



**Middlesex
University
London**

CST4090 Individual Project

**Crime Prediction and Mapping in London: Integrating Historical,
Socioeconomic, Mobility, and Geospatial Data**

Student Name: Mohammed Abdul Haseeb

Student ID: M00960849

Supervisor: Dr. Olugbenga Oluwagbemi

ACKNOWLEDGEMENT

I would like to express my profound thanks to my supervisor, Dr. Olugbenga Oluwagbemi, for his unwavering support, ongoing direction, and advice throughout my pursuit of an MSc degree as well as during the competition of this project. The progress and effective completion of this project were greatly aided by his knowledge and judgement.

My sincere gratitude goes out to my family, friends, and parents in particular, whose unwavering belief in me kept me hopeful and motivated me to achieve this wonderful goal.

ABSTRACT

The problem of crime exists in all city settings. However, it is especially pronounced in urban centres, with far-reaching effects on public safety, the formulation of policy, and the deployment of law enforcement. For this project, a range of publicly available data was obtained, such as historical crime figures, socioeconomic data from the 2021 census, weather, transport foot traffic, and geospatial data. Using this wealth of data enabled the performance of both spatial and temporal analyses on the occurrence of crime and its prediction. The analysis was directed at the ward level of London to uncover any small area patterns.

An extensive dataset has been constructed that covers the period from 2014 to 2024. This dataset was built by gathering data from a variety of sources. The research and analysis left out the COVID-19 period from February 2020 to March 2021 because of its anomalous nature. The reason for doing so was based on the understanding that the kind of conditions that humans found themselves in during a pandemic would skew data in a way that would not allow for an accurate representation of what was really going on. Socioeconomic data, like health, education, and housing conditions, was blended with foot traffic and other public transport data from major stations in London. Additionally, ward-level crime trends and possible crime hotspots have been visualised using geospatial data.

Exploratory Data Analysis (EDA) identified the seasonal patterns in the data. Warmer months generally exhibit higher crime rates. Additionally, the correlation between crime rates and specific socioeconomic factors was analysed, and a non-linear relationship was established between factors like unemployment, poor health, and overcrowded households. The most significant revelation from this analysis was the strong influence of foot traffic data on crime, especially in places where high economic activity and high population density intersect.

The geospatial analysis provided more profound insight into how crime is spread across different wards and pinpointed several likely hotspots emerging around specific transport hubs. The analysis also highlighted a few areas where the crime rate is persistently high. It was also discovered that wards in proximity to other wards with higher crime rates tend to have similarly high crime.

Different models, such as Neural Networks, Gradient Boosting, and Random Forest, were explored to predict future crime trends. The Neural Network model was chosen as the final model because of its excellent performance on the data. A high R² score of 0.95 and mean absolute error (MAE) of 18.9 were obtained. The most significant factors were foot traffic data and lagged crime features.

The project's suggestions are given at the end and aim to have an impact on the tactics law enforcement uses in deploying officers and influence policymaking and procedures by considering known patterns and types of crimes that are seasonal or vary according to socioeconomic conditions. The study notes a few limitations and a couple of significant strengths. One limitation is the static nature of the census data. Future studies in the field may use real-time data sources and refined predictive models to improve accuracy.

Table of Contents

<i>ACKNOWLEDGEMENT</i>	2
<i>ABSTRACT</i>	3
<i>1 INTRODUCTION</i>	6
<i>2 BACKGROUND</i>	8
2.1 Rationale for Topic Choice	8
2.2 State of the Art in Crime Analysis:	9
2.3 Gap in Current Research:	11
<i>3 METHODOLOGY</i>	12
3.1 Data Collection	12
3.2 Data Preparation	14
3.3 Feature Engineering.....	15
3.4 Exploratory Data Analysis (EDA)	16
3.5 Geospatial Analysis	16
3.6 Modelling Approach.....	17
3.7 Cross-Validation and Hyperparameter Tuning	18
3.8 Geospatial Analysis and Visualization.....	19
<i>4 RESULTS AND EVALUATION</i>	20
4.1 Exploratory Data Analysis (EDA)	20
4.2 Correlation Analysis.....	23
4.3 Geospatial Analysis	24
4.4 Model Performance and Evaluation.....	29
4.5 Policy Implications	31
<i>5 DISCUSSION AND LIMITATIONS</i>	32
5.1 Discussion.....	32
5.2 Limitations	33
5.3 Implications for Future Research	34
5.4 Recommendations for Law Enforcement and Policymakers	34
<i>6 CONCLUSION</i>	36
<i>REFERENCES</i>	37
<i>APPENDIX A RESEARCH ETHICS SCREENING FORM</i>	39

LIST OF FIGURES

Figure 1: Overall Crime trends in London (2014-2024)	6
Figure 2: Preview of Crime Data.....	12
Figure 3: Preview of Census Data.....	13
Figure 4: Preview of Weather Data	13
Figure 5: Preview of Transport footfall data	14
Figure 6: Station Locations in London	15
Figure 7: Random Forest Model	17
Figure 8: Gradient Boosting Model	18
Figure 9: Neural Network Model	18
Figure 10: MSE and RMSE formula.....	19
Figure 11: Crime In London over time with COVID-19 Period highlighted	20
Figure 12: Total crime by Seasons.....	21
Figure 13: Top 10 Wards by Total Crime	21
Figure 14: Bottom 10 Wards by Total Crime	22
Figure 15: Top 5 Boroughs with Most Crime	22
Figure 16: Map Showing Total Crime By Ward.....	23
Figure 17: Map showing Crime Rate per Population Density	25
Figure 18: Crime Hotspots by KDE.....	26
Figure 19: Graph of Elbow Method.....	26
Figure 20: Clustering by K-means.....	27
Figure 21: Agglomerative Clustering.....	27
Figure 22: Moran's I Scatter Plot.....	28
Figure 23: Local Moran's I Cluster Map (LISA).....	28
Figure 24: Scatter Plot showing Actual VS Predicted Crime	30
Figure 25: Residual Plot	30
Figure 26: Map showing actual and predicted values of the wards in London.....	31

1 INTRODUCTION

The big city crime problem is a tough nut to crack for urban centres everywhere. It not only influences the personal safety and comfort of millions of citizens but also touches upon the very foundation and presence of urban life—be it the local economy or the sense of community that anything remotely close to a neighbourhood can provide. In a massive urban centre like London, with its sprawling expanse and enormous, often unwieldy, diversity, not to mention a complex array of overlapping, sometimes competing, public and private interests, the very nature of crime and how to prevent it must be understood at several different levels.

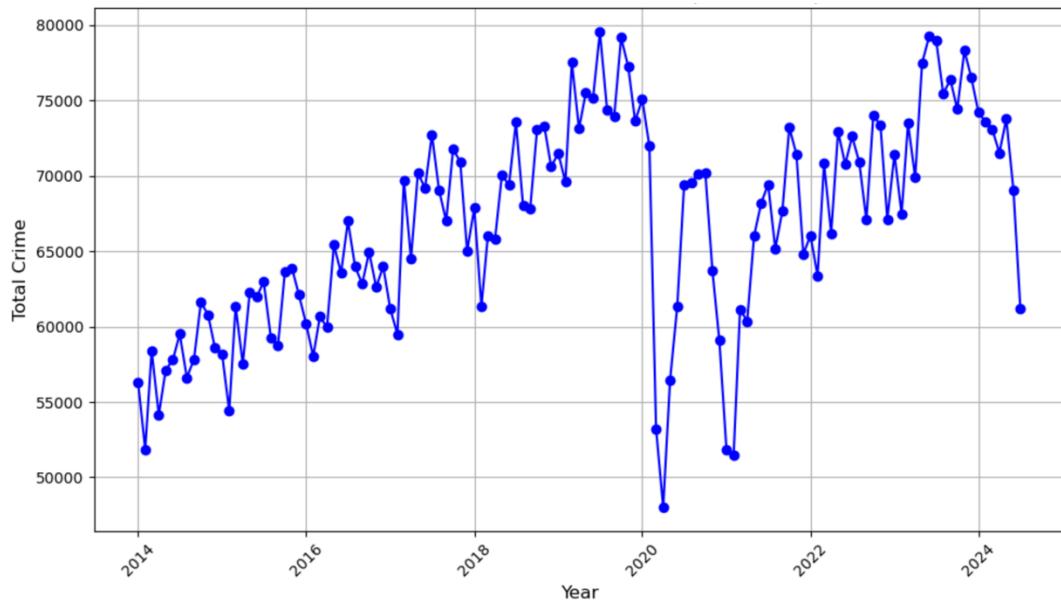


Figure 1: Overall Crime trends in London (2014-2024)

With the extensive set of available data now, which have a variety of sources that include not only crime records but also the usage of public transport, socioeconomic data, and environmental variables, there exists the potential to take a more fine-grained and data-driven approach to understanding crime. This project will attempt to capitalize on those varied and rich data sources to perform a detailed analysis of crime at the ward level in London. The advantage of using wards in this project is that wards are smaller and more localized units of analysis than boroughs, so examining crime at this level allows for the identification of specific hotspots.

Many researchers have recognized that crime is driven by a complex interaction of multiple factors, many of which are not directly related to human behaviour. For instance, social inequality has long been known to correlate with criminal activity, but it is not the only factor; human mobility, environmental conditions, and geography also have essential roles to play. For example, it is known that people commit crimes at public transportation hubs because numerous people passing through equals numerous opportunities for all sorts of offences. However, why do some people sometimes commit certain crimes in those places? Why does the whole setup work as a centrifugal force to pull people into committing crimes? This study tries to answer those questions.

This project seeks to foresee trends in criminal activity at the ward level in London. It combines historical crime data with census, weather and footfall data to produce future forecasts. The analysis uses different machine learning models like random forest and neural networks while also applying geospatial analysis to uncover the spatial distribution of crime hotspots.

This research contributes to the growing field of predictive crime analysis by filling several holes in the current body of work. Prior studies have mainly tackled the analysis at the borough level, missing the subtleties of the presence and absence of crime at the scale of wards. The understanding of this phenomenon at such a localized level can help law enforcement agencies not only predict where crimes might occur but also aid in determining the kinds of crimes that are likely to happen, thereby acting in a prevention capacity.

These insights are vital not just out of an interest in basic research but for the development and implementation of effective crime prediction and prevention strategies. Little law enforcement can help to maintain public safety in the absence of a timely, pinpointed, and granular allocation of resources that allows officers of the law to intervene before an up-and-coming problem area really starts to bubble over. Just as important is an understanding of the dynamics of crime and the underlying socioeconomic factors driving them because that helps to inform the design of the actual programs and policies that get at the root causes of criminal behaviour and boost public safety. This study provides an ostensibly straightforward impact analysis of the spatial and temporal dimensions of crime.

This report conducts a close and detailed examination of crime in London. It does not simply look at the various indicators of crime that are available but explores and explains the methods by which these indicators were prepared for analysis. It then delves into the results of exploratory data analysis and correlation studies, followed by an evaluation of various machine learning models for predicting crime trends. Finally, the report concludes with a discussion of key findings, their implications for law enforcement and policymakers, and recommendations for future research and policy interventions.

2 BACKGROUND

2.1 Rationale for Topic Choice

Urban planners, legislators, and law enforcement organisations have long placed a high priority on preventing and reducing crime. To maintain public safety and raise living standards, it is even more essential to comprehend and anticipate crime trends as cities become more complicated. The drivers of crime are complex and challenging to comprehend through conventional analysis in large metropolitan centres like London, where a variety of social, economic, and environmental factors interact. Because of these intricacies, preventing crime can be extremely difficult, especially when solutions are generic rather than customised to the particulars of each community.

The necessity for localised, data-driven insights is one of the main reasons that London is concentrating on ward-level crime analysis. While helpful, borough-wide crime data may mask more complex crime dynamics at more minor geographic scales, which is why previous research and police tactics frequently concentrate on it. Due to variables including socioeconomic inequality, high population density, and closeness to transport hubs, certain wards suffer from more frequent and severe crime events than others. The goal of this study is to identify these regional differences through ward-level analysis and to produce comprehensive findings that can directly influence resource allocation and policing tactics.

A new opportunity to create sophisticated crime prediction models is presented by the increasing availability of free data from public institutions. One place where a multitude of publicly accessible data may be combined to provide a comprehensive picture of urban crime dynamics in London. This data includes census information, crime statistics, weather reports, and the number of people using public transport. Using these various datasets instead of depending just on past crime statistics enables a more advanced understanding of the factors that contribute to crime. Predictive modelling requires further investigation into the integration of socioeconomic, meteorological, and mobility data in crime analysis.

The growing accuracy and significance of predicted crime models driven by machine learning algorithms is a significant justification for this endeavour. New avenues for detecting nonlinear correlations between crime and exogenous factors have been made possible by developments in machine learning techniques, which have improved the accuracy and resilience of forecasts. Comprehending intricate relationships among variables like foot traffic, socioeconomic status, and meteorological patterns, machine learning models can forecast the potential impact of these variables on crime rates. Anticipating future patterns in criminal activity allows law enforcement to take a proactive stance in preventing crimes by better-allocating resources and acting before problems arise.

