indicating if the surrounding region is a hot or cool place. A high positive Z-score shows that a region has a larger than average concentration of calls, making it a hotspot, whereas a negative Z-score indicates a cold spot. This statistical method enables researchers to identify locations that require additional attention or resources, resulting in more targeted and effective emergency response methods. The hotspot analysis results share the geographic location of calls and inform strategic decisions about resource allocation and emergency service deployment in the county.

Profiling And Prediction Of Non-Emergency Calls In New York City [3]

The paper introduces a comprehensive model for analyzing crime and criminal data, aiming to provide valuable insights for specialists, law enforcement agencies, and decision-makers. By leveraging data mining techniques, the model seeks to uncover patterns, trends, and relationships within the data, facilitating forecasting, mapping criminal networks, and identifying potential suspects. Through manual collection from police departments in Libya, both crime and criminal data were gathered to develop and test the proposed model. To ensure the accuracy and cleanliness of the data, various preprocessing techniques were applied, including cleaning, handling missing values, and removing inconsistencies.

Two primary data mining techniques are employed in the proposed model: clustering and association rules mining. Clustering utilizes the k-means algorithm to group crimes and criminals based on shared attributes. This process enables the identification of common characteristics and relationships, providing a deeper understanding of crime patterns and behaviors. Association rules mining, on the other hand, utilizes the Apriori algorithm to extract frequent patterns from the crime dataset. By identifying these patterns, decision-makers can take proactive measures to prevent crime based on the discovered associations.

The analysis of the crime and criminal data is conducted using WEKA mining software and Microsoft Excel. These tools offer robust capabilities for data preprocessing, analysis, and visualization, allowing for a comprehensive examination of the dataset. Through the utilization of these software tools, the study presents the results of both clustering and association rules mining. Clustering results in the grouping of crimes and criminals based on important attributes, while association rules highlight frequent patterns within the dataset. Visualizations, like confusion matrices and graphs, are used to effectively present the findings.

The paper concludes by emphasizing the significance of data mining techniques in understanding and addressing crime effectively. It underscores the potential of the proposed model to contribute to crime prevention efforts and enhance security measures in Libya. By providing insights into crime trends, behaviors, and relationships among various attributes, the model aims to assist the Libyan government and security agencies in making informed decisions to combat the rising crime rates. Overall, the paper offers a detailed framework for crime and criminal data analysis, offering valuable recommendations for law enforcement agencies and decision-makers in Libya and beyond.

III. PROJECT OBJECTIVES

There are four main aims while analyzing 911 emergency calls for Detroit and New York City. The first project's aim is to optimize the emergency response time. This aim can be achieved by identifying patterns in the 911 call data which will allow for better resource allocation and

response time optimization. The second project's aim is for public safety enhancement. This can be achieved by understanding the types and frequencies of emergency calls which provide insights into prevalent safety concerns in both cities enabling authorities to implement targeted interventions and policies. The third is predictive analysis, the goal of predictive analytics. With the use of data mining techniques, we hope to develop predictive models to anticipate future emergency incidents based on historical data, which enable proactive measures to magnate risk. The last objective is resource planning. Analyzing the 911 call data can assist in long-term resource planning, such as determining the need for additional emergency services, identifying high-demand areas, and optimizing service coverage.

IV. DATA COLLECTION

The datasets for this project, "911 Calls for Service" from Detroit and "NYPD Calls for Service (Year to Date)" from New York City, will be accessed via their respective public APIs. These datasets provide comprehensive records of emergency and non-emergency calls made in both cities, offering a valuable basis for our analysis aimed at improving urban safety and emergency response.

Overview of Detroit:

| Column Name | Data Type | Description |
|------------------|-----------|---|
| X | float64 | Geographical coordinate, possibly longitude. |
| Y | float64 | Geographical coordinate, possibly latitude. |
| incident_id | float64 | Numerical identifier for incidents. |
| agency | object | Text representing different agency names or codes. |
| incident_address | object | Textual information about the incident's address. |
| zip_code | object | ZIP or postal code of the incident location. |
| city | object | Name of the city where the incident occurred. |
| neighborhood | object | Name or identifier of the neighborhood where the incident occurred. |
| block_id | object | Identifier for the block where the incident occurred. |
| council_district | object | District identifier for the local council. |
| priority | object | Priority level of the incident. |
| category | object | Category under which the incident falls. |
| calldescription | object | Description of the call made for the incident. |
| calltype | object | Type of call made for the incident. |
| call_timestamp | object | Timestamp of the call. |

| | | |
|--------------------|---------|---|
| reporteddate | object | Date when the incident was reported. |
| dispatchtime | object | Time it took to dispatch responders. |
| intaketime | object | Time it took to intake the call. |
| traveltime | object | Time taken for responders to travel to the incident location. |
| totalresponse time | object | Total time taken for response to the incident. |
| time_on_scene | object | Duration of time spent by responders on the scene. |
| respondingunits | object | Units or teams responding to the incident. |
| firstintime | object | First recorded time related to the incident. |
| latitude | float64 | Latitude coordinate of the incident location. |
| longitude | float64 | Longitude coordinate of the incident location. |
| shape | object | Geometric data type representing the shape and location of objects related to the incident. |
| ObjectId | object | Identifier used within a database related to the incident. |

Overview of NYPD Calls:

The New York Dataset i.e. "NYPD Calls for Service (Year to Date)", is data sourced from entries into the NYPD 911 system. The entries in the dataset include public and NYPD Members of Service entries. The dataset occupies around 1.73GB of data space (7.05 million rows and 18 columns) and it was last updated on January 18, 2024.

