



**Northumbria
University
NEWCASTLE**

Dissertation

KF7029

MSc Computer Science and Digital Technologies Project

Student Name: Arka Mandol
Student ID: w23023023

Supervisor Name: Dr. Mehmet Fatih Tuysuz
Second Supervisor Name: Dr. Biju Issac

MSc Programme: Data Science

Dissertation Title:

Time Series Crime Analysis in the UK: Integrating Explainable AI and Causal Inference

Month and Year of Submission: 09/2024

Academic Publication Venue: IEEE Transactions on Artificial Intelligence

Declaration

I declare the following:

- 1) That the material contained in this dissertation is the end result of my own work and that due acknowledgment has been given in the bibliography and references to ALL sources, be they printed, electronic, or personal.
- 2) The Word Count of this Dissertation is7944
- 3) That unless this dissertation has been confirmed as confidential, I agree to an entire electronic copy or sections of the dissertation being placed on the eLearning Portal (Blackboard), if deemed appropriate, to allow future students the opportunity to see examples of past dissertations. I understand that if displayed on the eLearning Portal it would be made available for no longer than five years and that students would be able to print off copies or download.
- 4) I agree to my dissertation being submitted to a plagiarism detection service, where it will be stored in a database and compared against work submitted from this or any other Department or from other institutions using the service. In the event of the service detecting a high degree of similarity between content within the service, this will be reported back to my supervisor and second marker, who may decide to undertake further investigation that may ultimately lead to disciplinary actions, should instances of plagiarism be detected.
- 5) I have read the Northumbria University Policy Statement on Ethics in Research and Consultancy and I confirm that ethical issues have been considered, evaluated, and appropriately addressed in this research.

SIGNED: Arka Mandol

DATE: 12-09-2024

Time Series Crime Analysis in the UK: Integrating Explainable AI and Causal Inference

Arka Mandol
MSc Data Science
Northumbria University
Newcastle Upon Tyne, United Kingdom
arka.mandol@northumbria.ac.uk

Abstract—In this paper, I propose a novel framework that combines Long Short-Term Memory (LSTM) and Traditional models with Explainable AI (XAI) techniques and causal inference to examine UK crime time series data. The author used the LSTM model to predict during periods of crime and further enhanced interpretability by applying SHAP or LIME. SHAP gave an overall perspective, and LIME was per-instance. The study found that the extent of unemployment and education are major contributors to crime. Used Bayesian Networks for the inferential task of discovering causal mechanisms—or network configurations—between these factors and crime rates. The findings demonstrate how the integration of predictive modeling and causal inference can yield interpretable results for informed crime prediction. Our future work will include streaming external causal factors and combining with advanced causal models to further enhance the prediction accuracy, making our results more useful for policy.

Index Terms—crime prediction, LSTM, Explainable AI, SHAP, LIME, Bayesian Networks, causal inference, Random Forest ,Linear Regression.

CONTENTS

I	Introduction	7
I-A	Explainable AI and Causal Inference in UK Crime Time Series Data Analysis	7
I-B	Problem Statement	7
I-C	Research Objectives	7
I-D	Significance of the Study	8
II	Literature Review	8
II-A	Historical Context and Evolution	8
II-B	Deep Learning in Time Series Forecasting	8
II-C	Explainable AI (XAI) Techniques	8
	II-C1 The Need for Interpretability in AI	8
	II-C2 Local Interpretable Model-agnostic Explanations (LIME)	9
	II-C3 SHapley Additive exPlanations (SHAP)	9
II-D	Causal Inference in Crime Data Analysis	9
	II-D1 Importance of Causal Understanding in Crime Analysis	9
	II-D2 Bayesian Networks in Crime Analysis	9
II-E	Identified Gaps and Research Contributions	9
II-F	Dataset Overview	10
	II-F1 Data on Migrants	10
	II-F2 Data on Education Level	10
	II-F3 Data on Employment History	10
	II-F4 Data on Ethnic Groups	10
	II-F5 Data on Economic Status	10
	II-F6 Crime Information	10
	II-G Integration and Analytical Approach	10
III	Methodology	10
III-A	Crime Data Analysis and Forecasting using LSTM and Explainability Techniques	10
	III-A1 Data loading and Preprocessing	10
	III-A2 Time Series Visualization and Decomposition	10
	III-A3 LSTM Model Development and Training	11
	III-A4 Model Predictions and Evaluation	11
	III-A5 Incorporating External Factors for Robust Predictions	11
	III-A6 Model Interpretability with LIME and SHAP	11
III-B	Random Forest and Linear Regression, along with model interpretability techniques SHAP and LIME	11
	III-B1 Data loading and Preprocessing	11
	III-B2 Merging Datasets and Feature Engineering	12
	III-B3 Visualization of Correlation Matrix	12
	III-B4 Model Training and Train-Test Split	12
	III-B5 Model Evaluation	12
	III-B6 Feature Importance Visualization	12
	III-B7 SHAP Values for Global and Local Interpretation	12
	III-B8 LIME for Local Interpretation	12
III-C	Bayesian Network for Causal Inference on Crime Data	12
	III-C1 Imports and Setup	12
	III-C2 Crime Data Preparation	12
	III-C3 Data Cleaning and Integration	12
	III-C4 Correlation and Distribution Analysis	12
	III-C5 Data Transformation for Bayesian Modeling	13
	III-C6 Bayesian Network Structure Learning	13
	III-C7 Fitting and Evaluating Bayesian Network	13
	III-C8 Inference and Analysis	13

	III-C9	Visualization of Bayesian Network	13
III-D	LSTM Hotspot Mapping	13	
	III-D1	Data Collection and Preprocessing	13
	III-D2	Data Normalization	13
	III-D3	Sequence Creation for LSTM	13
	III-D4	Building and Compiling the LSTM Model	14
	III-D5	Model Training	14
	III-D6	Predictions and Inverse Transformation	14
	III-D7	Visualization with Heatmaps	14
	III-D8	Model Evaluation using RMSE	14
IV	Results and Discussion		14
IV-A	LSTM Crime Trend Prediction and Model Performance	14	
	IV-A1	LSTM Model Performance: Training and Testing	14
	IV-A2	Discussion: Impact of External Disruptions on LSTM Models	14
	IV-A3	Model Interpretability: SHAP vs. LIME	14
IV-B	LSTM Crime Hotspot Prediction	16	
	IV-B1	Model Performance: Root Mean Squared Error (RMSE)	16
	IV-B2	Crime Hotspot Predictions Visualization	16
	IV-B3	Practical Implications for Law Enforcement	16
	IV-B4	Limitations and Potential Improvements	16
IV-C	Crime Counts and Socio-Economic Factors Prediction	17	
	IV-C1	Model Performance: Linear Regression vs. Random Forest	17
	IV-C2	Feature Importance Analysis: SHAP and LIME	17
	IV-C3	SHAP vs LIME, which is Better?	17
	IV-C4	Discussion: Socio-Economic Factors and Crime Prediction	18
IV-D	Discussion: Socio-Economic Factors and Crime Prediction	18	
IV-E	Bayesian Network Analysis on Crime and Socio-Economic Factors	18	
	IV-E1	Influence of "Economy_Not_Employed_binned" on Crime	18
	IV-E2	Influence of "Education_Level_3_binned" on Crime	18
	IV-E3	Influence of "Migrant_Same_Address_binned" on Crime	18
	IV-E4	Influence of "Ethnic_White_British_binned" on Crime	19
	IV-E5	Influence of "Employment_Worked_12M_binned" on Crime	19
	IV-E6	Influence of "Month_Encoded" on Crime	19
	IV-E7	Discussion: Socio-Economic Factors Driving Crime	19
	IV-E8	Policy Implications	19
IV-F	Ethical Considerations	20	
V	Conclusion and Future Works		20
References			20
Appendix A: Dataset Images			21
Appendix B: Annotated Papers			22
Appendix C: LSTM Model Code			22
Appendix D: Bayesian Network Code			26
Appendix E: Feature Importance Analysis Code			33
Appendix F: LSTM Hotspot Mapping			38