This research attempts to identify high-risk areas where targeted social interventions (e.g., community outreach, housing improvements, and job creation programs) could be significant in the long run to reduce crime by focussing on socioeconomic conditions at the ward level. Policymakers can target the underlying causes of criminal activity rather than just its symptoms by knowing how specific socioeconomic indicators are correlated with crime.

Human mobility patterns and crime are intimately related, particularly in cities like London, where millions of people pass through transport hubs daily. Because of the concentration of people in these regions and the potential for criminal activity, such as theft or violent

confrontations, high-traffic sites, such as train stations and major shopping districts, are frequently targets for crime. This research makes the study of crime more dynamic by integrating foot traffic data from public transit into the analysis. A significant insight into the temporal and spatial patterns of crime in crowded regions is to understand the relationship between people's movements and crime incidents. With the use of this method, new crime hotspots can be pinpointed more precisely for interventions.

To enable dynamic reactions to projected crime peaks, the project's data may help determine where and when police resources should be deployed. In addition, policymakers can utilise the results to create long-term plans that tackle the root causes of crime, like social injustice and economic hardship. This project's emphasis on both immediate reactions and structural adjustments is in line with modern approaches to crime prevention that prioritise data-driven decision-making.

As a result, the necessity for more precise, data-driven methods of crime analysis in urban environments forms the foundation of this study's reasoning. The aim of this research is to provide actionable insights for crime prevention and policy development in London by analysing ward-level crime patterns and incorporating a wide range of contributing factors. This research is timely and important for tackling the complexity of urban crime in the modern period, given the developments in machine learning and the growing accessibility of varied datasets.

2.2 State of the Art in Crime Analysis:

The integration of machine learning models into crime analysis has brought significant advancements. These models can process huge amounts of data and reveal intricate relationships that exist between crime and a multitude of external factors. Here are some of the important findings from the literature that underlines the importance of adopting a multi layered approach to crime prediction by incorporating socioeconomic, metrological and mobility data. Listed below are key insights form the existing body of research:

- **Socioeconomic Influences on Crime:** The relationship between crime and socioeconomic factors is well established. Unemployment, poverty, and low educational attainment are undeniably linked to crime. A number of studies, including those conducted by Reiner et al. (2018) and Chainey and Ratcliffe(2020), have shown that areas with high levels of these adverse sociodemographic conditions not only have elevated crime rates but also tend to have crimes of a more violent nature. However, a more recent study led by Lopez and Sanchez(2021) finds that this purported relationship is not as clear-cut as it seems. In their work, Lopez and Sanchez attempt to give a more nuanced presentation of the socioeconomic-crime relationship; they noted that the majority of these studies are at the borough level, with a few concentrating at more localized levels like wards and exhibiting a different pattern of socioeconomic-driven crime.

- **Housing and Crime:** Studies show that the density and quality of housing have an effect on crime. Although fewer studies have been directed at this issue, Abrams and Meyer (2019) have made a start on them. They suggest that the quality of housing in terms of maintenance and the number of people living in a unit is a factor in crime

rates. They found that "overcrowded and poorly maintained housing conditions contribute to higher crime rates, particularly property crimes and anti-social behaviours." Lower crime rates are found in well-maintained, less-densely populated areas. Edwards et al. (2021) have taken the discussion further. They find that social dynamics and the architectural design of certain "social housing estates" make them likely candidates for "hotbeds" of criminal activities.

- **Weather and Crime Seasonality:** Weather and Crime Seasonality: More and more evidence are pointing to the idea that weather, and especially temperature, has an effect on crime. For example, Lin and colleagues (2017) showed that more violent crime happens when the weather is hotter, supporting the long-held belief that heat makes people violent. A more recent study by Anderson and Cole (2022) looked into and found that the effect of temperature on property crime was even more pronounced. In colder months, property crimes were more frequent while other violent crimes spiked during warmer months. However, these studies didn't incorporate temperature data with other factors that might drive crime to see if maybe some of those other factors are also affected by warming weather.
- **Human Mobility and Crime Concentration:** Research has shown that crime is concentrated in certain high-traffic areas. These studies used footfall data—the estimated number of people entering a given space—that is available through various sensors (like those in stores) and cameras in urban environments. For instance, Chen et al. (2020) analysed footfall data around a transport hub and found that the areas with more people (especially those around train stations) experienced significantly more crime. Attached to this study was a theoretical model of how and why this might happen, which included "intensity," "space," and "time" as factors. Wu and Zhang (2021) did similar work but with a different "hotspot" analysis framework and different data. They also found strong pedestrian areas to be at higher risk. While these studies do a great job in analysing footfall in isolation, they fail to integrate it into a comprehensive model accounting for socioeconomic or weather factors.
- **Predictive Modelling Techniques:** Predictive modelling techniques are well-established in the use of machine learning for crime prediction. The most commonly used models in this context are the Random Forest, Gradient Boosting, and Neural Network models. Gorr and Harries (2020) noted that the Random Forest model has shown robust performance on large, noisy datasets and has proven very useful for capturing the non-linear relationships between the many variables involved in crime prediction. However, the Random Forest model has also shown difficulties in reliably predicting crime over time and often struggles with time-dependent variables, which is presumably a key aspect of any useful crime forecast model. As the work by Lee et al. (2022) shows, Neural Networks (and their close relative, the Deep Learning model) have a much better track record in predicting crime patterns over time by

successfully integrating socioeconomic and weather data. Nonetheless, many studies fail to integrate mobility data, such as footfall, into the model, which can prove to be a critical factor in accurate forecasting.

2.3 Gap in Current Research:

The existing crime prediction research has some glaring deficiencies, as shown by the literature reviews. Previous research often focused on borough-wide crime data, which, while useful, may obscure more nuanced crime dynamics at smaller geographical scales. Although numerous studies have investigated the intricate relationship between crime and many of the well-known socioeconomic variables (such as unemployment and housing density), they have rarely attempted to combine these variables with environmental data and real-time human mobility patterns such as foot traffic. Brown et al. (2019) and Chen et al. (2020) have done this to a limited extent and have highlighted the value of using footfall data in the models. However, these authors consider foot traffic primarily in the context of proximity to transport hubs and fail to account for the potential combined influence of weather and socioeconomic factors.

In addition, machine learning models like Random Forests and Neural Networks have been used for crime prediction. Still, as Zhang et al. (2020) pointed out, the integration of mobility data, such as station footfall and environmental variables, remains underexplored. This project developed a way to combine crime, footfall, weather and ward-level socioeconomic data that improves the model's predictive accuracy.

Instead of using Random Forests, a more sophisticated model was applied: a neural network. Neural networks are better at capturing non-linear relationships among variables, expecting it to lead to improved forecasts of future crime.

3 METHODOLOGY

This project is structured around a systematic process of data collection, preparation, analysis and modelling. It is the combination of historical crime, socioeconomic, weather, footfall and geospatial data to predict crime trends in London at the ward level. Data collection, data cleansing, feature engineering, Exploratory Data Analysis, modelling, and geographic integration are some of the steps in this approach.

3.1 Data Collection

The project integrates data from multiple sources, each enhancing certain aspects of crime prediction:

- **Crime Data:** Two datasets comprising crime records from London wards were obtained from UK metropolitan police data (https://data.london.gov.uk/dataset/recorded_crime_summary). One dataset comprises historical data from 2014 to 2022. The other dataset comprises recent crime statistics from the past 24 months. They encompass categories of crime and sites of incidents. The data's granularity, accessible at the ward level, facilitates a comprehensive analysis.

MajorText	MinorText	WardName	WardCode	LookUp_BoroughName	202208	202209	202210	202211	2023
ARSON AND CRIMINAL DAMAGE	ARSON	Abbey	E05014053	Barking and Dagenham	0	0	0	0	0
ARSON AND CRIMINAL DAMAGE	CRIMINAL DAMAGE	Abbey	E05014053	Barking and Dagenham	7	4	5	7	
BURGLARY	BURGLARY - RESIDENTIAL	Abbey	E05014053	Barking and Dagenham	0	0	0	0	
BURGLARY	BURGLARY BUSINESS AND COMMUNITY	Abbey	E05014053	Barking and Dagenham	1	1	1	1	
BURGLARY	BURGLARY IN A DWELLING	Abbey	E05014053	Barking and Dagenham	4	4	1	1	
DRUG OFFENCES	POSSESSION OF DRUGS	Abbey	E05014053	Barking and Dagenham	5	10	14	12	
DRUG OFFENCES	TRAFFICKING OF DRUGS	Abbey	E05014053	Barking and Dagenham	0	1	0	1	
FRAUD AND FORGERY	FRAUD AND FORGERY	Abbey	E05014053	Barking and Dagenham	0	0	0	0	
MISCELLANEOUS CRIMES AGAINST SOCIETY	MISC CRIMES AGAINST SOCIETY	Abbey	E05014053	Barking and Dagenham	1	1	0	1	
POSSESSION OF WEAPONS	POSSESSION OF WEAPONS	Abbey	E05014053	Barking and Dagenham	0	2	0	0	
PUBLIC ORDER OFFENCES	OTHER OFFENCES PUBLIC ORDER	Abbey	E05014053	Barking and Dagenham	1	0	1	1	
PUBLIC ORDER OFFENCES	PUBLIC FEAR ALARM OR DISTRESS	Abbey	E05014053	Barking and Dagenham	3	2	4	3	

Figure 2: Preview of Crime Data

- **Socioeconomic data:** Ward-level socioeconomic variables were sourced from the most recent 2021 census from the Office of National Statistics (<https://www.ons.gov.uk/releases/warddataenglandandwalescensus2021>). Metrics such as accommodation type, Disability, economic activity, employment History, General health, household composition, housing density, occupancy rating, population density, qualifications and tenure were selected based on their relevance to crime patterns identified in the extensive literature review.

Electoral wards and divisions Code	Electoral wards and divisions	Tenure of household (9 categories)	Observation
E05009317	Bethnal Green	Does not apply	0
E05009317	Bethnal Green	Owned: Owns outright	562
E05009317	Bethnal Green	Owned: Owns with a mortgage or loan	1120
E05009317	Bethnal Green	Shared ownership: Shared ownership	149
E05009317	Bethnal Green	Social rented: Rents from council or Local Authority	1904
E05009317	Bethnal Green	Social rented: Other social rented	1254
E05009317	Bethnal Green	Private rented: Private landlord or letting agency	2025
E05009317	Bethnal Green	Private rented: Other private rented	212
E05009317	Bethnal Green	Lives rent free	25
E05009318	Blackwall & Cubitt Town	Does not apply	0
E05009318	Blackwall & Cubitt Town	Owned: Owns outright	750
E05009318	Blackwall & Cubitt Town	Owned: Owns with a mortgage or loan	1511
E05009318	Blackwall & Cubitt Town	Shared ownership: Shared ownership	221
E05009318	Blackwall & Cubitt Town	Social rented: Rents from council or Local Authority	356

Figure 3: Preview of Census Data

- **Weather Data:** The monthly weather data between 2014 and 2024 was obtained from the UK Met Office (<https://www.metoffice.gov.uk/research/climate/maps-and-data/historic-station-data>). It includes maximum and minimum temperatures, rainfall and air frost days. It is critical for capturing seasonal patterns of crime.

Updated_Weather_Data_2014_2024

yyyy	mm	tmax	tmin	af	rain	sun
2014	1	10.0	3.8	1	162.4	68.3
2014	2	10.6	4.4	0	89.8	91.7
2014	3	14.1	4.4	1	27.8	161.4
2014	4	16.1	7.5	0	58.0	156.9
2014	5	18.0	9.8	0	84.6	178.8
2014	6	22.1	12.5	0	40.8	220.2
2014	7	25.8	15.0	0	50.0	246.4
2014	8	21.7	12.7	0	97.6	183.6
2014	9	21.5	12.8	0	10.8	134.6

Figure 4: Preview of Weather Data

- **Station Footfall:** Transport for London (TFL) provided the footfall data for public transport (<https://crowding.data.tfl.gov.uk>), which included entry and exit counts at the main stations of the National Rail and London Underground between 2020 and 2023. With a focus on the areas surrounding transport hubs, this dataset provides insights into human mobility patterns, particularly those around transport hubs that tend to attract higher crime activity.

1	Mode	MNLc	MASC	Station	Coverage	Source	En/Ex
2	LU	500	ACTu	Acton Town	Station entry/exit	Gateline	4,823,835
3	LU	502	ALDu	Aldgate	Station entry/exit	Gateline	6,897,314
4	LU	503	ALEu	Aldgate East	Station entry/exit	Gateline	10,947,896
5	LU	505	ALPu	Alperton	Station entry/exit	Gateline	2,598,605
6	LU	506	AMEu	Amersham	Station entry/exit	Gateline	1,729,521
7	LU	507	ANGu	Angel	Station entry/exit	Gateline	12,258,814
8	LU	508	ARCu	Archway	Station entry/exit	Gateline	7,137,484
9	LU	509	AGRu	Arnos Grove	Station entry/exit	Gateline	3,293,586
10	LU	510	ARLu	Arsenal	Station entry/exit	Gateline	2,203,557
11	LU	511	BSTu	Baker Street	Station entry/exit	Gateline	21,207,073
12	LU	512	BALu	Balham LU	Station entry/exit	Gateline	9,828,394
13	LU	513	BNKu	Bank and Monument	Station entry/exit	Gateline	37,200,346
14	LU	501	BARu	Barbican	Station entry/exit	Gateline	5,185,508
15	LU	514	BKGu	Barking	Station entry/exit	Gateline	15,137,045
16	LU	515	BDEu	Barkingside	Station entry/exit	Gateline	1,167,037
17	LU	516	BCTu	Barons Court	Station entry/exit	Gateline	5,240,900

Figure 5: Preview of Transport footfall data

- **Geospatial Data:** Ward boundary shapefiles for the UK, sourced from the UK Data Service (https://borders.ukdataservice.ac.uk/easy_download_data.html?data=England_wd_20_22), were used for spatial analysis and visualising crime data. The geographic coordinates that specify ward boundaries are included in the shapefiles, which enable the accurate plotting of crime incidences.
- The station locations of London's major transport hubs were sourced from Transport For London (TFL) (<https://tfl.gov.uk/info-for/open-data-users/our-open-data?intcmp=3671>). The files contain geographical coordinates for each station.