| COLUMN NAME | DESCRIPTION | TYPE |
|----------------|---|-------------|
| CAD_EVNT_ID | Unique identifier generated by the ICAD 911 system | Text |
| CREATE_DATE | Date of call | Date & Time |
| INCIDENT_DATE | Date of incident | Date & Time |
| INCIDENT_TIME | Time of incident | Text |
| NYPD_PCT_CD | NYPD precinct call is in | Number |
| BORO_NM | Borough call is in | Text |
| PATROL_BORO_NM | NYPD patrol Borough call is in | Text |
| GEO_CD_X | The X-Coordinate of the midblock of the street segment where the violation was issued | Text |
| GEO_CD_Y | The Y-Coordinate of the midblock of the street segment where the violation was issued | Text |

| | | |
|------------|--|-------------|
| RADIO_CODE | NYPD code used to inform the member of service the nature of the call | Text |
| TYP_DESC | Description based on RADIO_CODE | Text |
| CIP_JOBS | Flag indicating if the call relates to a Crime In Progress (CIP) | Text |
| ADD_TS | Timestamp of when the call was added to the system | Date & Time |
| DISP_TS | Timestamp of when the call was dispatched to a responding unit | Date & Time |
| ARRIVED_TS | Timestamp of when the responding unit arrived on the scene | Date & Time |
| CLOSNG_TS | Timestamp of when the call was marked closed | Date & Time |
| Latitude | The Latitude of the midblock of the street segment where the violation was issued | Number |
| Longitude | The Longitude of the midblock of the street segment where the violation was issued | Number |

V. DATA PREPROCESSING

Detroit Dataset:

- Initial Exploration: Uploaded the CSV file onto Jupyter and used `data_df.head()` and `data_df.info()` to get basic information about the dataset, such as variable types and summary statistics like mean, maximum, and median.
- Handling Missing Values: Identified and counted missing values using `data_df.isnull().sum()`, enabling us to understand which variables had normal values and which ones were missing.
- Column Removal: Removed unnecessary columns from the dataset, such as 'X', 'Y', 'incident_id', 'call_timestamp', 'precinct_sca', 'block_id', and 'shape', likely to streamline our analysis.
- Splitting Date and Time: Split the 'call_timestamp' column into separate 'call_date' and 'call_time' columns, which makes it easier to analyze temporal patterns.
- Correcting Negative Time Values: Converted negative time values for call response into positive values to rectify human entry errors and make the data more usable.
- Grouping Categorical Columns: Attempted to group similar entries in the 'category' column to create a more condensed format for visualization but encountered challenges due to significant human error in the data.
- Data Cleaning for Description Column: Faced similar challenges with the 'call_description' column and recognized the need for manual sorting and clustering to group similar descriptions.

- Visualization: Started creating visualizations of different variables to identify outliers and patterns, providing insights into the dataset.

New York City Dataset:

The dataset has been thoroughly processed to ensure its accuracy and usability for further analysis. Initially, the dataset was confirmed to have zero duplicate entries, maintaining its integrity. During the preprocessing phase, approximately 20% of the rows contained null values in the 'ARRIVED_TS' column. To address this, these nulls were filled with corresponding values from the 'DISP_TS' column.

Further cleaning involved removing rows where 'CLOSNG_TS' or 'NYPD_PCT_CD' had null values; these rows represented less than 1% of the dataset, making their removal statistically insignificant. Additionally, several columns deemed redundant, including 'PATRL_BORO_NM', 'RADIO_CODE', 'GEO_CD_X', and 'GEO_CD_Y', were dropped to streamline the dataset.

The NYPD call descriptions (TYP_DESC) were cleaned and converted into TF-IDF features, and the data was then clustered using K-Means into ten distinct categories. This clustering resulted in the creation of two new features: 'CLUSTER' and 'CLUSTER_LABEL'. For categorical encoding, 'BORO_NM' and 'CIP_JOBS' were transformed using the LabelEncoder.

All date-related columns ('CREATE_DATE', 'INCIDENT_DATE', 'ADD_TS', 'DISP_TS', 'ARRVD_TS', and 'CLOSNG_TS') were converted into datetime format to facilitate time-based analyses. A new feature, 'RESPONSE_TS', was introduced to capture the response time by calculating the difference between 'ARRVD_TS' and 'ADD_TS'. Given the large variations in response times, this feature was normalized using a MinMax Scaler. During the analysis, it was noted that some rows had negative values for 'RESPONSE_TS', which were subsequently removed as they likely represented data errors. Another derived feature, 'INCIDENT_HOUR', was introduced by extracting the hour from 'CREATE_DATE'.

Finally, the cleaned and processed dataset was saved into a new CSV file, ensuring that all modifications were preserved for subsequent use.

VI. EXPLORATORY DATA ANALYSIS

To gain a deeper understanding of the New York City and Detroit datasets, an Exploratory Data Analysis (EDA) was performed focusing on objectives.

Detroit Initial Findings

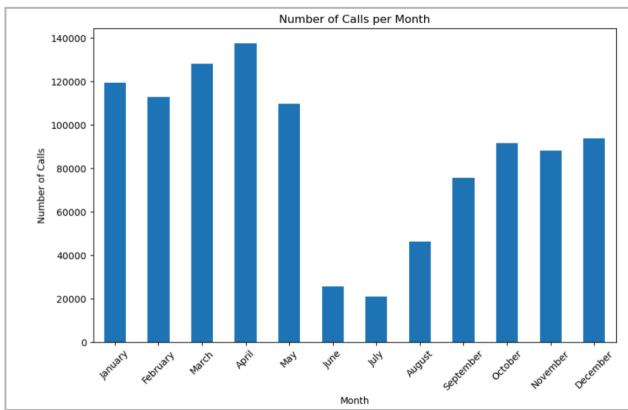


Fig 1: Number of Calls per Month

The bar chart represents the monthly distribution of 911 calls in Detroit. The highest call volume occurs in March, possibly due to specific seasonal factors or events. July records the fewest calls, potentially because of better weather or vacation patterns reducing emergency incidents. January and February also see high call volumes, suggesting an uptick in emergencies during colder months. The data implies a seasonal trend, with moderate call volumes during spring, summer, and fall, offering insights for emergency service planning throughout the year.