I. INTRODUCTION

A. Explainable AI and Causal Inference in UK Crime Time Series Data Analysis

A gradual tendency is a more significant incorporation of Artificial Intelligence (AI) into crime prediction to help ensure public safety and law enforcement activities are done with more accuracy and speed [1]. This is a field that has recently been penetrated by the researchers. Traditional crime analytical methods such as statistical models and regression analyses to the use of Artificial Intelligence (AI) still encounter problems despite the fact that we live in complex realities behind samples, which are affected by various socio-economic, environmental, and temporal factors [2]. The emergence of machine learning, and more particularly, deep learning techniques such as Long Short-Term Memory (LSTM) networks can improve the predictive ability of crime analysis in style. LSTMs, which are able to learn long-range dependencies in sequential data, are particularly suited to issue forecasts of crime trends based on previous data [3]. But still, the models' results are very often "black boxes" even those with quite high accuracy. The practitioners and policymakers cannot easily understand the generated predictions, do not rely on them, and cannot act upon them.

Explainable AI (XAI) methods are a necessity to solve this issue and, thus, to improve the transparency and interpretability of machine learning models. For instance, Local Interpretable Model-agnostic Explanations (LIME) and SHapley Additive exPlanationS (SHAP) are such methods that give insights into the decision-making process of models when predicting the target variable [4]. Interpretability is of great importance in crime analysis. The reason is that the decisions based on AI predictions can have serious societal consequences [4].

The need for more accurate and timely insights to support law enforcement and public safety activities has led to a growing trend of integrating Artificial Intelligence (AI) into crime prediction [1]. Conventional crime analysis techniques like statistical models and regression analyses often struggle in capturing the complicated nature of crime data, which is impacted by different socio-economic, environmental, and temporal factors [2]. The emergence of machine learning, especially deep learning methods such as Long Short-Term Memory (LSTM) networks, has created new opportunities for improving the predictive power of crime analysis. LSTM models, known for their capability to identify long-term dependencies in sequential data, are a good choice for forecasting crime trends from past data [3]. However, despite their accuracy, these models frequently function as "black boxes," making their predictions challenging

to understand and rely on by practitioners and policy-makers.

To address this issue, explainable AI (XAI) methods have become essential for improving the interpretability and transparency of machine learning models. Methods like Local Interpretable Model-agnostic Explanations (LIME) and SHapley Additive exPlanationS (SHAP) provide insights into how models make decisions, thereby helping analysts comprehend the fundamental elements influencing predictions [4]. In crime analysis, where decisions based on AI predictions can have major societal implications, interpretability is very essential [4].

Another vital aspect of crime data analysis is comprehending the causal relationships between different factors and crime rates. Although predictive models like LSTM can predict trends, they do not inherently reveal the underlying reasons for these trends. Causal inference techniques, like Bayesian Networks and Structural Equation Modeling (SEM), provide an organized way to unravel the complex interactions between various variables that contribute to crime [5]. Combining causal analysis with predictive modeling can give a more holistic comprehension of crime dynamics, facilitating informed decision-making [6].

B. Problem Statement

Despite advancements in AI and machine learning, there are still a number of issues with the current crime prediction models [7]. The main issue is that complex models like LSTM provide predictions that are challenging to understand [8], which restricts the actual use of these models in law enforcement [9]. Furthermore, recent models fail to consider the causal relationships between different factors that impact crime, bringing about predictions that may be accurate but not insightful [10]. One way to fix this issue is to integrate XAI methodologies with that of causal inference [6], but this combination has not been thoroughly researched in the field of crime data analysis. Furthermore, despite a growing body of research on the use of AI to crime prediction, there is a noticeable lack of knowledge regarding the difficulties of crime data in the United Kingdom in the literature [11]. A customized approach to predictive modeling and causal analysis is important because of the varied and dynamic nature of crime in the UK, which is impacted by social, economic, and regional factors [10]. This paper tries to bridge this gap by focusing on the analysis of UK crime time series data employing LSTM models together with XAI and causal inference approaches.

C. Research Objectives

The main aim of this research is to create and evaluate a predictive framework that integrates LSTM models with XAI methods and causal inference approaches to

analyze UK crime data. Specifically, this paper will concentrate on the following main objectives:

- To utilize LSTM models trained on historical crime data to forecast trends in crime throughout the United Kingdom.
- To compare the effectiveness of LIME and SHAP in offering interpretable insights for the predictions that LSTM and Traditional models make.
- To find and analyze the causal relationships between various socio-economic indicators and crime rates using Bayesian Networks or SEM.

D. Significance of the Study

This study is anticipated to make several significant contributions to the field of AI and crime analysis. By combining XAI approaches with LSTM models, the research intends to improve the interpretability of crime predictions and make them more useful for law enforcement agencies. The combination of causal inference techniques will yield insights that go beyond simple prediction and a deeper comprehension of the variables influencing crime. These contributions are not only pertinent for academic study, but also hold practical significance for enhancing public safety strategies in the UK [12]. Furthermore, the emphasis on UK crime data, which has distinct regional and socioeconomic features, offers a level of specificity that is often lacking in general AI crime prediction studies. The result of this research could inform policy decisions and resource allocation in law enforcement, which might ultimately result in safer communities.

II. LITERATURE REVIEW

A. Historical Context and Evolution

Crime prediction has long been an area that social sciences and criminology have been interested in [13] [14]. It has traditionally depended on statistical techniques like regression analysis and spatial grouping [14]. Early models concentrated on finding patterns in historical crime data to predict future occurrences. For example, it is common practice to anticipate crime hotspots by doing spatial analysis of crime data using techniques like Kernel Density Estimation (KDE) and hot spot mapping [15].

The transition from traditional statistical methods to machine learning techniques marked a major change in the domain. Machine learning has shown promise in increasing the accuracy of crime forecasts, especially in its capacity to model intricate and non-linear relationships. Research shows that machine learning techniques like Random Forest, Support Vector Machines (SVM), and Gradient Boosting Machines (GBM) can help predict crime hotspots [16].

B. Deep Learning in Time Series Forecasting

Deep Learning models, especially Long Short-Term Memory, have become increasingly popular these days in time series forecasting than other machine learning techniques because they can capture temporal connections in sequential data [17] [18]. [19] put forward LSTM models, a type of Recurrent Neural Networks (RNNs), to address the reducing gradient problem in standard RNNs. These models are particularly suited for tasks that require the use of long-term memory, making them perfect for predicting crime trends that rely on historical patterns [20].

Recent literature has shown promising results in the application of LSTM models in crime prediction. Research by [21] utilized LSTM networks to predict crime rates in smart cities, demonstrating that LSTM models could outperform traditional ARIMA models in terms of accuracy by achieving a validation predictive accuracy of 87% on a large dataset of historical crime data. Similarly, [22] demonstrates that an LSTM encoder-decoder architecture produced greater performance in forecasting crime numbers throughout 83 regions of Costa Rica, with the model performing especially well in areas with higher crime frequency. These results reveal promising outcomes when employing LSTM models to crime prediction.

Though LSTM models are very predictive, they have been criticized for their "black box" nature that makes their internal workings and decision-making processes opaque. This lack of interpretability is a huge obstacle to their widespread use in critical domains like crime prediction, where grasping the factors driving predictions is important for actionable insights [23].