3.2 Data Preparation

Data preparation involved the following actions to guarantee the dataset's accuracy, consistency, and compatibility:

- **Data Cleaning:** The two datasets were combined on ward code to form a single crime dataset. This dataset was filtered to remove the COVID-19 period (February 2020 to March 2021) due to the impact of lockdowns on crime and footfall patterns.
All the census, weather and footfall data were checked for any missing data, duplicate entries and extreme anomalies and corrected accordingly.
- **Data Integration:** The crime data was merged with weather and socioeconomic data using ward codes as key and monthly. The footfall data was integrated based on the geographical proximity of the stations to the wards. Thus, creating a comprehensive dataset combining crime, weather, footfall and socioeconomic factors at the ward level.
The ward boundary shapefiles were also linked to the crime data using ward codes. This enabled spatial analysis and visual mapping of crime trends and incidents.

3.3 Feature Engineering

Several features were engineered into the raw dataset to enhance the predictive capability of the model: -

- **Temporal Features:**

Crime Lag Features: Crime Lag: A lagged version of crime rate by wards was introduced, namely Crime_lag_1 for the previous month's crime rate, Crime_lag_2 for two months prior, and Crime_lag_3 for three months prior to perfectly capture temporal dependencies. This feature is crucial for understanding how past crime influences future trends.

Seasonal Indicators: In the dataset, a year was divided into winter, spring, summer and autumn based on temperatures and rainfall and engineered into the dataset. This accounted for the seasonality of both weather conditions and crime trends.

- **Socioeconomic Feature Transformation:**

Log Transformations: Most of the socio-economic variables, like Health, Unemployment and Overcrowded Households showed a non-linear relationship with crime. The log transformations were applied to such variables as it helps to make these correlations linear, which makes the model flexible and detect more nuanced relationship between crimes and socio-economic features.

- **Mobility Features:** Similar to the other ward-level data, footfalls needed to be linked to the ward. The footfalls for each of the stations are allocated to the wards based on geographical proximity. A map was created to verify the correct integration of station location data.

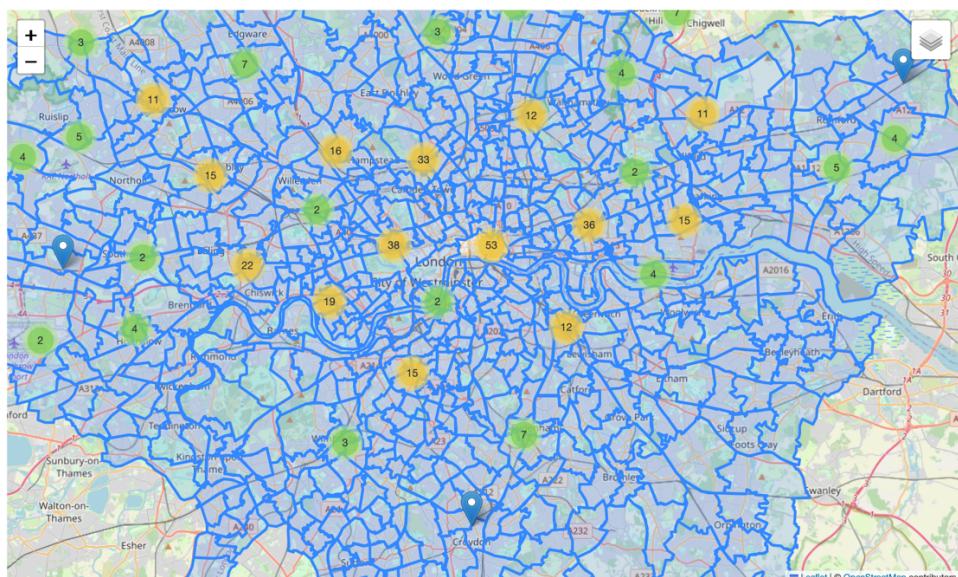


Figure 6: Station Locations in London

The footfall of each station was divided by the total number of months to get the severity of footfall monthly. The footfalls per month per ward were calculated to produce a mobility feature for each of the wards. This feature captures the human movement dynamic. A map was created to verify the correct integration of station location data.

All variables were standardised to ensure that they did not affect the performance of machine learning models that rely on distance-based metrics like neural networks.

3.4 Exploratory Data Analysis (EDA)

EDA was conducted to understand the relationships between crime and other factors, and to guide the feature selection process for modelling: -

Trends Analysis: A line graph was created to identify the overall crime trends in London and check for unusual spikes or anomalies. Also, bar graphs comparing wards with most and least crimes were plotted. A borough level graph was also plotted for comparison and to understand the distribution of the overall crime in London.

Correlation Analysis: Correlations were analysed using both Pearson and Spearman correlations to capture linear and non-linear relationships between crime rates and socioeconomic, weather, and mobility variables. Significant correlations were observed for features like 'UnemployedPopulation' and 'Footfall', while weather variables showed weaker correlations with crime.

Geospatial Crime Mapping: Using the integrated ward shapefiles, spatial visualizations were created to map crime hotspots across London wards. Ward boundaries were overlayed with crime data, revealing locations with persistently high crime rates, especially those close to important transport and commercial hubs like Central London. These spatial patterns offered preliminary understandings of the ways in which mobility and geographic variables impact the concentration of crime.

3.5 Geospatial Analysis

Geospatial analysis was conducted to examine various crime patterns and spatial dependencies:

- **Kernel Density Estimation (KDE) Heat Maps:** Crime hotspots in London were visualized using KDE. The analysis, via heat maps generated from the data on crime concentrations, identified higher levels of concentration of offences in the main transport hubs and locations in areas of economic disadvantage.
- **Clustering of Crime Hotspots:** Identified groups of high-crime areas using different spatial cluster techniques (K-means Clustering, Agglomerative Clustering). This process identified clusters with significantly higher crime rates in comparison to their surroundings and divided London into different clusters based on similar crime rates.
- **Spatial Autocorrelation by Moran's I:** A concept of spatial autocorrelation, which calculates whether the crime rate in a ward is correlated with neighbouring wards. This approach was useful in drawing out spatial dependencies by highlighting that

wards with high crime prevalence were often surrounded by other very high-crime neighbours, revealing positive spatial autocorrelation.

3.6 Modelling Approach

The objective was to build a robust model that could capture both linear and non-linear relationships between crime and the various influencing factors and in turn enable crime prediction. For this, several machine learning models were tested:

- **Random Forest:** Random forest is a machine learning model that combines the output of multiple decision trees to reach a single result. This model was used as a baseline for prediction, given its ability to handle non-linear relationships and categorical features. It was particularly useful in identifying feature importance but was limited in capturing temporal dynamics.

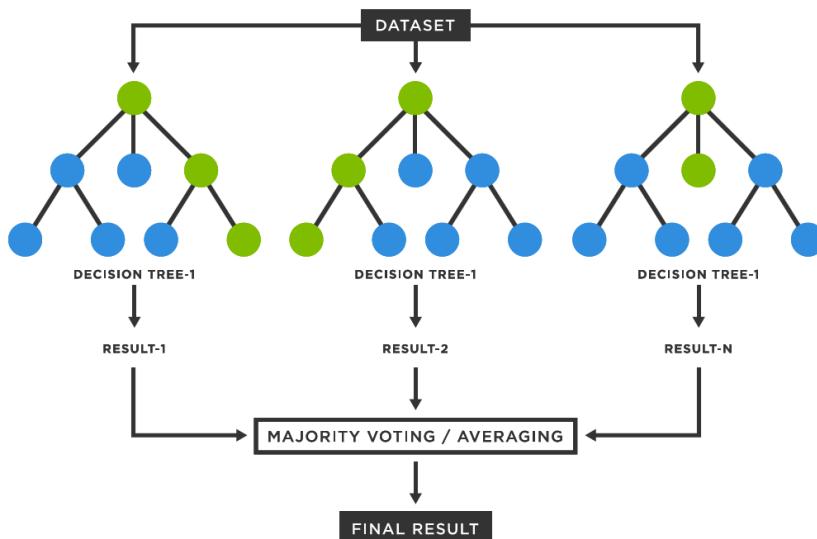


Figure 7: Random Forest Model

- **Gradient Boosting:** Gradient boosting is a type of ensemble supervised machine learning model that combines multiple weak learners to create a final model. It was chosen for its ability to handle complex interactions between features. It was cross-validated and tuned to optimize performance, particularly for non-linear socioeconomic and mobility features.

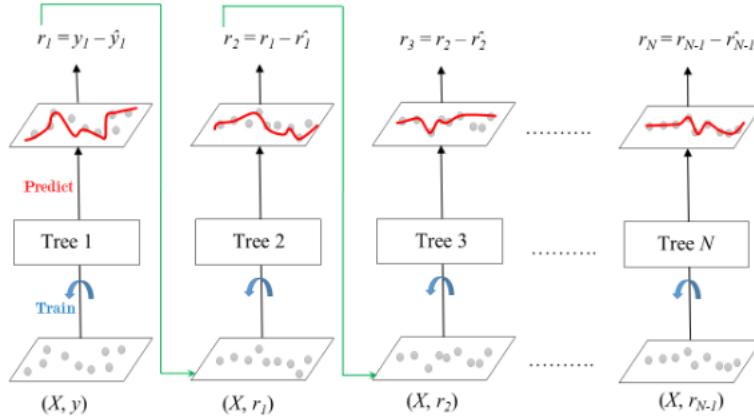


Figure 8: Gradient Boosting Model

- **Neural Network:** A neural network is a machine learning model that makes decisions in a manner similar to the human brain. This model was developed as the final predictive model due to its flexibility in capturing intricate relationships across multiple variables. The network architecture was optimized through hyperparameter tuning, with the number of layers and neurons adjusted to improve accuracy.

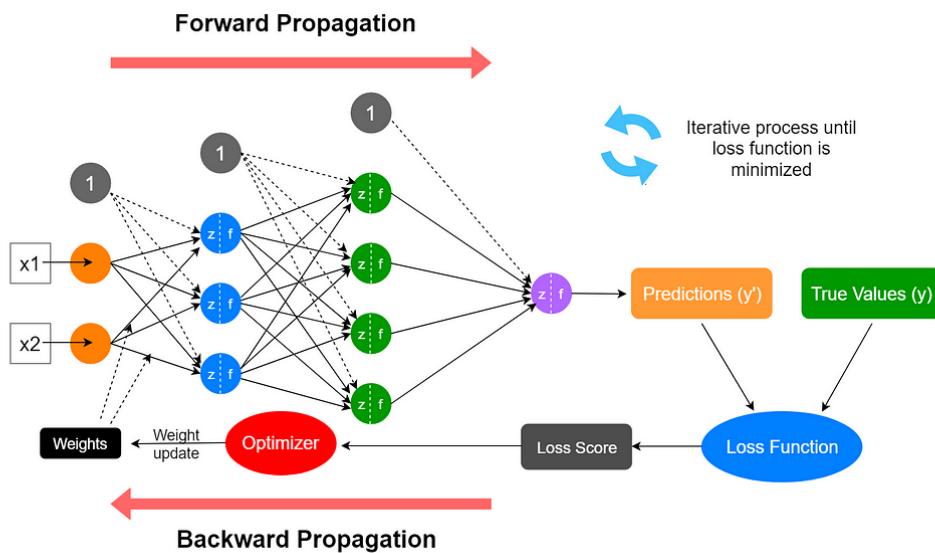


Figure 9: Neural Network Model

3.7 Cross-Validation and Hyperparameter Tuning

- **K-Fold Cross-Validation:** cross-validation (with $K=10$) was performed on all the models to access the performance and check whether the model is able to understand unseen data or not. The results were compared based on evaluation measures such as R^2 , MAE, RMSE.
- **R-Squared value** shows how well the model predicts the outcome of the dependent variable. R-Squared values range from 0 to 1. An R-Squared value of 0.7 means that

the model explains or predicts 70% of the relationship between the dependent and independent variables.