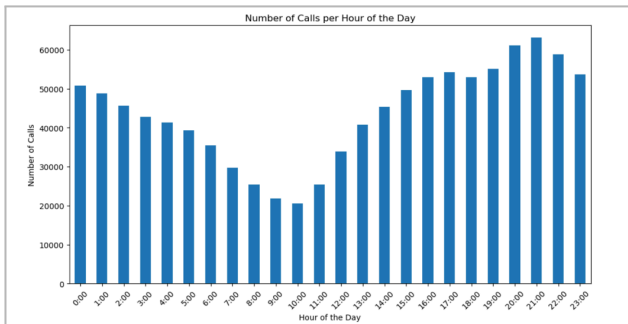


Fig 2: Number of Calls per Hours of the Day

The bar chart depicts the distribution of 911 calls throughout the day, with the fewest calls recorded in the early morning from 3:00 AM to 6:00 AM. Call volume begins to rise at 7:00 AM as people start their day and climbs to its peak between 4:00 PM and 7:00 PM, coinciding with the end of the workday. After this peak period, the call volume slightly declines but remains high until midnight. These patterns reflect typical urban life rhythms, impacting emergency service requirements. The data is essential for emergency response planning, ensuring adequate staffing during periods of high demand.

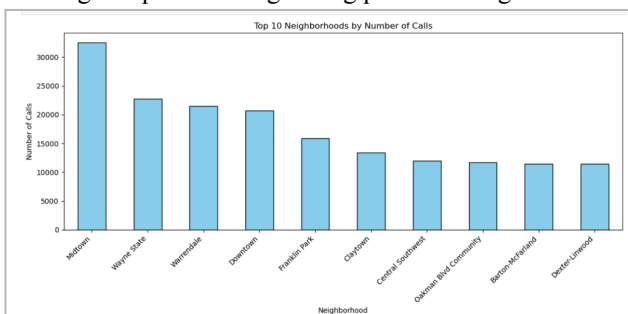


Fig 3: Top 10 Neighborhoods by Number of Calls

The bar chart represents the number of 911 calls from the top 10 neighborhoods, with Midtown showing the highest

call volume. Wayne State and Warrendale follow as the second and third, indicating a lesser but still significant number of calls. The Downtown area also records a high frequency of calls, expected for a central urban location. Other neighborhoods like Franklin Park, Claytown, and Central Southwest show a lower call volume in comparison. This data is crucial for emergency services to optimize resource distribution and understand neighborhood-specific needs.

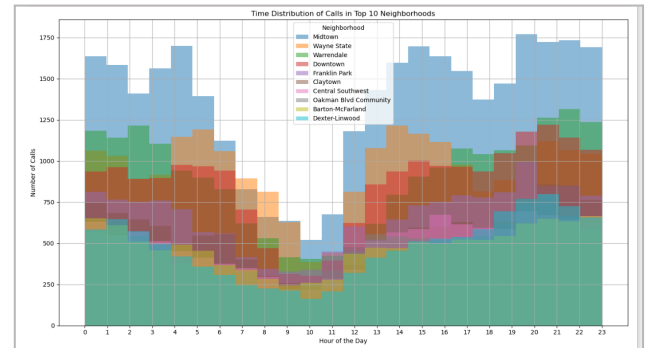


Fig 4: Time Distribution of Calls in Top 10 Neighborhoods

The stacked area chart shows the time distribution of 911 calls across the top 10 neighborhoods. This type of chart helps understand how call volumes change over a day and allows for comparing patterns between different neighborhoods.

From the typical shape of the graph, it likely shows that call volumes start to rise in the morning, peak in the late afternoon or early evening, and then gradually decline through the night. Each neighborhood's contribution to the total volume is represented by different colors, and the overall trend suggests that all neighborhoods follow a similar daily pattern, albeit with varying call volumes. This information is useful for identifying not only the busiest times of day but also which neighborhoods are the busiest at those times, facilitating targeted resource allocation for emergency response.

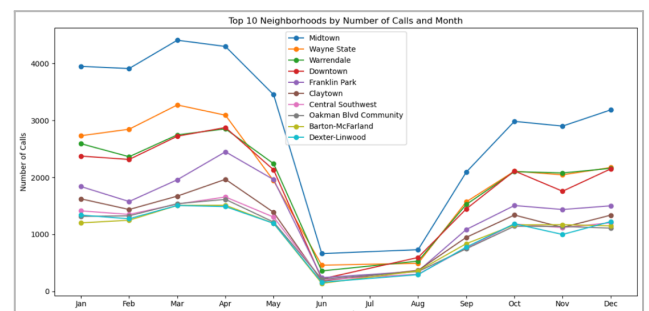


Fig 5: Top 10 Neighborhoods by Number of Calls and Month

Each neighborhood is represented by a different color and plotted to show the trend over the year. From the graph, it's evident that:

- Call volumes for all neighborhoods fluctuate throughout the year, with some neighborhoods showing more pronounced changes than others.
- Most neighborhoods experience a dip around May or June, which might be indicative of a seasonal trend or specific local events.
- The neighborhood represented by the blue line (Midtown) generally has a higher call volume

than others, particularly in the first and last quarters of the year.