C. Explainable AI (XAI) Techniques

1) *The Need for Interpretability in AI:* The growing complexity of AI models, particularly deep learning networks, has resulted in escalated concerns about their interpretability. In the realm of crime prediction, where AI-driven judgments can have huge societal impacts, the capability to clarify and justify predictions is essential. This demand for interpretability has led to the development of Explainable AI (XAI), which aims to make AI models more transparent and easier for people to understand [12].

Although the terms interpretability and explainability are frequently used interchangeably in the field of AI, they have different meanings: interpretability relates to comprehending the internal logic and mechanics of the model, whereas explainability is concerned with understanding the intuition behind a model's outputs and cause-and-effect relationships. As highlighted by [12], interpretability alone is insufficient without explainabil-

ity, both are necessary in high-stakes applications like crime prediction to ensure openness and accountability.

There are two primary approaches to Explainable AI (XAI): ante-hoc and post-hoc methods. After a model is developed, post-hoc XAI uses methods to explain its predictions. For complicated models such deep learning networks, it uses techniques like GradCAM, LIME, and SHAP [24]. These methods might not, however, always faithfully capture the model's decision-making process. Ante-hoc XAI, on the other hand, concentrates on creating models that are naturally interpretable from the start, including decision trees and linear models, which offer transparent insights into their predictions. Ante-hoc models allow users to understand the logic behind the model's decisions, even though they may trade some predictive accuracy in the name of openness [25]. Depending on the required balance between accuracy and interpretability, one can choose amongst these strategies.

2) *Local Interpretable Model-agnostic Explanations (LIME)*: LIME, developed by Ribeiro, Singh, and Guestrin in 2016 [26], is a widely used technique that approximates the decision boundary of a complicated model locally around a given prediction. LIME creates a simple, interpretable model (often a linear model) that replicates the behavior of the complicated model around the instance by perturbing the input data and tracking the changes in the model's predictions. With this method, users can comprehend which features had the most influence and why the model predicted a certain outcome [27].

3) *SHapley Additive exPlanations (SHAP)*: SHAP, developed by Lundberg and Lee in 2017 [27], is another prominent technique for explaining how AI models make predictions. Based on principles of cooperative game theory, SHAP calculates the contribution of each feature to a prediction by considering all possible combinations of features. This method gives a consistent and theoretically based way to rank the value, making it particularly suitable for understanding complex models.

SHAP has been utilized in different domains, including crime prediction, to provide insight on how models reach their decisions. In a recent study, researchers employed SHAP, to analyze and evaluate the importance of various features within clustered groups of data. By combining k-means clustering, PCA, and a LightGBM classifier during the preparation of data, the study has grouped features into multiple clusters before using SHAP. The simulation results show that SHAP values play an important role in developing accurate and meaningful representations of each cluster, showing the most significant features influencing the model's predictions [28]. However, despite how good they are, both have limitations. While SHAP's processing complexity may

be unaffordable for very big datasets or extremely complicated models like LSTMs, LIME's local explanations may occasionally be deceptive if the model acts nonlinearly in other regions of the input space [28].

D. Causal Inference in Crime Data Analysis

1) *Importance of Causal Understanding in Crime Analysis*: Predictive models are useful for predicting crime trends based on historical data, but they frequently don't provide details about the underlying causes of these trends. Understanding the causal relationships between various factors and crime rates is pivotal for effective crime prevention measures. For instance, knowing that greater unemployment rates lead to higher crime rates can guide policymakers in designing targeted interventions. Causal inference techniques offer a body for studying these relationships, moving beyond simple correlation to comprehend causality. This distinction is crucial in crime analysis, where interpositions that rely only on correlations may result in ineffective or counterproductive policies.

2) *Bayesian Networks in Crime Analysis*: Bayesian Networks are robust graphical models that encode probabilistic relationships between variables, making them better for analyzing complex interactions in crime data. These networks are quite beneficial for data analysis, especially when it comes to modeling and understanding the causal processes that can result in illegal activities. By encoding variables and their dependencies using a directed acyclic graph (DAG), Bayesian Networks allow researchers to investigate and predict the effects of several factors impacting crime.

Bayesian networks have been used in several studies to analyze crime data. For instance, [29] showed how these networks can capture small correlations between offender characteristics and crime scene evidence by modeling criminal profiling from limited data using Bayesian networks. The study demonstrated how well Bayesian Networks could manage vagueness and missing data in criminal investigations. Similarly, [30] used a spatiotemporal Bayesian model to show trends in burglaries, investigate correlations between crimes and socioeconomic indicators, and examine distribution patterns, demonstrating how useful the model is for figuring out causal relationships in safety-related data.

E. Identified Gaps and Research Contributions

Even with the advances in causal inference, XAI, and AI, there are still several gaps in the literature. Firstly, although LSTM models have showed potential in the field of crime prediction, their interpretability issues contain their real world use. Secondly, there is a shortage of research on the application of XAI approaches to causal inference methods in the study of crime data,

especially when it comes to UK crime data. Finally, more targeted research is required to address the difficulties in analysing crime statistics in the UK while considering particular socioeconomic and regional aspects.

By creating a comprehensive framework that combines LSTM models, XAI approaches, and causal inference methods to analyse UK crime time series data, this research aims to close these gaps. A more understandable and practical crime prediction model, insights into the causal connections underlying crime patterns, and a tailored method for handling the inherent challenges of UK crime data analysis. This framework will be judged based on its capacity to support both academic research and real-world applications in public safety by offering precise forecasts, comprehensible insights, and a thorough grasp of the causative variables impacting crime.

F. Dataset Overview

1) *Data on Migrants:* The ONS's migration dataset, which records UK and international migration, is essential for studying demographic changes and local crime rates.

2) *Data on Education Level:* This ONS dataset categorises people by their highest education level, helping analyse educational disparities' impact on crime and apply explainable AI to find patterns.

3) *Data on Employment History:* ONS data classifies people into employed, unemployed, or outside the labour force, allowing causal inference to examine crime rates and employment movements.

4) *Data on Ethnic Groups:* This ONS dataset divides the UK population by ethnicity, allowing crime statistics to be analysed across ethnic groups and causal connections to be drawn.

5) *Data on Economic Status:* The ONS's economic status dataset looks at people's economic activity, which is important for using causal analysis to investigate the connection between economic disparity and crime.

6) *Crime Information:* The main dataset for crime prediction and causal inference in this project is the "UK Police Street Crime" dataset from Kaggle, which covers street-level offences from 2018 to 2021.

G. Integration and Analytical Approach

All datasets are combined in this project to provide a comprehensive examination of the factors that contribute to crime in the UK. The objective is to identify and quantify how demographic, educational, economic, and ethnic factors impact crime using causal inference methodologies. Applying explainable AI (XAI) techniques will ensure that the models developed are understandable and accurate, enabling stakeholders to comprehend the fundamental causes of crime and make informed decisions.

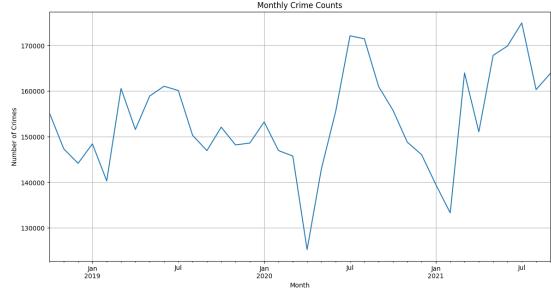


Fig. 1. TimeSeriesVisualizationofCrimeData

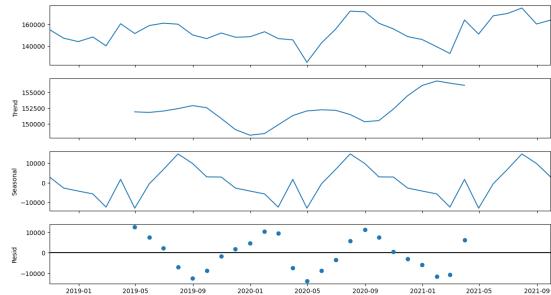


Fig. 2. Time Series Decomposition

III. METHODOLOGY

A. Crime Data Analysis and Forecasting using LSTM and Explainability Techniques

This section presents the approach utilized to forecast crime data employing interpretable machine learning algorithms and a Long Short-Term Memory (LSTM) model.