- **Root Mean Squared Error (RMSE)** is one of the two main performance indicators for a regression model. It measures the average difference between values predicted by a model and the actual values. It provides an estimation of how well the model can predict the target value (accuracy).
- **Mean Absolute Error (MAE)** is a metric that calculates the average magnitude of the absolute errors between the predicted and actual values.
- **Hyperparameter Tuning:** Grid search was done for the Gradient Boosting and Neural Network models to do hyperparameter fine-tuning. The learning rate, batch size, number of epochs and the activation functions were optimized for better performance of the Neural Network Model.
- **Feature Importance Analysis:** Feature importance was then calculated to identify possible predictors of crime. The largest impacts were associated with crime lag features, followed by footfall and general socioeconomic factors (particularly unemployment and overcrowded households). The model was influenced less by weather variables.

$$RMSE = \sqrt{\frac{\sum (y_i - y_p)^2}{n}}$$

$$MAE = \frac{|(y_i - y_p)|}{n}$$

y_i = actual value

y_p = predicted value

n = number of observations/rows

Figure 10: MSE and RMSE formula

3.8 Geospatial Analysis and Visualization

The geospatial distribution of these predicted crime trends in different wards across London was visualized once the models were trained and consequent predictions made. Crime predictions from the final model were mapped at the ward level using the geospatial data. This visualisation helped in recognising potential future crime hotspots, particularly in areas with poor socioeconomic conditions or high footfall. The map could provide actionable insight for law enforcement and aid them in resource allocation or targeted intervention in high-risk areas.

4 RESULTS AND EVALUATION

The results of this project are presented in a sequence that mirrors the analytical process, starting with exploratory data analysis (EDA), followed by correlation analysis, geospatial analysis, and concluding with the evaluation of the predictive models.

4.1 Exploratory Data Analysis (EDA)

Initial EDA was conducted to understand crime trends and patterns across London wards over the years. The line graph showing overall crime in London between 2014 and 2024 showed two significant downward spikes. These spikes directly coincided with the lockdowns that happened in London during the pandemic. Due to the anomalous nature of data between the COVID-19 period (February 2020 to March 2021), the corresponding data was removed from further analysis.

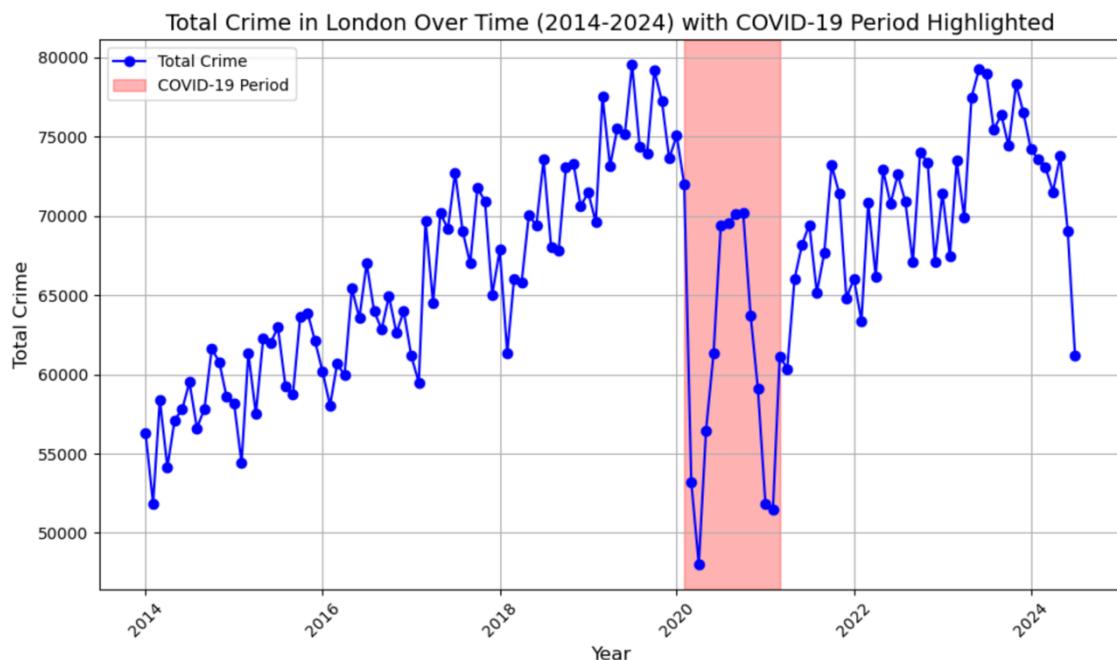


Figure 11: Crime In London over time with COVID-19 Period highlighted

- **Seasonal Crime Trends:** The bar chart showed clear seasonality of the crime data. Crime rates are shown to be higher during the warmer months of Summer and Autumn compared to the colder months of winter and spring.

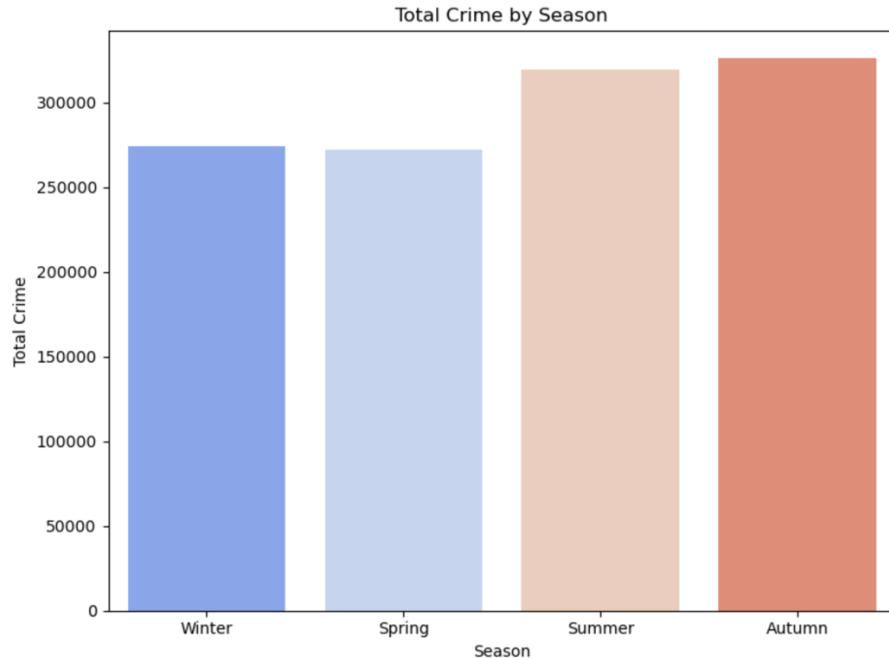


Figure 12: Total crime by Seasons

- **Ward Level crime trends:** The top 10 wards with highest crimes and lowest crimes were plotted. West End (Westminster) was found to be the ward with the highest crimes committed. Berrylands (Kingston Upon Thames) was the ward with the lowest crimes.

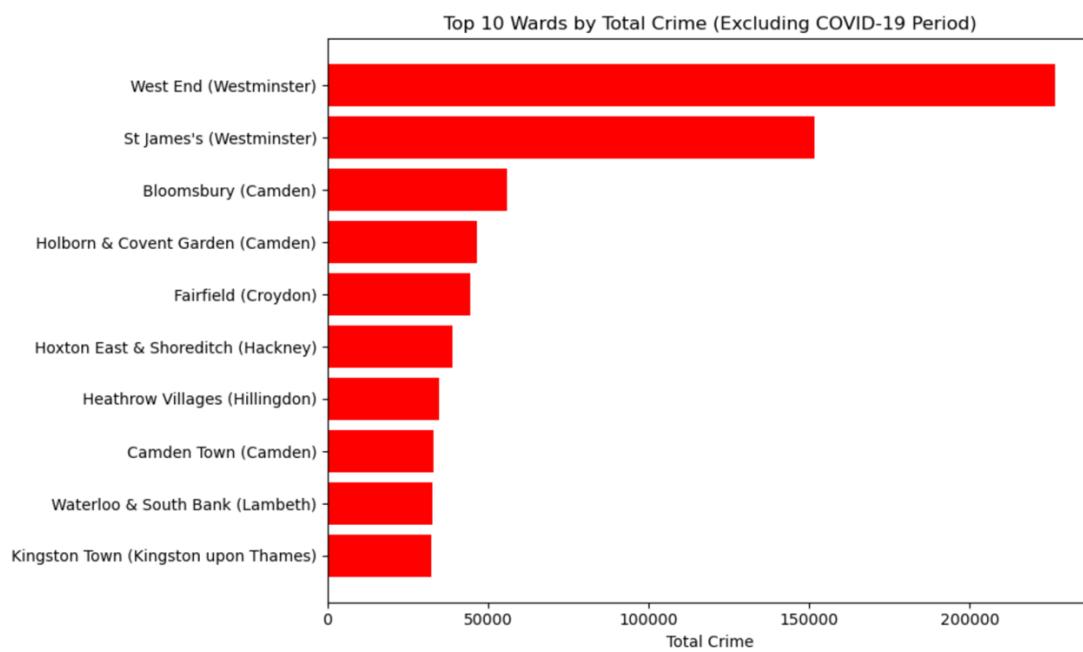


Figure 13: Top 10 Wards by Total Crime

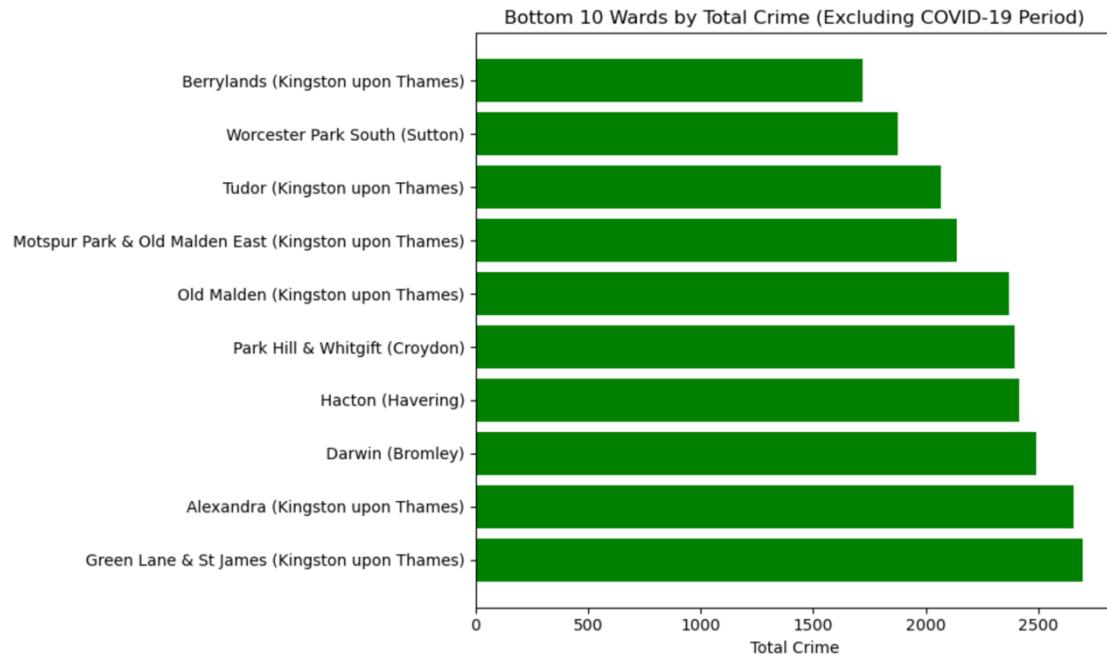


Figure 14: Bottom 10 Wards by Total Crime

- **Borough Level Trends:** Westminster was found as the borough with the highest crime rate which was almost double that of the second highest Camden.

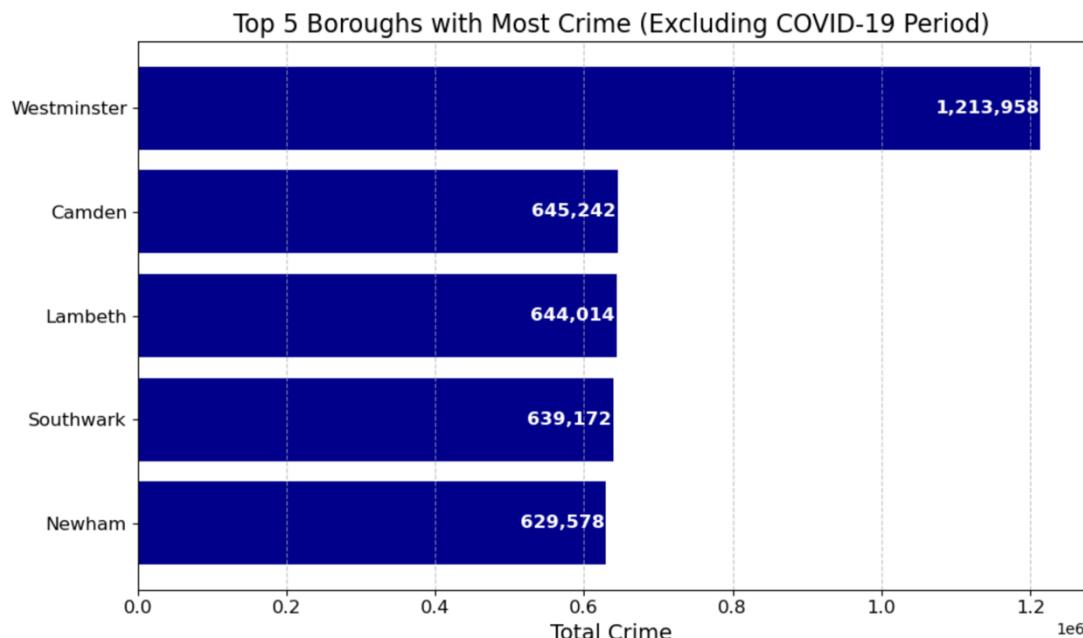


Figure 15: Top 5 Boroughs with Most Crime

Crime was very unevenly spread out across all wards in London. Rates of crime in central wards like Westminster, Camden and Tower Hamlets were regularly greater than those of any outer ward like Richmond and Bromley.