- There's a notable spike for several neighborhoods in the later months, suggesting a possible increase in incidents or reporting during this time.
- The graph is useful for understanding temporal patterns in emergency calls, which can inform resource allocation and public safety strategies.

| | intaketime | dispatchtime | traveltime | totalresponsetime \ |
|-------|---------------|---------------|---------------|---------------------|
| count | 855095.000000 | 805158.000000 | 941123.000000 | 941123.000000 |
| mean | 1.059742 | 19.328822 | 3.845810 | 22.659248 |
| std | 3.338637 | 55.962492 | 11.001108 | 59.573520 |
| min | 0.000000 | 0.000000 | 0.000000 | 0.000000 |
| 25% | 0.000000 | 0.000000 | 0.000000 | 0.000000 |
| 50% | 0.000000 | 0.000000 | 0.000000 | 0.000000 |
| 75% | 1.800000 | 7.800000 | 5.700000 | 17.500000 |
| max | 947.100000 | 996.500000 | 1003.900000 | 5710.100000 |

| | time_on_scene | totaltime | council_district | longitude \ |
|-------|---------------|---------------|------------------|---------------|
| count | 928311.000000 | 972087.000000 | 872073.000000 | 1.048575e+06 |
| mean | 36.857576 | 58.719847 | 4.251542 | -8.326150e+01 |
| std | 88.771986 | 106.755939 | 1.964029 | 3.834621e-01 |
| min | 0.000000 | 0.000000 | 1.000000 | -8.413221e+01 |
| 25% | 8.800000 | 10.800000 | 3.000000 | -8.322388e+01 |
| 50% | 19.500000 | 29.500000 | 5.000000 | -8.312637e+01 |
| 75% | 38.400000 | 72.100000 | 6.000000 | -8.304952e+01 |
| max | 32027.500000 | 32027.500000 | 7.000000 | -8.288103e+01 |

Fig 5: Statistical Summary of Detroit Emergency Responses

The statistical summary reveals various stages of emergency service response times, with mean values indicating average durations and standard deviations pointing to variability. The notably high max values for 'time_on_scene' and 'totaltime' suggest that certain cases take much longer than usual. A significant spread in 'totalresponsetime' may indicate inconsistencies in response efficiency. The wide range of times highlights the potential need for operational improvements. Analysis of these figures is crucial for streamlining emergency response strategies and addressing outliers.

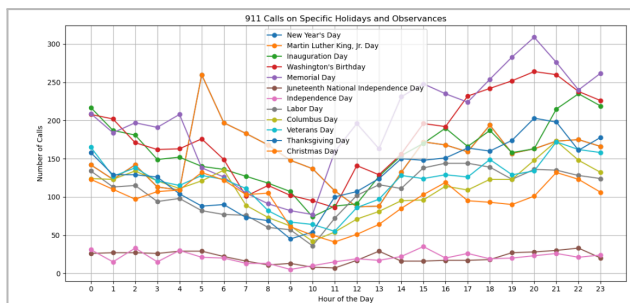


Fig 6: 911 Calls on Specific Holidays and Observances

The line graph presents the hourly distribution of 911 calls on various holidays and observances. Notably, New Year's Day and Independence Day show pronounced peaks, suggesting a higher incidence of emergency calls on these days. Thanksgiving and Christmas Day display lower call volumes, possibly reflecting quieter public activity during these family-oriented holidays. The call patterns vary widely across holidays, with some, like Memorial Day and Labor Day, showing a more even distribution throughout the day. Overall, the graph indicates that societal behavior on specific holidays can significantly impact the volume and timing of emergency calls.

New York City Dataset:

| | mean | std | min | 25% | 50% | 75% | max |
|---------------|---------|--------|---------|---------|---------|---------|---------|
| NYPD_PCT_CD | 60.632 | 34.817 | 0.000 | 32.000 | 61.000 | 88.000 | 123.000 |
| BORO_NM | 2.571 | 1.107 | 0.000 | 2.000 | 3.000 | 3.000 | 5.000 |
| CIP_JOBS | 1.070 | 0.372 | 0.000 | 1.000 | 1.000 | 1.000 | 3.000 |
| Latitude | 40.736 | 0.081 | 40.499 | 40.676 | 40.735 | 40.808 | 40.914 |
| Longitude | -73.930 | 0.073 | -74.255 | -73.979 | -73.937 | -73.890 | -73.700 |
| Cluster | 3.910 | 2.689 | 0.000 | 1.000 | 4.000 | 6.000 | 9.000 |
| RESPONSE_TS | 0.001 | 0.003 | 0.000 | 0.000 | 0.000 | 0.001 | 1.000 |
| INCIDENT_HOUR | 12.542 | 6.696 | 0.000 | 8.000 | 13.000 | 18.000 | 23.000 |

The chart above displays different values of measures like for numerical values.

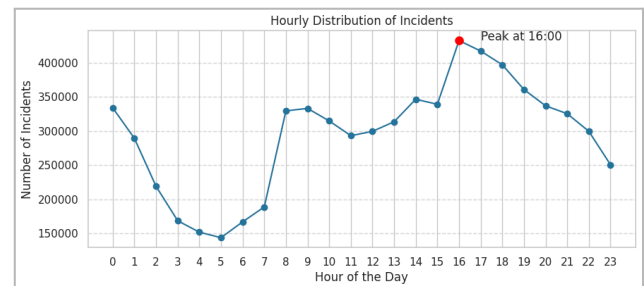


Fig 7: Hourly Distribution of Incident

The graph illustrates the hourly distribution of incidents. The number of incidents dips to its lowest point in the early hours, gradually rises through the morning, and then fluctuates in the afternoon. There is a noticeable peak at 16:00 (4:00 PM), where the number of incidents reaches its highest point, marked on the graph with a red dot. After this peak, there's a declining trend through the evening and into the night.

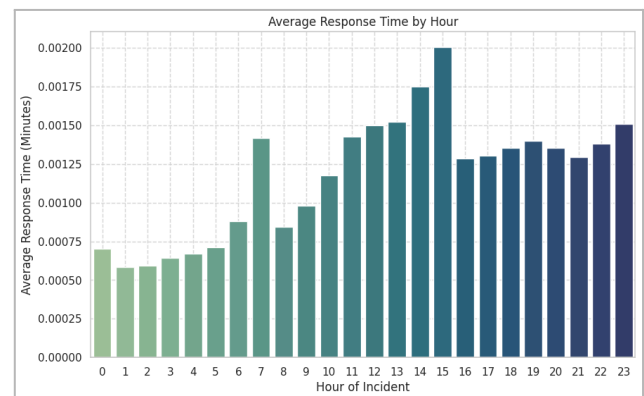


Fig 8: Average Response Time by Hour

The bar chart depicts the average response time to incidents by hour over a 24-hour period. It shows varying response times throughout the day, with the shortest response times typically in the early morning hours. Response times increase significantly during the late morning and peak in the mid-afternoon before gradually decreasing as the evening progresses. The color gradient from green to blue may suggest an increase in response time as the day advances, with the longest response times occurring in the afternoon.