1) *Data loading and Preprocessing:* UK crime dataset was loaded and then filtered by "Violence and sexual offenses" and checked for missing values.

Time Series Transformation: To aggregate crime counts, the data was resampled to a monthly frequency and the 'Month' column was transformed into a datetime format to capture broader trends and seasonal variations, simplifying complexity and aligning with typical crime reporting frequencies.

Resampling: To capture broader trends and seasonal fluctuations, monthly aggregate was selected.

2) *Time Series Visualization and Decomposition:* Line plots were used to analyze the cleaned time series data to understand general crime patterns.

The `seasonal_decompose` method from `statsmodels` split the time series into trend (long-term crime changes), seasonal (recurring yearly patterns), and residual (unexplained noise) to uncover deeper data insights [31].

$$X_t = T_t + S_t + R_t \quad (1)$$

Where: - X_t is the observed time series. - T_t is the trend component. - S_t is the seasonal component. - R_t is the residual component.

Decomposition helps in the identification of long-term trends and seasonal patterns that can inform proactive crime prevention intervention tactics [32].

3) LSTM Model Development and Training: The LSTM model was used because of its ability to capture long-term dependencies in time series data [33]. The following steps were used:

Data Preparation: To forecast the number of crimes for the upcoming month, the data was transformed into a supervised learning format using 12-month sequences. To ensure robust training, features were normalized to the range [0, 1] using MinMaxScaler.

Model Architecture: To avoid overfitting, dropout layers (rate 0.3) and two LSTM layers, each with 50 units, were employed. The estimated crime count was obtained by adding a fully connected dense layer.

$$C_t = f_t \times C_{t-1} + i_t \times \tilde{C}_t \quad (2)$$

Where: - C_t is the cell state. - f_t is the forget gate. - i_t is the input gate. - \tilde{C}_t is the candidate cell state.

Model Compilation and Training: The adaptive learning rate and Mean Squared Error (MSE) loss function of the model were compiled using the Adam optimizer. To improve learning and provide frequent weight updates, it was trained for 500 epochs with a batch size of 1.

Incorporating lag features, seasonal features, and temporal covariates can improve the LSTM model's ability to capture long-term dependencies and seasonal crime trends [34].

4) Model Predictions and Evaluation: Predictions were evaluated using MSE, Root Mean Squared Error (RMSE), and Mean Absolute Error (MAE):

- **Mean Squared Error (MSE):** Measures the average squared difference between predicted and actual values [35].

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (3)$$

Where: - n is the number of observations. - y_i is the actual value. - \hat{y}_i is the predicted value.

- **Root Mean Squared Error (RMSE):** Provides an interpretable metric in the same units as the data.

$$\text{RMSE} = \sqrt{\text{MSE}} \quad (4)$$

Where: - MSE is the Mean Squared Error.

- **Mean Absolute Error (MAE):** Measures the average absolute difference between predicted and actual values.

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |Y_i - \hat{Y}_i| \quad (5)$$

Where: - n is the number of observations. - Y_i is the actual value. - \hat{Y}_i is the predicted value.

To improve model generalization and mitigate overfitting, L1/L2 regularization, dropout, and data augmentation methods like synthetic data generation or bootstrapping can be implemented to enhance the model's ability to learn diverse patterns [36], [37].

5) Incorporating External Factors for Robust Predictions: To improve robustness, external factors such as economic shifts, public events, and policy changes can be integrated as additional features in the LSTM model to capture their temporal effects on crime rates [16].

6) Model Interpretability with LIME and SHAP: Two interpretability techniques were used to improve the predictability of LSTMs:

LIME: LIME highlighted which months (time steps) had the greatest influence on a particular forecast by approximating the LSTM model with a simpler one for individual predictions. For expeditious criminal assessments, this method provides rapid, actionable information.

SHAP: This method provided both local and global interpretability by quantifying the contribution of each time step to the forecasts. To ensure consistency and fairness in feature attribution, SHAP summary plots were used to identify the most influential months for each forecast.

$$\phi_i = \sum_{S \subseteq N \setminus \{i\}} \frac{|N|!}{|S|!(|N| - |S| - 1)!} [f(S \cup \{i\}) - f(S)] \quad (6)$$

Where: - ϕ_i is the SHAP value for feature i . - S is a subset of features excluding i . - N is the total set of features. - $f(S)$ is the model output for subset S .

Both methods offered vital information on how historical crime data affected model projections, facilitating quick decision-making and guaranteeing impartiality in research.

B. Random Forest and Linear Regression, along with model interpretability techniques SHAP and LIME

In addition to SHAP and LIME for interpretability, this section describes how to create and evaluate crime prediction models using Random Forest and Linear Regression.

1) Data loading and Preprocessing: UK Police crime data for 2021, filtered for "Violence and sexual offenses." The 'Month' column was converted to datetime, and the 'LSOA name' was renamed. The data was combined by month and renamed column.

Socio-economic indicators were standardized by geographical identifiers and merged with the crime data using an outer join. SHAP and LIME were applied to the

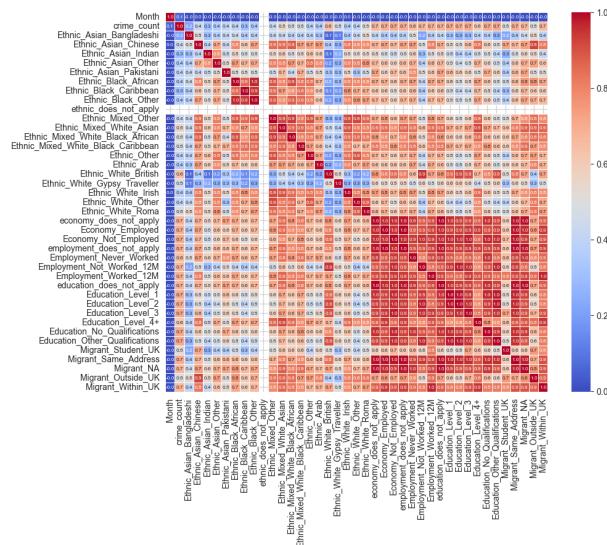


Fig. 3. Correlation Matrix of All Data Variables

models for interpretability to identify key socioeconomic factors influencing crime predictions.

2) Merging Datasets and Feature Engineering: The socioeconomic and criminal databases were combined, with relevant columns being renamed and redundant ones being eliminated. Among the features engineered were:

Feature Selection: Choosing the socioeconomic features that have the strongest correlation with crime count.

Feature Scaling: Apply StandardScaler to standardize a subset of features.

Time Features: To capture seasonality, convert the 'Month' column into a category variable.

3) Visualization of Correlation Matrix: The most important variables were highlighted in the correlation matrix and heatmap that were made to illustrate the connections between socioeconomic characteristics and crime rates.

4) Model Training and Train-Test Split: Models were trained on the training data and tested on the reserved data to assess performance. The data was divided into 70% training and 30% testing subsets.

a) Model Selection: Random Forest Regressor:
The Random Forest Regressor is a popular choice for complicated datasets due to its exceptional capacity to capture non-linear relationships and feature interactions.

Linear Regression: Linear Regression was selected as a benchmark model because it is straightforward and can produce comprehensible linear connections between characteristics and crime counts.

5) *Model Evaluation:* To evaluate the models, two key metrics were utilized:

Mean Squared Error (MSE) was used, It facilitates the quantification of prediction error.