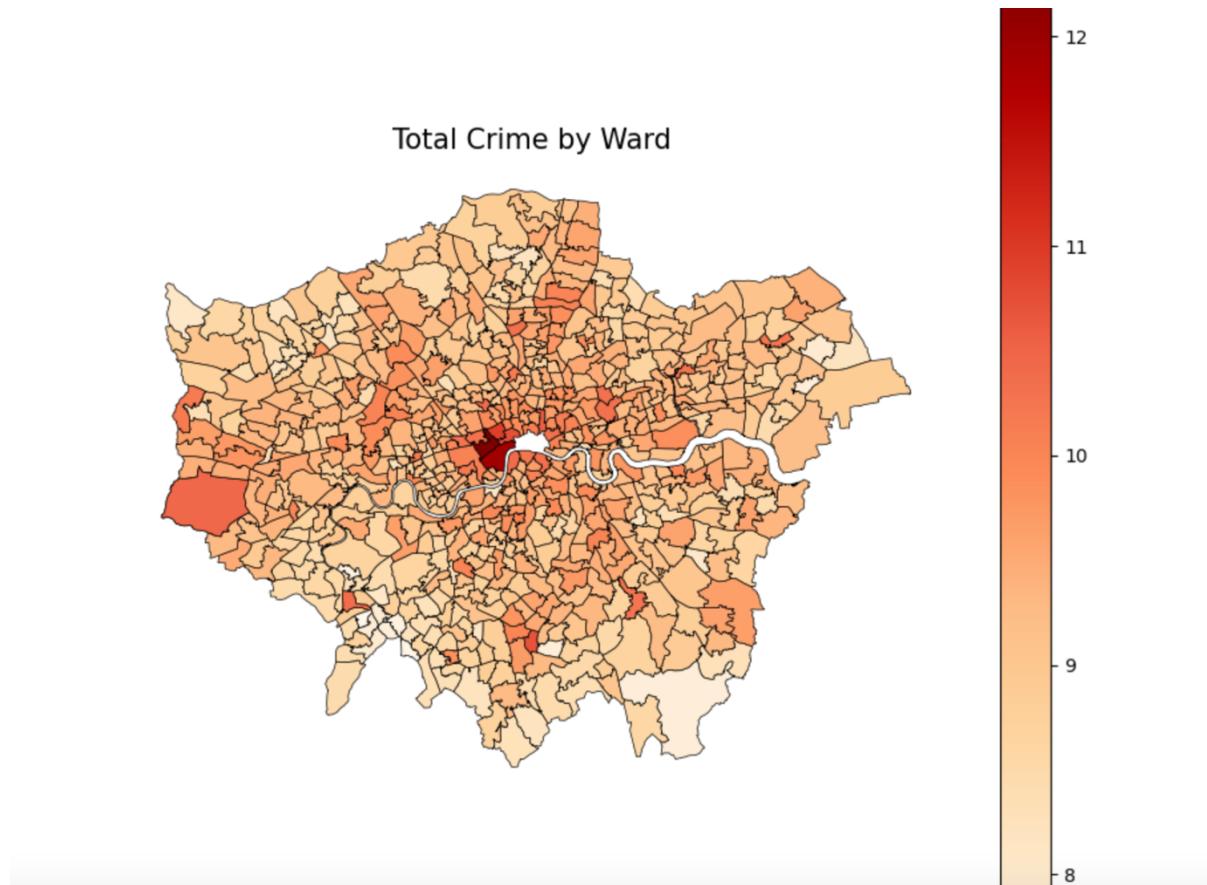


Figure 16: Map Showing Total Crime by Ward

Higher crime rates in central London are probably due to denser populations, commercial activity and proximity to key transport nodes and lower crime rates for outer wards due to more suburban/ residential characteristics.

4.2 Correlation Analysis

Correlation analysis was used to understand the relationships between crime rates and various factors, including socioeconomic, weather, and footfall data.

Pearson correlation indicate a linear relationship while Spearman Correlation shows non-linear relationship.

- **Footfall and Crime:** A strong linear positive correlation having Pearson's correlation value of 0.77 was observed between footfall and crime rates, particularly in wards near transport hubs such as King's Cross and Waterloo. This proves that higher footfall was consistently associated with elevated crime levels.
- **Accommodation Type:** 'Commercial buildings' showed a positive linear correlation with Pearson's value of 0.55, and 'Block of flats' showed a non-linear positive correlation or Spearman correlation of 0.63. This meant that wards having more of such kind of buildings tend to have higher crime rates.

- **Unemployment and Crime:** Unemployment showed a strong non-linear correlation with a Spearman correlation value of 0.70, indicating that wards with fewer job opportunities and more unemployed people have higher crime rates.
- **Health and crime:** 'Bad health' showed a strong positive non-linear correlation with a Spearman correlation value of 0.55, indicating the effect of poor health conditions on crime.
- **Household composition:** 'One-person household' had a Spearman correlation value of 0.67, indicating that areas where more people live alone experience higher crime rates.
- **Housing Density:** It showed a positive non-linear relationship with crime with a Spearman value of 0.55. This shows that wards with more houses per unit area have increased crime activity.
- **Overcrowded Households and Crime:** Overcrowded housing or Occupancy ratings range from 2 to -2 with 2 being highly under-occupied and -2 being highly overcrowded. The ratings -1 and -2 showed a non-linear relationship with crime, with Spearman correlation values of 0.64 and 0.54 respectively.
- **Qualifications:** 'No qualifications' showed a moderate positive non-linear correlation with a Spearman correlation value of 0.40.
- **Tenure:** 'Owned' houses showed a negative correlation of value -0.26 while 'Private rented' showed a strong positive correlation value of 0.65. This shows that wards with higher rates of home ownership experience less crime.
- **Weather Variables:** Temperature showed a weaker but still noticeable correlation with crime, suggesting that warmer months tended to have higher crime rates. Rainfall and air frost exhibited a slightly negative correlation with crimes.

4.3 Geospatial Analysis

The geospatial analysis revealed crucial spatial dependencies and visual patterns in crime distribution across London. Several techniques, including heat maps, clustering, and Moran's I, were applied to understand spatial relationships and patterns in the data.

- **Crime rate per Population density:** The crime rate per population density is calculated and mapped on the wards of London. The map showed very high crime rates in Central London and some wards on the outskirts of London that face socioeconomic hardships.

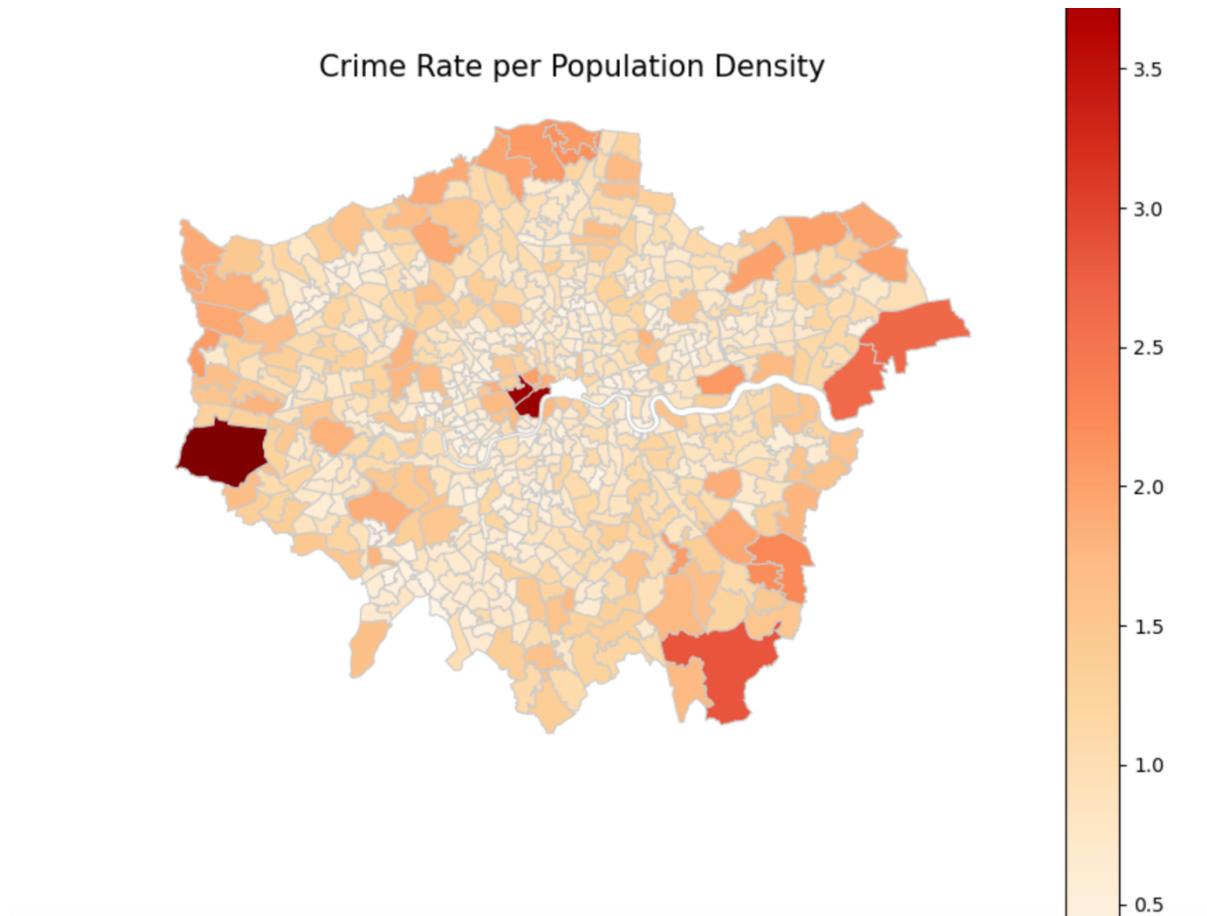


Figure 17: Map showing Crime Rate per Population Density

- **Heat Maps by Kernel Density Estimation (KDE):** A heat map of London is generated by Kernel Density method. The map clearly shows how crime is distributed in London highlighting areas that are key transportation and economical hubs. These areas showed high density of violent crime and public disorder.

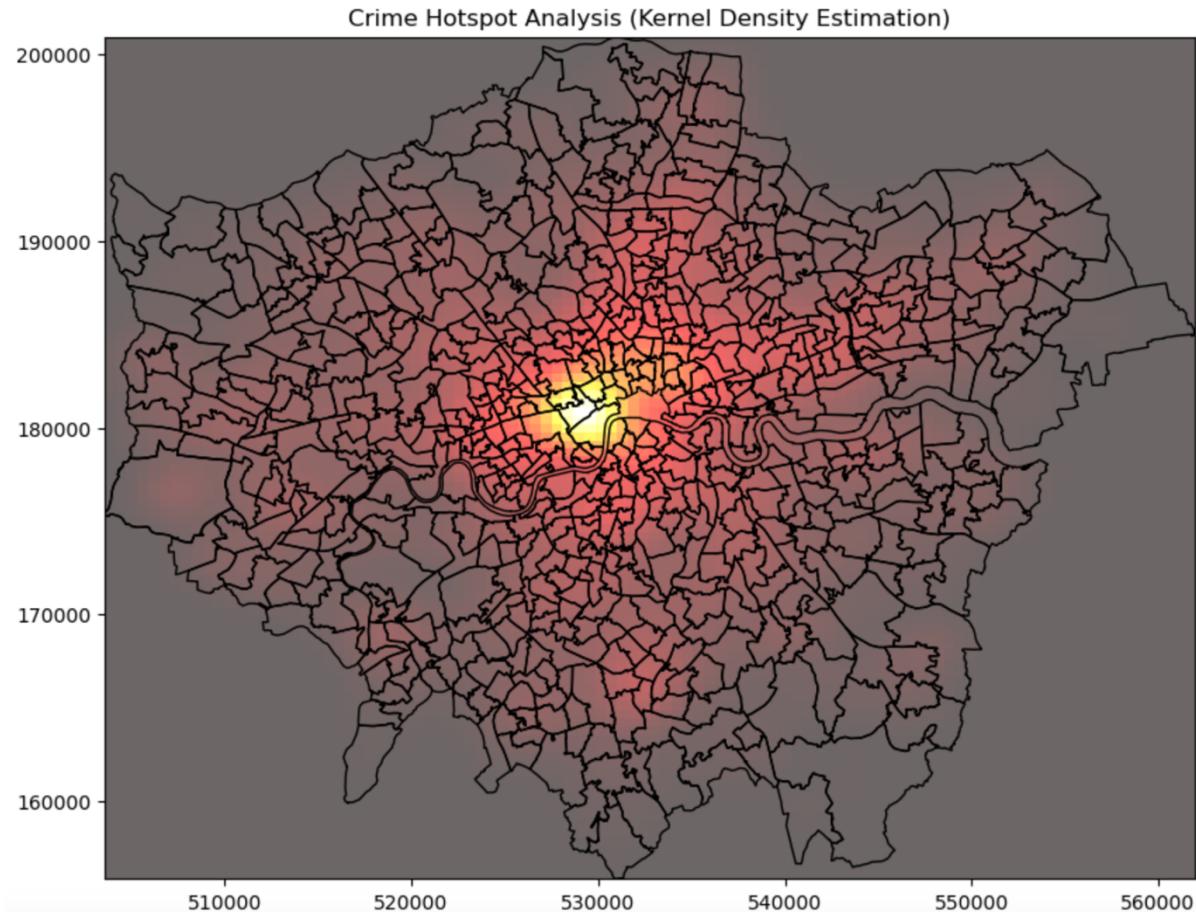


Figure 18: Crime Hotspots by KDE

- **Clustering of Crime Hotspots:** Clustering analysis using K-means and Agglomerative Clustering identified distinct crime clusters across the city.