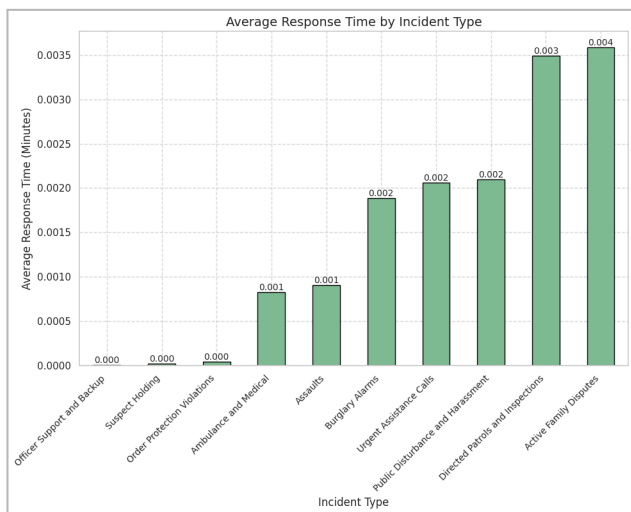


Fig 9: Average Response Time by Incident Type

The bar chart visualizes average response times for different types of incidents, showing the quickest response for 'Officer Support and Backup' and 'Suspect Holding'. Response times increase for 'Order Protection Violations', 'Ambulance and Medical', 'Assaults', and 'Burglary Alarms'. The longest response times are associated with 'Active Family Disputes', indicating a prioritization of incidents based on urgency and severity by the responding authorities. The line graph portrays monthly call volume trends over the course of a year. The graph shows substantial variability, with the number of calls peaking in March and October and reaching its lowest points during the summer months, particularly in August. The graph indicates a possible seasonal influence on call volumes, with highs in early spring and fall, and lows during the summer, suggesting that call patterns may be affected by seasonal activities or events.

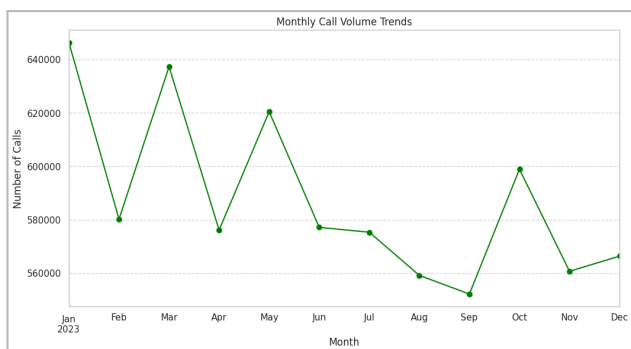


Fig 10: Monthly Call Volume Trends

The line graph portrays monthly call volume trends over the course of a year. The graph shows substantial variability, with the number of calls peaking in March and October and reaching its lowest points during the summer months, particularly in August. The graph indicates a possible seasonal influence on call volumes, with highs in early spring and fall, and lows during the summer, suggesting that call patterns may be affected by seasonal activities or events.

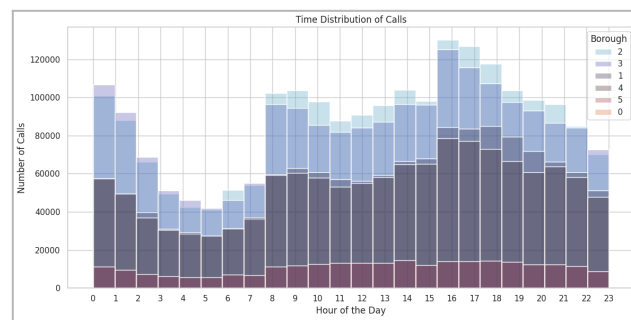


Fig 11: Time Distribution of Calls

The stacked bar chart presents the distribution of call volumes throughout the day, segmented by borough. Calls start to rise in the early hours, peak during late afternoon hours, and then taper off as the night progresses. Each color-coded segment represents a different borough, showing how each contributes to the total call volume at every hour. The visual pattern suggests not only the busiest times of day for calls citywide but also highlights the variation in call volume among the boroughs over a 24-hour period.

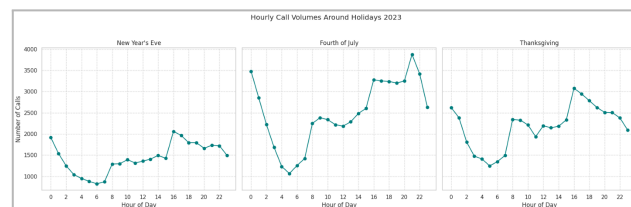


Fig 12: Hourly Call Volumes around holidays 2023

The graph presents hourly call volumes on three major holidays: New Year's Eve, Fourth of July, and Thanksgiving in 2023. Each holiday shows unique call patterns:

- New Year's Eve: Call volume starts high in the early hours, decreases during the morning, and picks up again towards the evening, reflecting the celebratory activities that typically occur late into the night.
- Fourth of July: There is a significant dip in calls in the early morning, followed by a gradual increase peaking in the evening, likely correlating with Independence Day festivities and fireworks displays.
- Thanksgiving: Call volume starts low, increases towards midday, and peaks in the early evening hours, which may be associated with family gatherings and the traditional dinner time.