R-squared (R^2): Represents the proportion of the variance in the dependent variable that is predictable from the independent variables.

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{\sum_{i=1}^n (\hat{y}_i - y_i)^2} \quad (7)$$

Where: - n is the number of observations. - y_i is the actual value. - \hat{y}_i is the predicted value. - \bar{y} is the mean of actual values.

The test data was used to evaluate both models, and the models' MSE and R² scores were compared to ascertain which one performed better.

6) *Feature Importance Visualization*: The linear link between characteristics and crime counts was demonstrated by plotting absolute values of the coefficients that were computed and visualized for Random Forest and Linear Regression Feature Importance.

7) *SHAP Values for Global and Local Interpretation:* SHAP values were used for interpretability in both models. TreeExplainer revealed global feature importance in Random Forest and shap.LinearExplainer gave local and global insights for Linear Regression.

8) *LIME for Local Interpretation*: LIME provided local insights for both models, simplifying Random Forest predictions and clarifying individual feature contributions in Linear Regression, making each prediction easier to understand.

C. Bayesian Network for Causal Inference on Crime Data

The process for building and evaluating a Bayesian network to determine the causal links between socioeconomic characteristics and crime rates is described in this section.

1) Imports and Setup: Main Python libraries utilized for data manipulation are pandas and numpy, for visualization are matplotlib and seaborn, and for Bayesian Network development pgmpy & networkx was used.

2) *Crime Data Preparation:* Converted 'Month' to datetime format and filtered for relevant 2021 crime data, focusing on "Violence and sexual offenses."

3) *Data Cleaning and Integration:* Aggregated crime data by location, integrated socioeconomic data (migration, employment, etc.), and cleaned for modeling.

4) Correlation and Distribution Analysis: Visualized correlation relationships with a heatmap and analyzed key variable distributions using histograms to identify trends for Bayesian modeling.

$$r = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum(x_i - \bar{x})^2} \sqrt{\sum(y_i - \bar{y})^2}} \quad (8)$$

Where: - x_i and y_i are individual data points. - \bar{x} and \bar{y} are the means of the x and y variables.

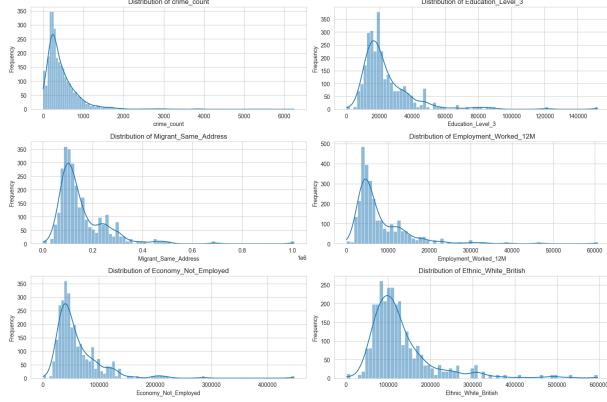


Fig. 4. Distribution of key variables related to crime.

5) *Data Transformation for Bayesian Modeling:* Logarithmic transformation was applied to normalize skewed variables, reducing the influence of outliers for more robust Bayesian Network modeling. Quantile binning then simplified the log-transformed data, categorizing variables like crime counts into "Low," "Medium," and "High," aiding discrete relationship modeling. Additionally, the 'Month' column was encoded categorically to capture seasonal crime trends.

$$y' = \log(y + 1) \quad (9)$$

Where: - y is the original variable. - y' is the log-transformed variable.

6) *Bayesian Network Structure Learning:* Using the K2 score method to guide the search, the Hill-Climb Search algorithm was used to discover the structure of the Bayesian Network. 'Month_Encoded' was one of the variables that the BIC score rejected owing to its complexity penalty, hence K2 was selected because it permitted the inclusion of all pertinent variables.

7) *Fitting and Evaluating Bayesian Network:* The Bayesian Estimator was used to fit the learnt structure to the data. BIC and K2 scores were used to assess the model:

- BIC: model fit is measured, and complexity is penalized.
- K2: maximizes data fit while lessening the penalty for complexity.

The dual evaluation aids in achieving a trade-off between complexity and model accuracy.

8) *Inference and Analysis:* Key factors contributing to high crime rates were determined by inference utilizing the Variable Elimination approach, including socioeconomic status and employment.

$$P(C | E) = \frac{P(E | C) \cdot P(C)}{P(E)} \quad (10)$$

Where:

- $P(C | E)$ is the probability of the crime count being in a specific bin given evidence E .
- $P(E | C)$ is the probability of observing the evidence given the crime count.
- $P(C)$ is the prior probability of the crime count bin.
- $P(E)$ is the probability of observing the evidence.

9) *Visualization of Bayesian Network:* The networkx package was used to visualize the topology of the Bayesian network in a circular form, with larger nodes and thicker edges signifying important interactions.

D. LSTM Hotspot Mapping

1) *Data Collection and Preprocessing:* The UK's crime dataset was utilized, which includes details of crime occurrences that contain geo-coordinates latitude and longitude as well as timestamps in monthly intervals. The data was imported from a CSV file and processed for time-series modeling.

- **Date Formatting and Sorting:** Converting the 'Month' column to datetime format, and sorting the dataset chronologically was critical so that the temporal order was maintained. This step guarantees proper time-series forecasting by the LSTM model.
- **Aggregating Crime Data:** By grouping the dataset by 'Month', 'Latitude', and 'Longitude', a spatial-temporal dataset was created, with each row showing the crime count for a specific location and time. This aggregation step is important for preparing the data, which helps the LSTM model to understand spatial and temporal crime patterns [38].

2) *Data Normalization:* MinMax() scaling normalized crime counts between 0 & 1, ensuring consistent feature scaling for stable & efficient LSTM model training.

$$x'_i = \frac{x_i - x_{\min}}{x_{\max} - x_{\min}} \quad (11)$$

Where:

- x_i is the crime count for location i .
- x_{\min} and x_{\max} are the minimum and maximum crime counts in the dataset.

3) *Sequence Creation for LSTM:* The dataset underwent the process of transforming it to a series of input-output pairs, using a sliding window approach. The model used a sliding window of the past three months' crime counts to predict the next month, with input for each time step derived.

$$X_t = [C_{i,t-2}, C_{i,t-1}, C_{i,t}] \quad \text{to predict} \quad y_t = C_{i,t+1} \quad (12)$$

Where: - X_t is the feature vector containing past crime counts. - $C_{i,t-2}, C_{i,t-1}, C_{i,t}$ are the crime counts for location i in previous months. - $y_t = C_{i,t+1}$ is the crime count for location i at time $t + 1$.

4) *Building and Compiling the LSTM Model:* Using Keras, the LSTM model was built with 50 units for capturing temporal crime patterns and a single-neuron dense layer was used to predict the next month's crime count.

5) *Model Training:* A balanced approach between computational efficiency and model efficiency was adopted for the model training utilizing 10 epochs and batch size of 512. The decision of a smaller number of epochs was adopted so as to prevent overfitting while ensuring that the model learns from the data well.

6) *Predictions and Inverse Transformation:* Subsequently, the model predicted the crime counts for the next year's of the training period. In the case where the data were normalized before training, thus the predicted values were converted back to the original scale through inverse MinMax scaling.

$$p_i = p'_i \cdot (x_{\max} - x_{\min}) + x_{\min} \quad (13)$$

Where: - p'_i is the normalized prediction. - x_{\min} and x_{\max} are the minimum and maximum crime counts in the dataset.

7) *Visualization with Heatmaps:* To visualize the predicted crime hotspots, a heatmap was created using the Folium library. A crime count was plotted on the map for each longitude and latitude. Each location's predicted crime count is shown by the intensity of the heatmap.