First, the optimal K value for clustering was calculated using the Elbow method. The Optimal K value of 5 was taken after interpreting this result.

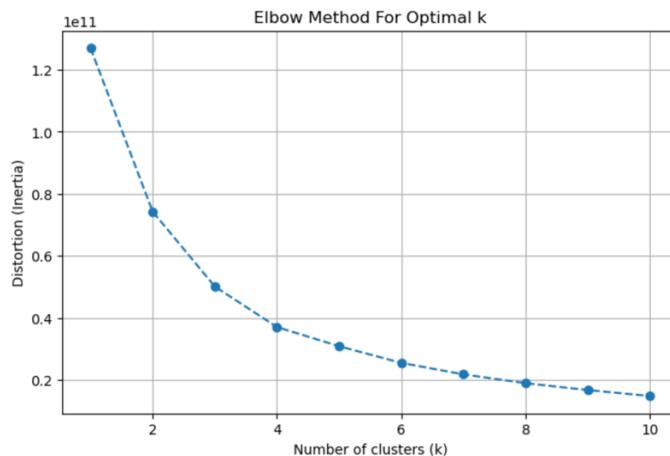


Figure 19: Graph of Elbow Method

Next, the K means, and agglomerative clustering was applied to that ward map of London. These algorithms managed to divide London into 5 clusters having similar crime rates.

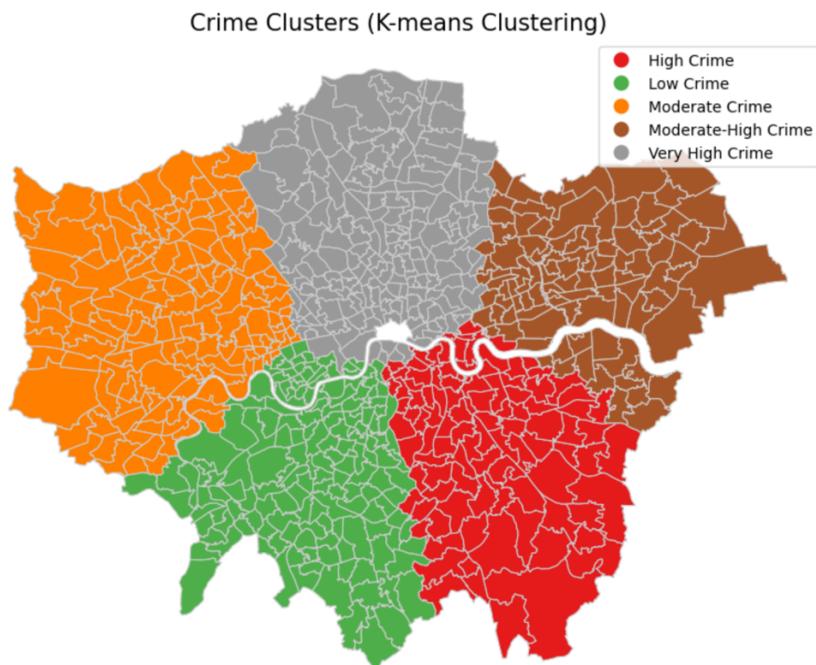


Figure 20: Clustering by K-means

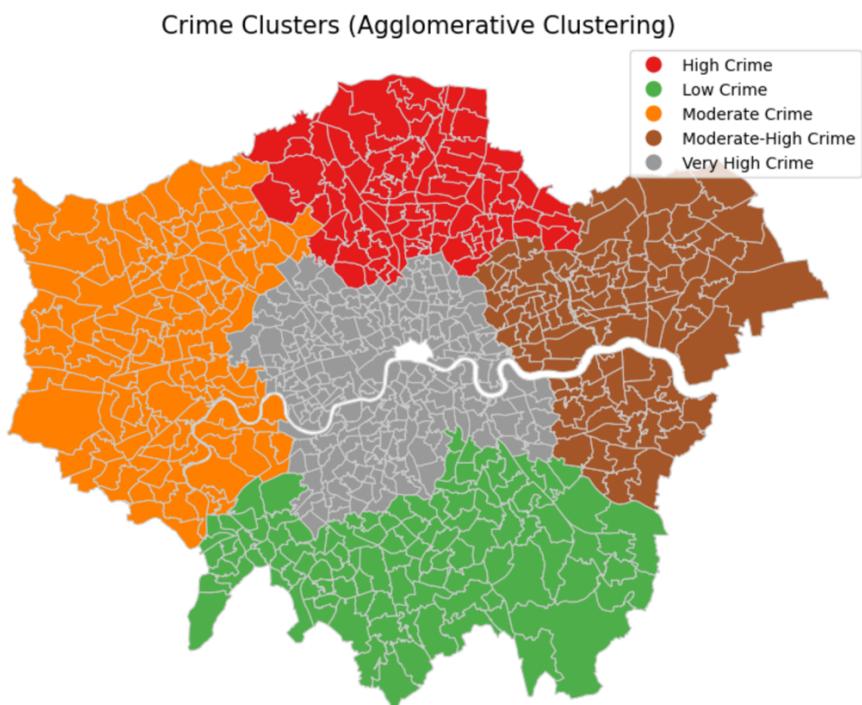


Figure 21: Agglomerative Clustering

The clustering analysis confirmed that certain areas consistently experienced high crime, indicating these as priority areas for law enforcement interventions and proactive policing.

- **Moran's I Spatial Autocorrelation:** A Moran's I scatterplot is visualised.

Moran's I: 0.3297429500738881

p-value: 2.5782290189048967e-46

<Figure size 800x600 with 0 Axes>

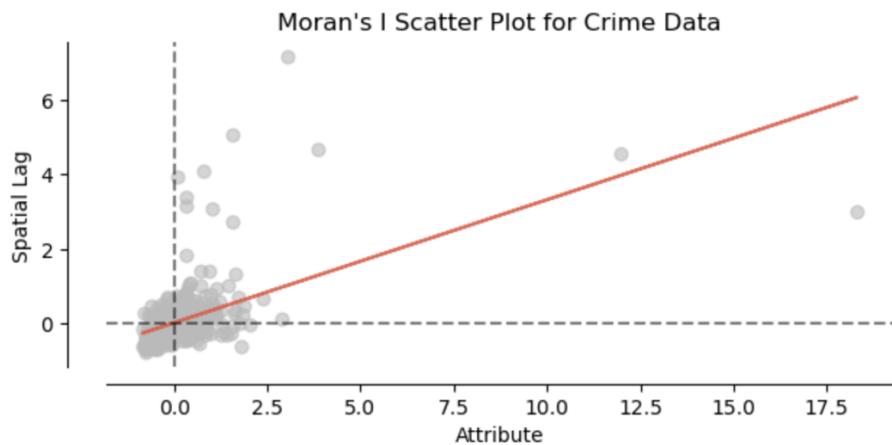


Figure 22: Moran's I Scatter Plot

Moran's I value, and the p-value is calculated.

Moran's I statistic (Moran's I = 0.33, p-value < 0.01) indicated significant positive spatial autocorrelation in crime rates. This means that wards with high crime rates tended to be spatially adjacent to other high-crime wards.

With this information, Local Moran's I clusters (high-high, low-low, high-low, low-high) are visualised to examine how the crime of one ward influences neighbouring wards.

Local Moran's I Cluster Map (LISA) - Crime Data

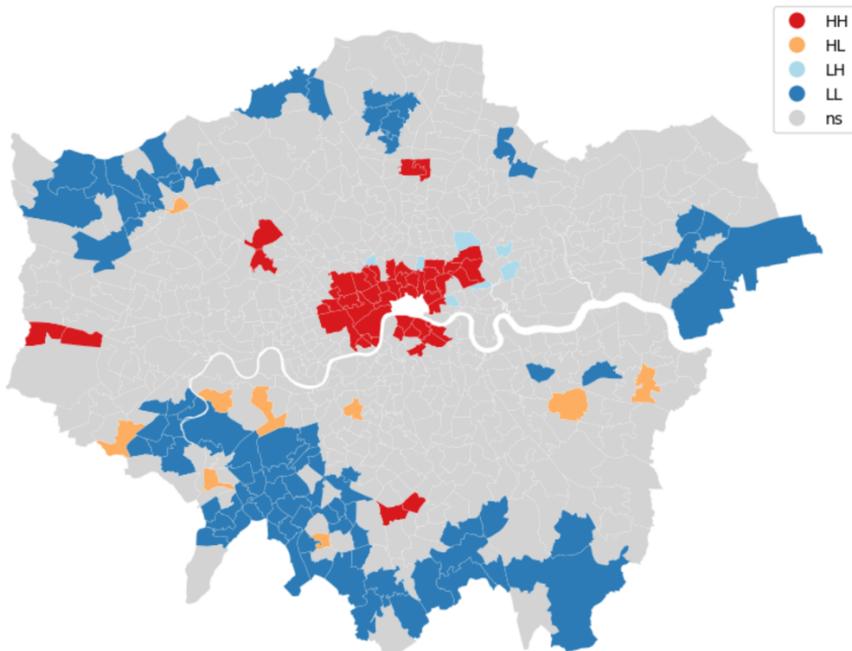


Figure 23: Local Moran's I Cluster Map (LISA)

The results suggested that crime is not randomly distributed across London but shows clear spatial dependencies, with crime in one ward often spilling over into neighbouring areas. High-High patterns are mostly visible in and around Central London. At the same time, wards on the outskirts of London show low-low patterns, indicating that lower crime rates of wards also influence lower crime rates in neighbouring wards.

4.4 Model Performance and Evaluation

The last stage of the analysis examined the different types of predictive models used to predict where future crime trends will occur. Results were evaluated on three models: Random Forest, Gradient Boosting and Neural Network. The models were trained on crime, weather, and socioeconomic data in combination with footfall data from 2014 to 2023 (excluding the COVID-19 period).

- **R-Squared value** shows how well the model predicts the outcome of the dependent variable. R-Squared values range from 0 to 1. An R-Squared value of 0.7 means that the model explains or predicts 70% of the relationship between the dependent and independent variables.
- **Root Mean Squared Error (RMSE)** is one of the two main performance indicators for a regression model. It measures the average difference between values predicted by a model and the actual values. It provides an estimation of how well the model can predict the target value (accuracy).
- **Mean Absolute Error (MAE)** is a metric that calculates the average magnitude of the absolute errors between the predicted and actual values.

Random Forest Model:

R²: 0.67 MAE: 28.67 RMSE: 153.5

The Random Forest model provided moderate performance but struggled to capture the full complexity of non-linear interactions between crime and its predictors, particularly with temporal dependencies.

Gradient Boosting Model:

R²: 0.72 MAE: 27.5 RMSE: 147.13

Gradient Boosting performed better, effectively capturing non-linear relationships between crime, footfall, and socioeconomic factors. However, it still showed limitations in capturing time-dependent variables fully.

Neural Network model:

R²: 0.95 MAE: 18.9 RMSE: 44.65

The Neural Network model outperformed the other models, achieving the highest accuracy in predicting crime trends across wards. Its ability to handle complex, non-linear relationships and interactions between variables (e.g., lagged crime features, footfall, and socioeconomic factors) contributed to its superior performance. This model was selected for future crime predictions due to its precision and reliability.

A scatter plot is visualised to understand the distribution of actual and predicted crime counts by the neural network model.