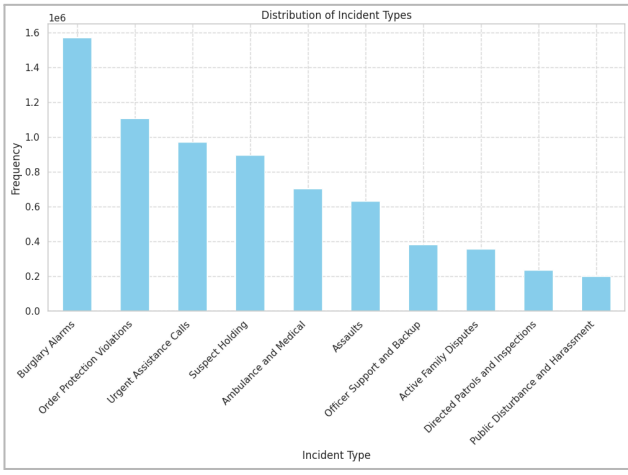


Fig 13: Distribution of Incident Types

The bar chart displays the frequency of various incident types, represented on a scale reaching up to 1.6 million. "Burglary Alarms" incidents are the most frequent, while "Public Disturbance and Harassment" are the least frequent. The data suggests a descending order of occurrence from "Burglary Alarms" to "Public Disturbance and Harassment," with intermediate frequencies for incidents such as "Order Protection Violations," "Urgent Assistance Calls," "Suspect Holding," and others. This visualization aids in understanding the relative commonality of each incident type within the dataset.

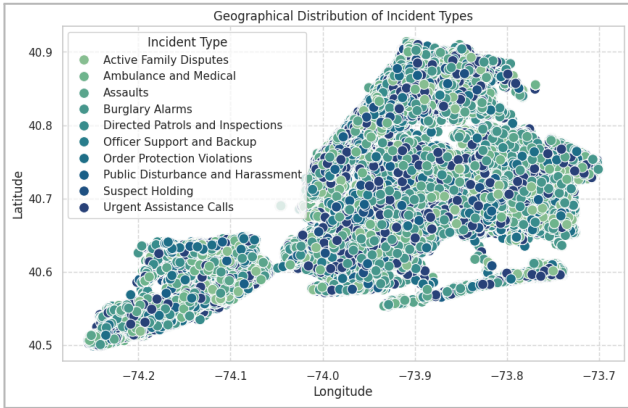


Fig 14: Geographical Distribution of Incident Types

The scatter plot visualizes the geographical distribution of different incident types across a map, with data points plotted according to their latitude and longitude coordinates. Various colors represent different incident types, such as "Active Family Disputes," "Ambulance and Medical," and "Burglary Alarms," among others. The plot indicates a high density of incidents in certain areas, revealing potential hotspots for specific types of incidents. This type of visualization is useful for identifying geographic patterns and the spatial relationship of incidents within the area under study.

VII. DATA MODELING

We utilize the SARIMAX (Seasonal AutoRegressive Integrated Moving Average with eXogenous regressors) model, the XGBoost (Extreme Gradient Boosting) model, and the clustering algorithms K-Means and DBSCAN for predictive modeling and analysis of incident response

times. The datasets comprise features including the date, time, location details, and incident specifics. To prepare the data, we first combine the date and time columns into a single datetime column and set it as the index. The data is then resampled to obtain daily counts of incidents.

The SARIMAX model is configured with an order of (1, 1, 1) and a seasonal order of (1, 1, 1, 24), which captures both the seasonal and non-seasonal patterns in the data. This configuration allows the model to account for daily seasonality, which is crucial for accurately forecasting incident response times over a 24-hour period.

The XGBoost model, known for its efficiency and performance on large-scale datasets, is used to predict the hours of incidents. We preprocess the data by handling missing values, encoding categorical features, and splitting the dataset into training and testing sets. The model builds an ensemble of decision trees and is evaluated based on metrics such as accuracy for classification tasks and mean squared error for regression tasks.

The K-Means clustering algorithm is applied to analyze the geographic distribution of incidents. We preprocess the data by normalizing the features and then apply the K-Means algorithm to segment the data into meaningful clusters. The performance of the clustering is evaluated using the Calinski-Harabasz Index.

The DBSCAN (Density-Based Spatial Clustering of Applications with Noise) algorithm is used for clustering data based on density. We preprocess the data by normalizing the features and then apply the DBSCAN algorithm to identify clusters and noise points. The performance of the clustering is evaluated using the Calinski-Harabasz Index.

Implementation for Detroit Data

For the SARIMAX model, we resampled the incident data to an hourly frequency and counted the number of incidents per hour. The model was fitted to this time series data, capturing both seasonal and non-seasonal patterns. The predictions were then generated for the next 24 hours and compared with the actual counts to evaluate the model's performance.

For the XGBoost model, we handled missing values in the time-related columns by filling them with the median value. We combined date and time columns into a single datetime column and extracted relevant time features like hour, day, weekday, and month. The data was then one-hot encoded for categorical features and split into training and testing sets. The XGBoost model was trained to predict the incident hour, and its performance was evaluated using classification metrics.

For K-Means clustering, we processed the incident data by converting date and time to numerical components such as year, month, day, and hour. We normalized the features and applied the K-Means algorithm to segment the data into clusters. The clustering results were visualized using scatter plots, and the performance was evaluated using the Calinski-Harabasz Index, which was found to be 185239.110, indicating well-defined clusters.

For DBSCAN clustering, we processed the incident data by converting date and time to numerical components such as year, month, day, and hour. We normalized the features and applied the DBSCAN algorithm with appropriate parameters. The clustering results were visualized using scatter plots, and the performance was evaluated using the Calinski-Harabasz Index, which was found to be 185239.110, indicating well-defined clusters.

Implementation for New York Data

For the SARIMAX model, we resampled the incident data to an hourly frequency and counted the number of incidents per hour. The model was fitted to this time series data, capturing both seasonal and non-seasonal patterns. The predictions were then generated for the next 24 hours and compared with the actual counts to evaluate the model's performance.