8) *Model Evaluation using RMSE:* The model's effectiveness was measured by Root Mean Squared Error (RMSE), which reflects the average prediction error. This metric is crucial for understanding how well the model predicts actual crime occurrences.

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2} \quad (14)$$

The closer the model's predictions are to the actual values, the lower the RMSE score and the better the model performance [39]. This metric is particularly valuable in evaluating the accuracy of forecasted crime counts, providing a clear measure of the model's predictive power and guiding further refinement.

IV. RESULTS AND DISCUSSION

A. LSTM Crime Trend Prediction and Model Performance

The LSTM model, trained on past monthly crime data, leverages historical trends to forecast future crime patterns.

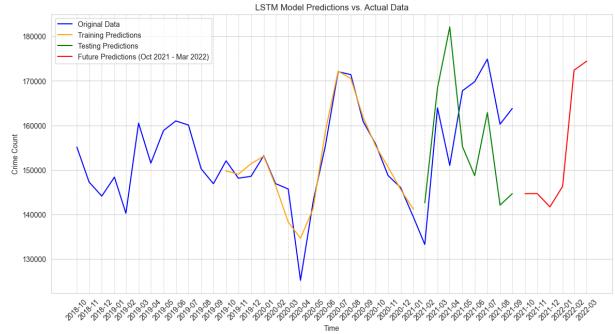


Fig. 5. LSTM Model Predictions vs. Actual Crime Data (Train, Test, and Future Predictions)

1) *LSTM Model Performance: Training and Testing:* During the training phase, the LSTM model demonstrated good performance, accurately capturing temporal changes in crime. The model was able to fit the training data accurately, as evidenced by its Train Mean Squared Error (MSE) of 11,264,257.07 and Train Root Mean Squared Error (RMSE) of 3,356.23. On the test set, however, the model's performance declined, with Test MSE of 314,326,459.66 and Test RMSE of 17,729.25, indicating a notable rise in error when exposed to unseen data.

The significant rise in errors during the testing period can probably be attributed to external interruptions like the COVID-19 pandemic, which brought unusual patterns of criminality [40]. The pandemic caused significant changes in the UK's crime rates, including sudden declines and increases that the model found difficult to account for.

2) *Discussion: Impact of External Disruptions on LSTM Models:* The pandemic drastically affected crime rates introducing noise & unpredictability, the LSTM model struggled to handle effectively. This begs the questions of how resilient classic time series models like LSTM are to such disturbances and if model resilience may be enhanced by including external influences (such economic shifts or public events).

The model's predictions and actual crime counts during the test phase differed, as shown by the visualization in Figure 6. This was especially true during the pandemic months when lockdowns caused a dramatic decline in crime rates.

3) *Model Interpretability: SHAP vs. LIME:* Enhancing the interpretability of the LSTM crime trend predictions using Explainable AI (XAI) techniques—specifically, SHAP and LIME—was one of the main goals of this work. These techniques were employed to clarify the ways in which various time lags (prior crime data points) influenced the predictions made by the LSTM model.

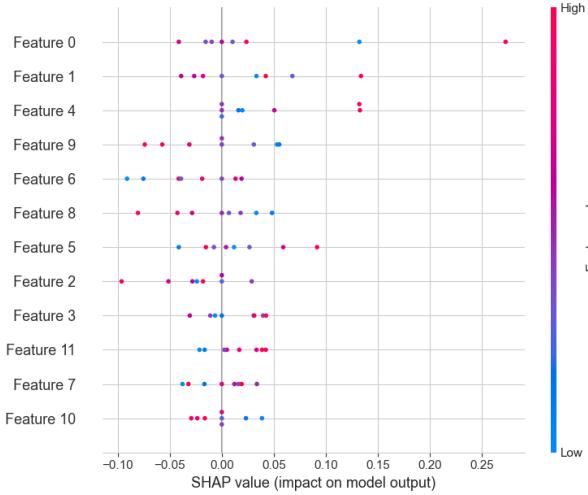


Fig. 6. SHAP Summary Plot Showing Global Feature Importance in LSTM Predictions

a) SHAP: Global and Local Explanations: SHAP was used to give both global and local interpretability by calculating the contribution of all features (time lag) to the model’s predictions. Figure ?? SHAP summary plot emphasizes how crucial recent crime data points—like t-0, t-1, and t-4—are to the model’s ability to make predictions. These characteristics consistently had the highest SHAP values, suggesting that the most important factors for predicting future crime trends are current criminal incidents. The impact of earlier time lags (t-7 to t-9) on the predictions was decreasing.

Understanding the global trends in how the LSTM model uses historical crime data is made possible by SHAP’s capacity to consistently and reliably deliver feature importance throughout the whole dataset. For example, SHAP disclosed that the most recent time point, t-0, positively impacted crime predictions. This implies that a high crime rate in the current month would probably lead to a higher prediction for the following month.

b) LIME: Instance-Specific Explanations: For each prediction, instance-specific explanations were produced using LIME. Here, LIME modifies the input data and tracks how slight modifications impact the model’s output, yielding a breakdown of feature contributions for each prediction.

For instance, LIME breaks clearly how recent crime counts (t-0, t-1) contribute to an LSTM prediction for a particular instance in Figure ???. LIME demonstrated how recent crime spikes affected the model’s output by capturing the local significance of characteristics like t-0 and t-1. To make short-term crime estimates, the LSTM model relies heavily on current data points; older time points (such as t-8 and t-9) contributed less to this.

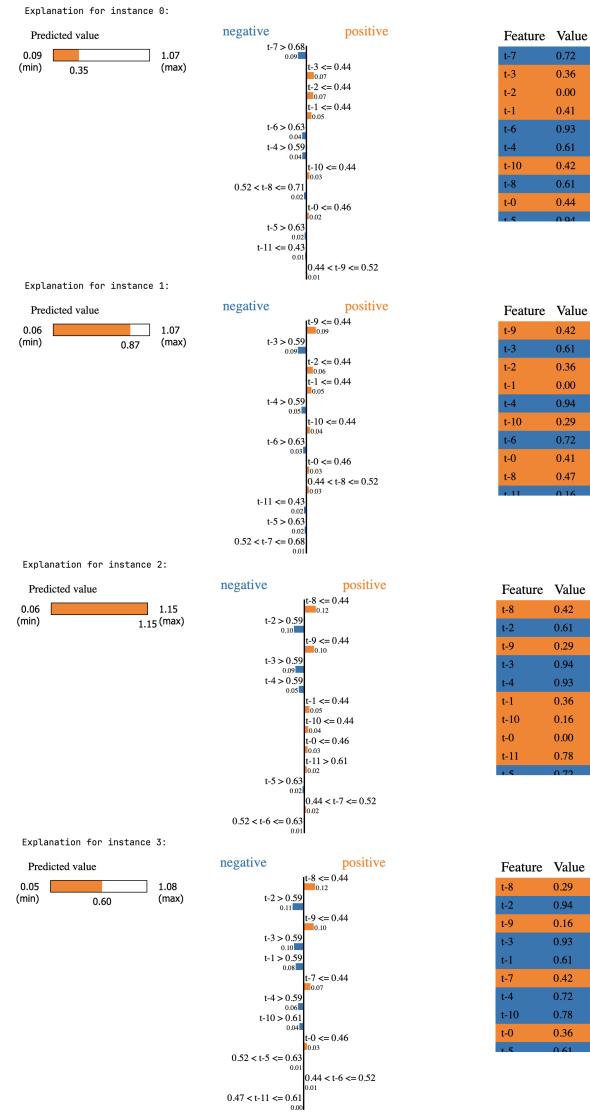


Fig. 7. LIME Explanation for a Specific LSTM Prediction (First 4 Instances)

Although LIME’s instance-specific insights are helpful in comprehending specific predictions, small variations between runs may result from the randomness induced by disturbance. Because of this, LIME is better suited for local interpretability than for identifying broad patterns in the dataset.