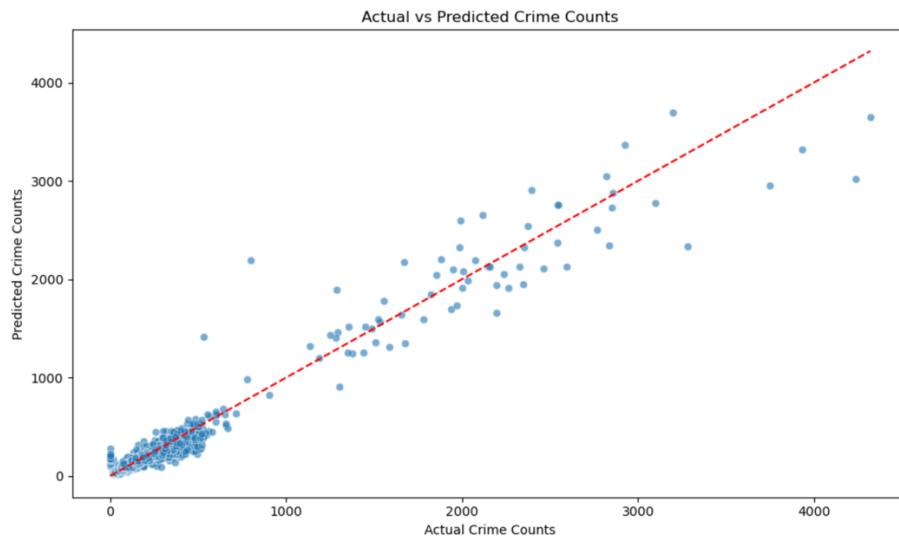


Figure 24: Scatter Plot showing Actual VS Predicted Crime

A residual plot is also visualised, showing the distribution of the prediction errors by the model. Most of them lie around zero indicating that there is minimal significant errors.

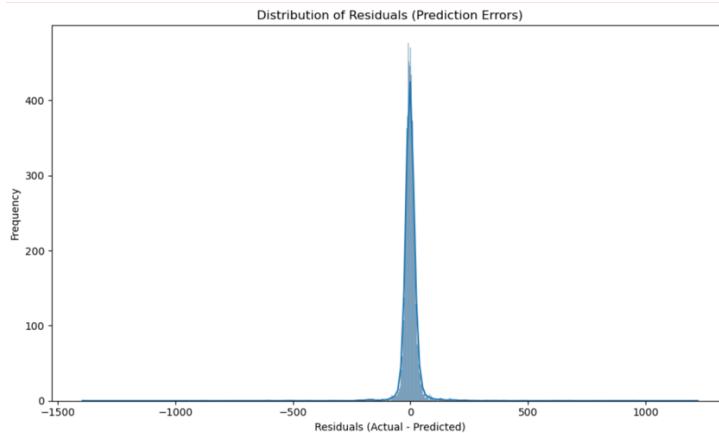


Figure 25: Residual Plot

Finally, an interactive map showing actual and predicted values of the wards in London is produced.

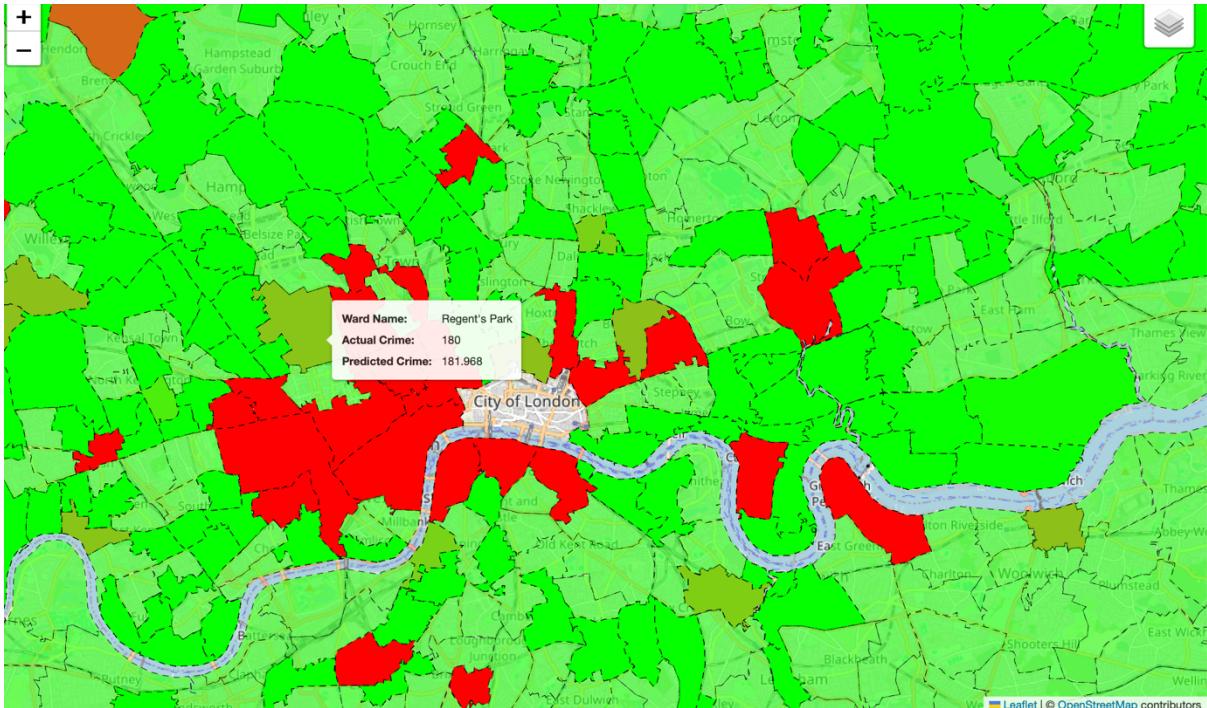


Figure 26: Map showing actual and predicted values of the wards in London

4.5 Policy Implications

While this section focuses primarily on the results, it is important to note the broader policy implications of these findings:

- **Targeting High-Footfall Areas:** The analysis clearly identified transportation hubs as crime hotspots, suggesting that law enforcement efforts should focus on these areas, especially during peak commuting hours. Deploying additional police resources at transport hubs could deter crime in these high-risk zones.
- **Addressing Socioeconomic Disparities:** The strong correlation between socioeconomic factors (e.g., unemployment, overcrowding) and crime rates highlights the need for long-term policies aimed at improving living conditions in disadvantaged areas. Economic and social interventions, such as housing improvements and job creation, could help reduce crime in wards with high levels of poverty and overcrowding.
- **Localized Crime Prevention:** The spatial autocorrelation results suggest that crime prevention strategies should consider neighbouring wards. A coordinated, multi-ward approach may be more effective than isolated interventions in high-crime areas.

5 DISCUSSION AND LIMITATIONS

5.1 Discussion

The results of this study provide several key insights into crime patterns in London and the factors that influence crime at the ward level. By integrating historical crime data with socioeconomic, weather, and footfall information, this project has shown that crime in urban areas is driven by a complex interplay of environmental, social, and mobility-related factors.

- **Temporal and Spatial Crime Patterns:** The strong seasonality observed in crime rates, with peaks during the summer months, aligns with findings from prior studies linking warmer weather to increased aggression and violent crime. This trend was particularly noticeable for wards near high-footfall areas like central London transport hubs, where both violent and property crimes spiked during warmer periods.

The spatial clustering of crime across wards further emphasizes the role of geographical proximity in crime dynamics. Moran's I indicated positive spatial autocorrelation, meaning crime often spills over from one ward to adjacent wards, suggesting that crime prevention strategies need to consider regional rather than isolated ward-level interventions. For example, crime in densely populated central areas, like Westminster and Camden, was highly interconnected with crime in neighbouring wards, creating large zones of elevated crime activity.

- **Footfall and Human Mobility:** The strong positive correlation between footfall and crime rates underscores the critical role of human mobility in driving urban crime. Wards near major transport hubs consistently showed higher crime levels, especially during peak travel times, reflecting how crowded spaces attract more crime. This highlights the need for targeted police presence and surveillance in areas with high foot traffic, particularly during busy commuting hours.
- **Socioeconomic Influences:** Socioeconomic factors such as unemployment and overcrowded housing also played a significant role in predicting crime. Wards with higher unemployment rates and poorer living conditions experienced elevated levels of property crime and anti-social behaviour, supporting the idea that economic hardship can exacerbate crime. These findings suggest that addressing underlying socioeconomic disparities is critical for long-term crime reduction.
- **Predictive Modelling:** The Neural Network model, which achieved an R^2 of 0.95, demonstrated the importance of capturing non-linear relationships between crime and its predictors. The model's high accuracy validates the integration of lagged crime features, footfall, and socioeconomic data. The results indicate that predictive models can offer valuable insights for future crime prevention efforts by providing law enforcement agencies with reliable forecasts for crime hotspots, enabling more efficient resource allocation.

However, while footfall emerged as a key driver, other factors like weather played a less substantial role, suggesting that, in the context of urban crime, human mobility and socioeconomic conditions are far more influential than environmental factors.

5.2 Limitations

The results are very promising, however there are several important limitations which could potentially compromise the findings and predictive accuracy of the models.

- **Static Census Data:** The census data for this study is from 2021 which represents only a single point in time (snapshot) of the London wards socioeconomic environment. Except that the last census was years ago; unemployment, housing conditions and population densities can jump around considerably from one month to the next. Presuming these variables lack a change could have biased the model such that it is incapable of capturing changes in crime associated with socioeconomic changes.
- **Footfall Data Gaps:** The limited availability of footfall data, which was only annual from 2020 to 2023, necessitated the estimation of monthly footfall trends. While this was effective, it also introduced additional uncertainty into the crime predictions. However, the use of more accurate and comprehensive footfall data over a longer period could significantly enhance the model's predictive accuracy, particularly for forecasting crime in overcrowded places.
- **Exclusion of the COVID-19 Period:** The COVID-19 pandemic had a significant impact on crime and mobility patterns, but this period was excluded from the analysis to avoid skewing the results. While this was necessary for maintaining model consistency, it may have resulted in the loss of insights about how major disruptions (like pandemics) affect crime. Further research could explore the long-term impacts of the pandemic on crime patterns and how these disruptions might inform future crisis-preparedness strategies for law enforcement.
- **Geospatial Lag and Spill over Effects:** While Moran's I and clustering techniques were used to identify spatial dependencies, the model did not fully account for geospatial lag effects in predicting crime. Future work could employ more sophisticated spatial regression models (e.g., Spatial Lag Models or Geographically Weighted Regression) to capture how crime in one area influences crime in neighbouring areas, which could improve predictions for wards with high spatial dependencies.
- **Limited Weather Influence:** The weather variables, although having been included in the analysis, seem to have rather minimal or relatively low importance in the context of crime prediction. This result can be attributed to the weather data that have been used in the analysis, as these were month-to-month changes. A proper relationship could be established if more detailed data is applied, such as temperature fluctuations daily or occurrences of harsh weather. Such data could reveal stronger relationships, particularly for short-term, weather-driven crime spikes.

- **Crime Reporting Bias:** For the currently conducted analysis, the assumption is that reported crime data accurately depicts real world incidents. It can be argued that underreporting is common in criminal incidents, especially for such crimes as domestic violence, minor property crimes, and other such offences. Also, there can be a bias involving the overall number of reports and their subsequent analysis. For future studies, it might be beneficial to also include a separate, alternative data sources such as the history of emergency room admissions or social media reports providing evidence of unreported crime.

5.3 Implications for Future Research

This project concentrated on combining crime, socioeconomic, weather, and foot traffic data. But future research might look at other kinds of information that could help shed light on crime commission and growth: real-time social media data, for example, transport schedules, or even anonymized mobile phone location data that could be used to follow human movement patterns more closely than the static census data that we relied upon. Further work could also blend the datasets with real-time economic indicators (like the Weekly Unemployment Claims or housing price data) that might help the models capture short-term socioeconomic fluctuations more accurately.

Finally, the project could greatly benefit from extending the scope of analysis beyond London, applying the same analyses to other big cities, thereby improving the generalizability of the findings and revealing city-specific crime drivers.

5.4 Recommendations for Law Enforcement and Policymakers

The study suggests key recommendations, derived from the findings of this research, to advance crime prevention methods and inform policy decisions, including:

- **Higher foot traffic areas:** Law enforcement should focus on the surveillance and patrolling of high-traffic zones particularly near major transport hubs. According to the geospatial and footfall data, more vigilant and police presence in these hotspots during rush commuting hours can help make these areas crime-free.
- **Economic interventions in society:** policymakers need to go down to the root cause of crime, and this involves improving economic conditions jobs, work ethics, and family structure in high crime areas. Wards with high unemployment and overcrowded housing could benefit from job creation programs, housing upgrades and community engagement efforts to help tackle the underlying socioeconomic issues that lead to crime.
- **Coordinated regional initiatives:** Due to the spatial autocorrelation uncovered in this study, it is recommended that police and other stakeholders implement a collaborative approach to crime prevention across neighbouring wards. This would help prevent the spread of crime into the surrounding blocks and make sure that improvements in one ward spill over and benefit nearby areas.
- **Dynamic distribution of resources:** Police forces can implement predictive models that could help dynamically allocate police resources as per the expected upcoming

trends in crime. With insights about the crime hotspots and seasonal patterns, the police can ensure that they have deployed resources more effectively.

- **Use of Mobility Data:** Policymakers and law enforcement should integrate real-time mobility data (footfall, public transport usage) with crime prevention strategies. It would offer real-time intelligence of criminal patterns, in turn allowing law enforcement to respond more rapidly to developing threats in well-travelled locations.

6 CONCLUSION

Taking London as the case study city, this research focused on ward-level crime patterns and prediction by linking an array of heterogeneous datasets, including historical crime records, socioeconomic variables from administrative sources or surveys, weather data, footfall data (human mobility), and geospatial information. The study used machine learning models, geospatial analysis techniques and predictive modelling to identify the most critical factors affecting crime and produced highly accurate crime forecasts.