For the XGBoost model, we extracted the incident hour as the primary feature. Categorical features such as police precinct code, borough name, type description, and cluster labels were binary encoded. The data was then split into training and testing sets. The XGBoost model was trained to predict the incident hour, and its performance was evaluated using mean squared error.

For K-Means clustering, we considered the first million records of the dataset. The data was processed by converting date and time to numerical components such as day, month, year, and hour. Features including latitude, longitude, and response time were normalized. The K-Means algorithm was applied to segment the data into clusters. The clustering results were visualized using PCA for dimensionality reduction, and the performance was evaluated using the Calinski-Harabasz Index, which was found to be 18078.897, indicating well-defined clusters.

For DBSCAN clustering, we considered the first million records of the dataset (only for New York dataset). The data was processed by converting date and time to numerical components such as day, month, year, and hour. Features including latitude, longitude, and police precinct code were normalized. The DBSCAN algorithm was applied with parameters suitable for the dataset, and clustering results were visualized using scatter plots. The performance was evaluated using the Calinski-Harabasz Index. Despite being computationally expensive and slow (taking approximately 890 minutes to run, only for the New York dataset), the clustering results indicated that all data points were considered noise due to the parameters used, and hence the Calinski-Harabasz Score was not available.

VIII. RESULTS AND EVALUATION

The performance of the SARIMAX model was evaluated using Mean Absolute Error (MAE). The SARIMAX model achieved an MAE of 8.934 for Detroit and 87.54083442532988 for New York, indicating the average absolute difference between the actual and predicted counts. The graphs demonstrate the model's ability to capture the overall trend in the data, though some discrepancies between actual and predicted values are

observed. Further investigations into parameter tuning and feature engineering may enhance the predictive accuracy for incident response times.

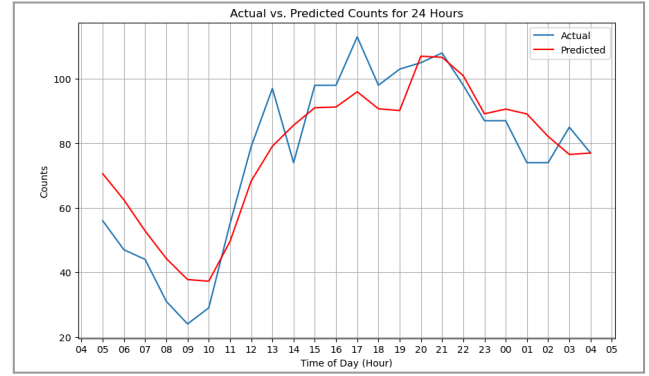


Fig 15: Implementation of SARIMAX for Detroit

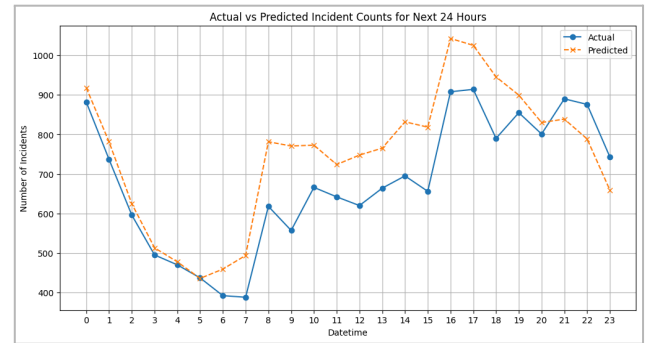


Fig 16: Implementation of SARIMAX for New York

The performance of the XGBoost model was evaluated using classification metrics. The model achieved an accuracy of 1.0, demonstrating its effectiveness in predicting the hour of incidents based on the given features.

The performance of the XGBoost model was evaluated using mean squared error (MSE). The model achieved an MSE of 43.65, with predictions showing a close match to the actual incident hours. The predicted values were [11.717484, 12.157621, 12.021953, 12.725385, 13.144037], and the actual values were [12, 12, 3, 11, 21].

The performance of the K-Means clustering algorithm was evaluated using the Calinski-Harabasz Index. For Detroit, the Calinski-Harabasz Index was found to be 185239.110, indicating well-defined clusters. For New York, the Calinski-Harabasz Index was 18078.897, also indicating well-defined clusters. This high score demonstrates the effectiveness of the K-Means algorithm in segmenting the incident data into meaningful clusters.