c) SHAP vs LIME, which is Better?: When the two methods were compared, SHAP proved to be the most effective way to understand the crime predictions made by the LSTM model. SHAP was better suited for comprehending the behavior of the model over the entire dataset because of its capacity to offer consistent, theoretically supported global and local insights. In order to explain time-series models like LSTM, where the link between time lags and future predictions might be

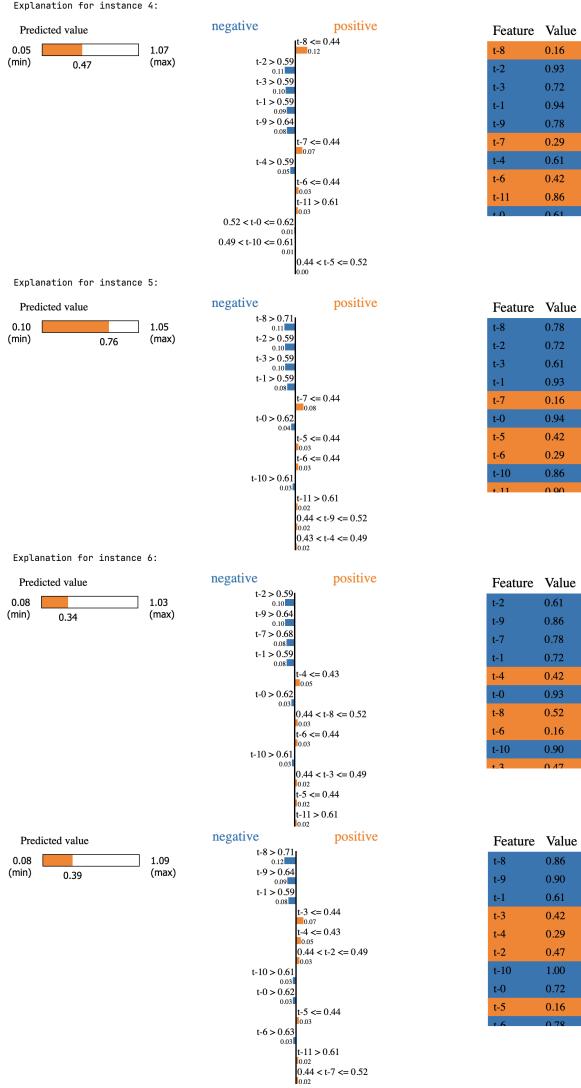


Fig. 8. LIME Explanation for a Specific LSTM Prediction (Last 4 Instances)

complex and non-linear, SHAP also makes sure that the significance values are constant across instances.

Here, SHAP should be preferred for interpretability as it provided clearer, more consistent insights into how time lags influenced the LSTM predictions.

B. LSTM Crime Hotspot Prediction

The findings of predicting crime hotspots in different UK regions using the LSTM model are shown in this section. Finding locations with a high likelihood of increased criminal activity was the main objective in order to help law enforcement organizations allocate resources more wisely.

1) Model Performance: Root Mean Squared Error (RMSE): The Root Mean Squared Error (RMSE) metric

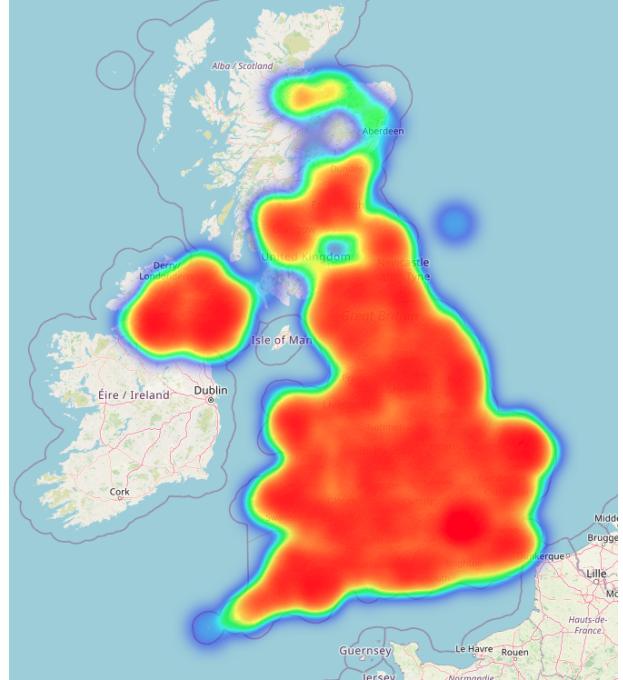


Fig. 9. Crime Hotspot Prediction Heatmap for the First Month

was used to assess the LSTM model's hotspot prediction performance. With an RMSE of 1.572, the model demonstrated a respectable degree of accuracy in estimating the number of crimes per site for the first month of the study. According to this RMSE value, there was an average deviation of 1.57 crimes per region between the model's predictions and the actual crime counts.

Although this performance is encouraging, there is still space for improvement, particularly in areas where crime rates are lower and the model shows somewhat higher prediction errors. This is in line with the difficulties time series models frequently encounter when forecasting results in low-variance areas.

2) Crime Hotspot Predictions Visualization: Figure 9 shows the expected crime hotspots for the first month. Areas with higher predicted crime activity are shown by red shading. The algorithm identified major hotspots as urban centers like London, Birmingham, and Manchester, as anticipated, which is in line with past trends in crime (Hart, 2020).

3) Practical Implications for Law Enforcement: Law enforcement organizations can take action based on the predictions made by the LSTM model, especially when it comes to allocating resources. Law enforcement can more efficiently allocate resources, speeding up reaction times and possibly averting crime in high-crime regions discovered by the model.

4) Limitations and Potential Improvements: Although the low RMSE suggests that the LSTM model for hotspot

prediction did quite well, there are a few issues that could be fixed in next iterations of the model:

- **Prediction Granularity:** At the moment, the model offers predictions at the regional level. More accurate insights might be obtained by predicting at the neighborhood level, for example, to increase the granularity of the projections.
- **External Factors and Long-Term Validation:** The model, relying solely on past crime data, overlooks external factors like economic shifts or public events, which could enhance accuracy. Longer validation periods may offer deeper insight into its predictive abilities.

C. Crime Counts and Socio-Economic Factors Prediction

This section compares crime forecasts from Linear Regression, a simple baseline, and Random Forest, which captures non-linear patterns, using SHAP and LIME to reveal the impact of socioeconomic and demographic factors.

1) Model Performance: Linear Regression vs. Random Forest: The Random Forest model, with an MSE of 7421.12 and R² of 0.9603, explained 96% of the crime count variance, handling complex non-linear relationships effectively. In contrast, the Linear Regression model, with a much higher MSE of 85258.72 and R² of 0.5433, struggled to capture these interactions, performing significantly worse.

TABLE I
PERFORMANCE METRICS FOR RANDOM FOREST AND LINEAR REGRESSION MODELS

Model	MSE	R ²
Random Forest	7421.12	0.9603
Linear Regression	85258.72	0.5433

2) Feature Importance Analysis: SHAP and LIME: The Random Forest and Linear Regression models' predictions were interpreted using both SHAP and LIME, revealing which socioeconomic and demographic variables had the biggest impact on crime rates.

a) SHAP Analysis: International and Regional Perspectives: Figure 10 shows SHAP values indicated that Employment_Worked_12M, which positively impacted crime rates, was the most significant feature for the Random Forest model. This indicates that areas with more employment during the previous 12 months may have had higher crime rates as a result of sporadic or insecure work circumstances.

Other notable characteristics, such as Economy_Not_Employed and Education_Level_3, also made a substantial contribution to the forecasts.