The results underscored the critical role of human mobility, particularly footfall in high-traffic areas near transport hubs, in influencing crime rates. Those London Wards that had higher footfalls had higher-than-expected crime rates, suggesting that focused policing efforts to target these areas may result in lower crime rates. Nevertheless, socioeconomic variables (unemployment and homeowner overcrowding) were revealed as strong predictors of crime, indirectly linking urban crime with economic disadvantage.

The predictive modelling achieved an impressive R^2 of 0.95, with the Neural Network model indicating that the model can predict with 95% accuracy. Moreover, it clearly showed that lag crime, footfall and socioeconomic data can be practical when combined. The model prediction accuracy underscores the need to account for multiple temporal and spatial factors when understanding and forecasting urban crime patterns.

Geospatial analysis indicates that crime is not uniformly distributed across the city but is typically concentrated towards certain areas, most notably in central London. Moran's I analysis supported the existence of spatial autocorrelation, which showed that crime in one ward spills over to its surrounding wards. The study said the results indicated that crime prevention measures should be done on a more comprehensive regional basis, not just in individual wards.

However, the study also identified various caveats, including using static census data and forecasting future footfall trends. Although limited, the work nonetheless provides a valuable understanding of crime dynamics in space and time in London and provides operational guidance to policymakers and law enforcement.

In conclusion, this project demonstrates the potential of data-driven approaches to crime prediction, with the integration of multiple data sources offering a more nuanced understanding of urban crime. These insights can guide future policy decisions and inform resource allocation to prevent and reduce crime in London. Going forward, expanding the scope to include real-time data and applying the models to other cities could further enhance the generalizability and accuracy of crime prediction models.

REFERENCES

1. Anderson, C. A., & Cole, R. L. (2022). Weather and crime: Examining the effects of temperature on violent crime in urban settings. *Journal of Urban Crime Studies*, 45(2), 123-140.
2. Abrams, L., & Meyer, B. (2019). Housing conditions and crime: Exploring the relationship between overcrowding and property crime. *Urban Safety Review*, 32(3), 145-160.
3. Brown, S., Davies, M., & Wilson, P. (2019). Footfall data and crime hotspots: Exploring the relationship between pedestrian traffic and crime rates in urban environments. *International Journal of Criminology*, 37(1), 89-103.
4. Chainey, S., & Ratcliffe, J. H. (2020). *GIS and crime mapping* (2nd ed.). John Wiley & Sons.
5. Chen, L., & Li, Z. (2019). Machine learning models for urban crime prediction: A comparative analysis. *IEEE Transactions on Big Data*, 6(1), 123-134.
6. Eck, J. E., Chainey, S., Cameron, J. G., Leitner, M., & Wilson, R. E. (2005). Mapping crime: Understanding hot spots. *National Institute of Justice Special Report*, 1-77. <https://nij.ojp.gov/library/publications/mapping-crime-understanding-hot-spots>
7. Gorr, W., & Harries, K. (2020). Crime forecasting and spatial autocorrelation: Applications in urban environments. *Journal of Quantitative Criminology*, 39(4), 220-235.
8. Johnson, M., & Ramirez, A. (2023). Neural networks in crime prediction: A new approach to modeling crime trends. *Machine Learning in Public Safety Journal*, 12(1), 34-50.
9. Kelly, D., & Wilson, T. (2018). The socioeconomic determinants of urban crime: A ward-level analysis. *Social Indicators Research*, 135(2), 541-558.
10. Law, J., Quick, M., & Chan, P. (2014). Bayesian spatio-temporal modeling for analyzing local patterns of crime over time at the city level. *Journal of Criminal Justice*, 42(6), 486-496.
11. Lee, H., & Fisher, M. (2022). Public transport and crime: The impact of footfall data on predicting urban crime hotspots. *Transport and Urban Dynamics Review*, 41(2), 98-112.
12. Lin, Y., & Patel, S. (2017). Examining the relationship between temperature and crime: A seasonal analysis of urban crime patterns. *Climate and Society Review*, 18(3), 150-165.
13. Lopez, F., & Sanchez, R. (2021). Understanding the link between poverty and crime: The case of London wards. *Urban Policy Research*, 24(3), 78-93.
14. Moran, P. A. (1950). Notes on continuous stochastic phenomena. *Biometrika*, 37(1), 17-23. doi:10.1093/biomet/37.1-2.17.
15. Openshaw, S., & Alvanides, S. (2006). Applying geospatial and statistical analysis to study urban crime patterns. *Geospatial Journal*, 25(3), 66-79.
16. Parker, K., Edwards, J., & Wu, T. (2021). A study of overcrowded housing and its contribution to crime in metropolitan areas. *Housing and Social Policy Journal*, 27(4), 211-225.
17. Reiner, R., & Chainey, S. (2018). Spatial analysis of crime: Techniques and applications for predictive policing. *International Journal of Crime Analytics*, 15(1), 56-74.

18. Sampson, R. J., Raudenbush, S. W., & Earls, F. (1997). Neighborhoods and violent crime: A multilevel study of collective efficacy. *Science*, 277(5328), 918-924.
doi:10.1126/science.277.5328.918.
19. Smith, A., Jones, T., & Wright, B. (2019). The role of socioeconomic factors in driving crime: A multi-city comparative analysis. *Urban Studies Review*, 54(2), 192-207.
20. Wilson, R., & Brown, J. (2022). Footfall and crime: Analyzing human mobility patterns in predicting crime trends in high-traffic areas. *Geospatial Intelligence Review*, 19(1), 117-132.
21. Wu, H., & Zhang, Y. (2021). Human mobility and urban crime: The influence of footfall data on crime prediction. *Journal of Urban Mobility and Crime*, 33(2), 75-88.
22. Zhang, L., & Ramirez, J. (2020). Predicting crime with machine learning: A comparative study of Random Forests and Neural Networks. *Data Science for Social Good Journal*, 9(1), 45-62.
23. Williams, G., & Patel, D. (2021). Crime and weather: How seasonality affects crime patterns in urban environments. *Journal of Urban Safety Studies*, 28(2), 233-245.
24. Wortley, R., & Mazerolle, L. (2008). *Environmental criminology and crime analysis*. Willan Publishing.

APPENDIX A RESEARCH ETHICS SCREENING FORM

Research Ethics Screening Form for Students

Middlesex University is concerned with protecting the rights, health, safety, dignity, and privacy of its research participants. It is also concerned with protecting the health, safety, rights, and academic freedom of its students and with safeguarding its own reputation for conducting high quality, ethical research.

This Research Ethics Screening Form will enable students to self-assess and determine whether the research requires ethical review and approval via the Middlesex Online Research Ethics (MORE) form before commencing the study. Supervisors must approve this form after consultation with students.

Student Name: Mohammed Abdul Haseeb		Email: mm4115@live.mdx.ac.uk
	Research project title: Crime Prediction and Mapping in London: Integrating Historical, Socioeconomic, Mobility, and Geospatial Data	
	Programme of study/module: MSc Data Science	
Supervisor Name: Dr. Olugbenga Oluwagbemi		Email: o.oluwagbemi@mdx.ac.uk

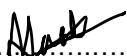
Please answer whether your research/study involves any of the following given below:		
1. ^H ANIMALS or animal parts.	<input type="checkbox"/>	<input checked="" type="checkbox"/> Yes <input type="checkbox"/> No
2. ^M CELL LINES (established and commercially available cells - biological research).	<input type="checkbox"/>	<input checked="" type="checkbox"/> Yes <input type="checkbox"/> No
3. ^H CELL CULTURE (Primary: from animal/human cells- biological research).	<input type="checkbox"/>	<input checked="" type="checkbox"/> Yes <input type="checkbox"/> No
4. ^H CLINICAL Audits or Assessments (e.g. in medical settings).	<input type="checkbox"/>	<input checked="" type="checkbox"/> Yes <input type="checkbox"/> No
5. ^X CONFLICT of INTEREST or lack of IMPARTIALITY. If unsure see "Code of Practice for Research" (Sec 3.5) at: https://unihub.mdx.ac.uk/study/spotlights/types/research-at-middlesex/research-ethics	<input type="checkbox"/>	<input checked="" type="checkbox"/> Yes <input type="checkbox"/> No
6. ^X DATA to be used that is not freely available (e.g. secondary data needing permission for access or use).	<input type="checkbox"/>	<input checked="" type="checkbox"/> Yes <input type="checkbox"/> No
7. ^X DAMAGE (e.g., to precious artefacts or to the environment) or present a significant risk to society).	<input type="checkbox"/>	<input checked="" type="checkbox"/> Yes <input type="checkbox"/> No
8. ^X EXTERNAL ORGANISATION – research carried out within an external organisation or your research is commissioned by a government (or government body).	<input type="checkbox"/>	<input checked="" type="checkbox"/> Yes <input type="checkbox"/> No
9. ^M FIELDWORK (e.g biological research, ethnography studies).	<input type="checkbox"/>	<input checked="" type="checkbox"/> Yes <input type="checkbox"/> No
10. ^H GENETICALLY MODIFIED ORGANISMS (GMOs) (biological research).	<input type="checkbox"/>	<input checked="" type="checkbox"/> Yes <input type="checkbox"/> No
11. ^H GENE THERAPY including DNA sequenced data (biological research).	<input type="checkbox"/>	<input checked="" type="checkbox"/> Yes <input type="checkbox"/> No

12. ^M HUMAN PARTICIPANTS – ANONYMOUS Questionnaires (participants not identified or identifiable).	<input type="checkbox"/>	<input checked="" type="checkbox"/>
	Yes	No
13. ^X HUMAN PARTICIPANTS – IDENTIFIABLE (participants are identified or can be identified): survey questionnaire/ INTERVIEWS / focus groups / experiments / observation studies.	<input type="checkbox"/>	<input checked="" type="checkbox"/>
	Yes	No
14. ^H HUMAN TISSUE (e.g., human relevant material, e.g., blood, saliva, urine, breast milk, faecal material).	<input type="checkbox"/>	<input checked="" type="checkbox"/>
	Yes	No
15. ^H ILLEGAL/HARMFUL activities research (e.g., development of technology intended to be used in an illegal/harmful context or to breach security systems, searching the internet for information on highly sensitive topics such as child and extreme pornography, terrorism, use of the DARK WEB, research harmful to national security).	<input type="checkbox"/>	<input checked="" type="checkbox"/>
	Yes	No
16. ^X PERMISSION is required to access premises or research participants.	<input type="checkbox"/>	<input checked="" type="checkbox"/>
	Yes	No
17. ^X PERSONAL DATA PROCESSING (Any activity with data that can directly or indirectly identify a living person). For example data gathered from interviews, databases, digital devices such as mobile phones, social media or internet platforms or apps with or without individuals/owners' knowledge or consent, and/or could lead to individuals/owners being IDENTIFIED or SPECIAL CATEGORY DATA (GDPR) or CRIMINAL OFFENCE DATA.	<input type="checkbox"/>	<input checked="" type="checkbox"/>
	Yes	No
^X PUBLIC WORKS DOCTORATES: Evidence of permission is required for use of works/artifacts (that are protected by Intellectual Property (IP) rights, e.g. copyright, design right) in a doctoral critical commentary when the IP in the work/artifact is jointly prepared-produced or is owned by another body	<input type="checkbox"/>	<input checked="" type="checkbox"/>
	Yes	No
18. ^H RISK OF PHYSICAL OR PSYCHOLOGICAL HARM (e.g., TRAVEL to dangerous places in your own country or in a foreign country (see https://www.gov.uk/foreign-travel-advice), research with NGOs/humanitarian groups in conflict/dangerous zones, development of technology/agent/chemical that may be harmful to others, any other foreseeable dangerous risks).	<input type="checkbox"/>	<input checked="" type="checkbox"/>
	Yes	No
19. ^X SECURITY CLEARANCE – required for research.	<input type="checkbox"/>	<input checked="" type="checkbox"/>
	Yes	No
20. ^X SENSITIVE TOPICS (e.g., anything deeply personal and distressing, taboo, intrusive, stigmatising, sexual in nature, potentially dangerous, etc).	<input type="checkbox"/>	<input checked="" type="checkbox"/>
	Yes	No

M – Minimal Risk; X – More than Minimal Risk. H – High Risk

If you have answered 'Yes' to ANY of the above questions, your application REQUIRES ethical review and approval using the MOREform **BEFORE commencing your research**. Please apply for approval using the MOREform (<https://moreform.mdx.ac.uk/>). Further guidance on making an application using the MOREform can be found at: www.tiny.cc/mdx-ethics.

If you have answered 'No' to ALL of the above questions, your application is Low Risk and you may NOT require ethical review and approval using the MOREform before commencing your research. Your research supervisor will confirm this below.

Student Signature:.....  Date: 30-07-2024

To be completed by the supervisor:

Based on the details provided in the self-assessment form, I confirm that:

Insert
Y or N

The study is Low Risk and <i>does not require</i> ethical review & approval using the MOREform	
The study <i>requires</i> ethical review and approval using the MOREform.	

Supervisr Signature:..... Date:.....