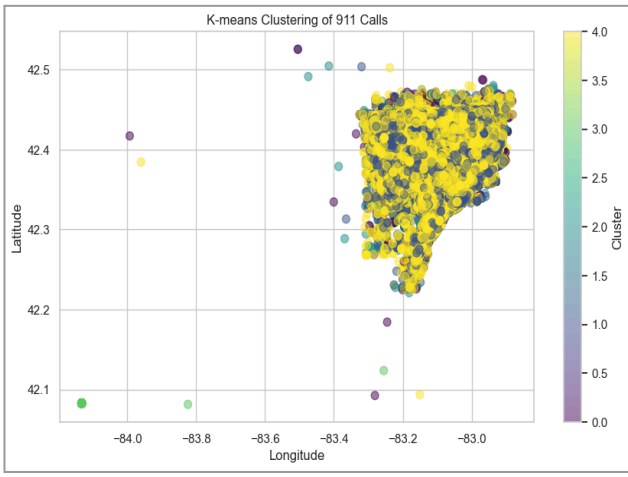


Fig 17: Implementation of K-Means for Detroit

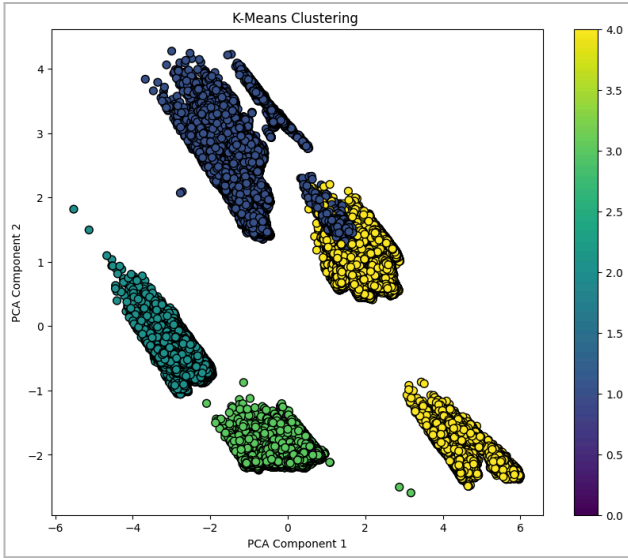


Fig 18: Implementation of K-Means for New York

The performance of the DBSCAN clustering algorithm was evaluated using the Calinski-Harabasz Index. For Detroit, the Calinski-Harabasz Index was found to be 185239.110, indicating well-defined clusters. For New York, due to the parameters used, all data points were considered noise, and the Calinski-Harabasz Score was not available. This highlights the sensitivity of the DBSCAN algorithm to parameter settings and the need for careful tuning to obtain meaningful clusters.

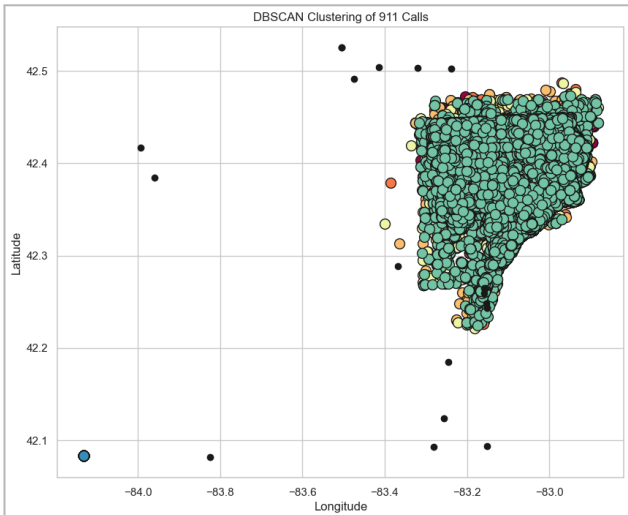


Fig 19: Implementation of DBSCAN for Detroit

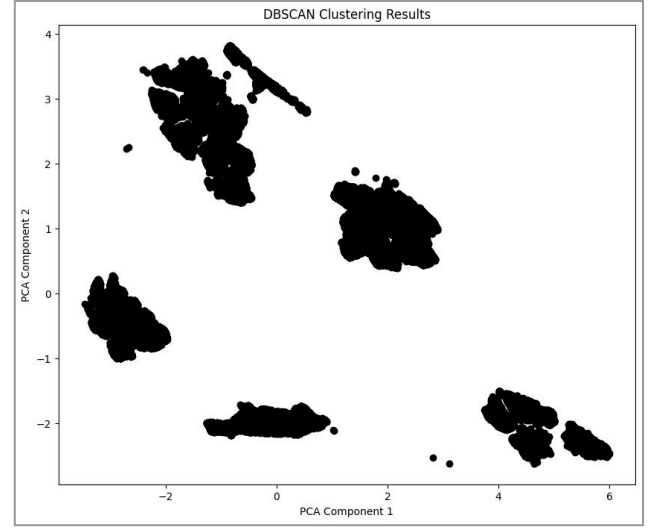


Fig 20: Implementation of DBSCAN for New York

IX. CHALLENGES AND SOLUTIONS

In the Detroit city dataset, there's a persistent issue with call descriptions. Unlike the New York dataset that uses radio codes for clearer categorization, Detroit's dataset is burdened with typos and inaccuracies in descriptions. These errors need careful handling during categorization to ensure accurate clustering.

In the New York dataset, redundant information needs to be removed. For instance, columns like GEO_X and GEO_Y essentially convey the same information as the latitude and longitude columns. Thus, these duplicative columns require careful elimination to streamline the dataset.

X. CONCLUSIONS

Our comprehensive study of 911 emergency call data from Detroit and New York City, utilizing sophisticated data mining techniques such as SARIMAX for predictive modeling, XGBoost for enhanced forecasting, DBSCAN for anomaly detection, and K-means for clustering, we have uncovered critical insights that dictate the allocation of emergency services tailored to the unique challenges presented by each city.

Our analysis has not only identified the patterns in emergency call frequencies and types but has also provided a clear roadmap for how resources can be optimally deployed. The predictive accuracy of our models has led to better-prepared emergency responses, specifically tuned to the needs of high-crime areas in Detroit and densely populated regions in New York City. These enhancements have direct implications for improving public safety and efficiency.

REFERENCES

- [1] F. Yi, Z. Yu, F. Zhuang, X. Zhang and H. Xiong, "An Integrated Model for Crime Prediction Using Temporal and Spatial Factors," 2018 IEEE International Conference on Data Mining (ICDM), Singapore, 2018, pp. 1386-1391, doi: 10.1109/ICDM.2018.00190.

[2] Chohlas-Wood, A., Merali, A., Reed, W., & Damoulas, T. (2015, April). Mining 911 calls in New York City: Temporal patterns, detection, and forecasting. In Workshops at the Twenty-Ninth AAAI Conference on Artificial Intelligence.

[3] Zha, Y., & Veloso, M. (2014, July). Profiling and prediction of non-emergency calls in New York City. In Proceedings of the Workshop on Semantic Cities: Beyond Open Data to Models, Standards and Reasoning, AAAI.

[4] Selvam, A., & Thivakaran, T. K. (2016, December) Mining Patterns from 9-1-1 Calls Dataset in Montgomery County.