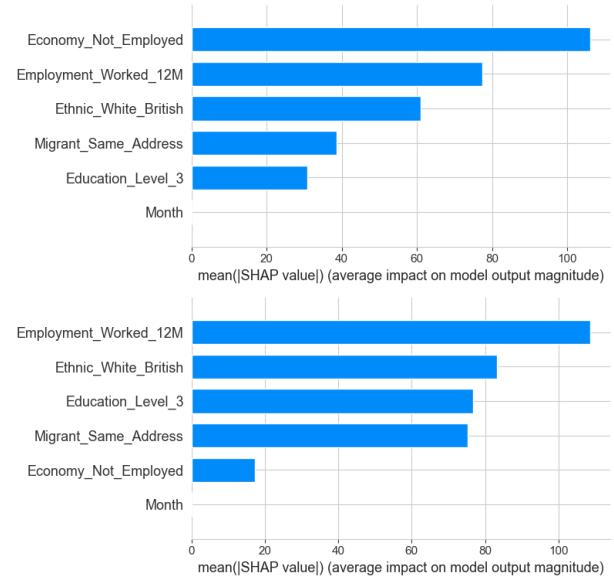


Fig. 10. SHAP Summary Plot for Feature Importance in Random Forest (top) and Linear Regression (bottom) Models

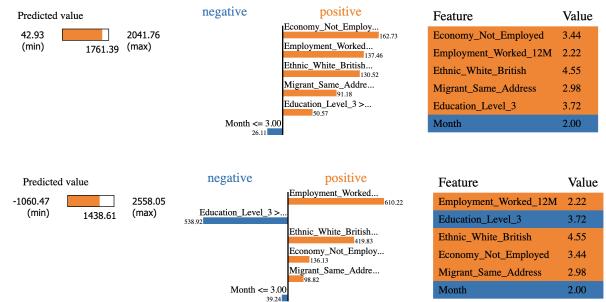


Fig. 11. LIME Explanation for Random Forest Prediction and Linear Regression

b) LIME Analysis: Instance-Level Explanations: For cases in both models, local explanations were produced using LIME. The contribution of features like Employment_Worked_12M and Ethnic_White_British to a specific Random Forest model prediction is displayed in Figure 11 by LIME. For instance, Education_Level_3 contributed negatively by lowering the projected crime count, whereas Employment_Worked_12M contributed favourably to a crime prediction of 1438.61.

3) SHAP vs LIME, which is Better?: SHAP and LIME both provided insightful analyses of the models' interpretability. SHAP, by offering both global and local feature importance, helps explain how socioeconomic factors shape crime rates, providing the clarity and consistency needed for stakeholders to trust model insights. LIME was very helpful in providing instance-specific explanations, including brief descriptions of the meth-

ods used to make individual predictions. Nevertheless, because feature importance might fluctuate somewhat between runs, its dependence on perturbation rendered it less trustworthy for identifying global patterns.

4) Discussion: Socio-Economic Factors and Crime Prediction: The Random Forest model's findings point to a number of important socioeconomic variables that affect crime rates. The robust significance of Employment_Worked_12M implies a close relationship between crime rates and employment patterns, especially in areas where there have been recent swings in employment. This result can be a reflection of the erratic nature of contract or temporary work, which can put financial strain on people and raise crime rates.

Higher education levels were found to be linked to decreased crime rates, whereas Education_Level_3 was found to have a protective impact against crime. This is consistent with other research that indicates increased education increases opportunities and lowers the risk of criminal activity.

The more conventional criminological theory that points to unemployment as a major contributing factor to crime is reflected in the Linear Regression model's emphasis on Economy_Not_Employed.

Policymakers can better target crime prevention initiatives by concentrating on communities with lower levels of education and insecure employment, according to these insights.

D. Discussion: Socio-Economic Factors and Crime Prediction

The Random Forest model's findings point to a number of important socioeconomic variables that affect crime rates. The robust significance of Employment_Worked_12M implies a close relationship between crime rates and employment patterns, especially in areas where there have been recent swings in employment. This result can be a reflection of the erratic nature of contract or temporary work, which can put financial strain on people and raise crime rates.

Higher education levels were found to be linked to decreased crime rates, whereas Education_Level_3 was found to have a protective impact against crime. This is consistent with other research that indicates increased education increases opportunities and lowers the risk of criminal activity.

The more conventional criminological theory that points to unemployment as a major contributing factor to crime is reflected in the Linear Regression model's emphasis on Economy_Not_Employed.

Policymakers can better target crime prevention initiatives by concentrating on communities with lower levels of education and insecure employment, according to these insights.

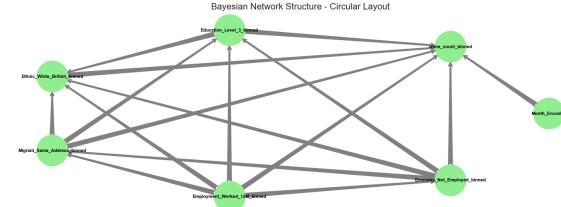


Fig. 12. Advance Bayesian Network Structure

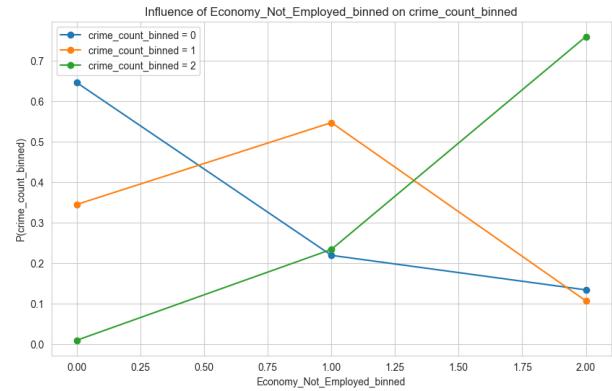


Fig. 13. Bayesian Network - Influence of "Economy_Not_Employed_binned" on Crime Counts

E. Bayesian Network Analysis on Crime and Socio-Economic Factors

The Bayesian Network analysis, which was done to determine the probability correlations between crime rates and several socioeconomic indicators such as unemployment, education, migratory stability, and ethnic composition, is shown in this section. These components and their conditional dependencies are represented graphically by the Bayesian technique, which aids in the discovery of causal correlations that can guide crime prevention strategies.

1) Influence of "Economy_Not_Employed_binned" on Crime: Figure 13 shows low crime rates dropped steadily as unemployment increased, but the likelihood of high crime rates increased in regions with high unemployment.

2) Influence of "Education_Level_3_binned" on Crime: Figure 14 illustrates how, according to conventional wisdom, higher education levels were associated with a higher likelihood of rising crime rates. This unexpected finding raises the possibility that additional variables, such as urbanization, are present in these areas.

3) Influence of "Migrant_Same_Address_binned" on Crime: Figure 15 shows that areas with greater migrant stability—that is, people who lived at one residence for a prolonged amount of time—were also more likely to have higher rates of crime.

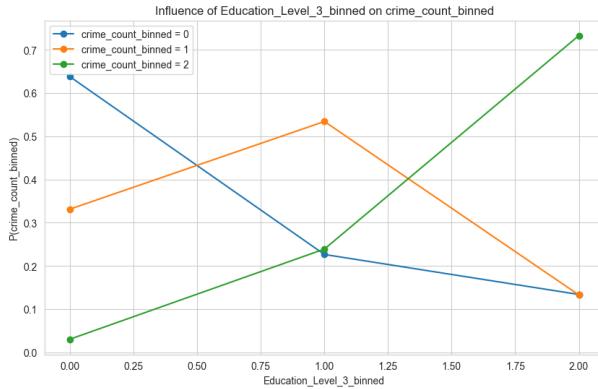


Fig. 14. Bayesian Network - Influence of "Education_Level_3_binned" on Crime Counts

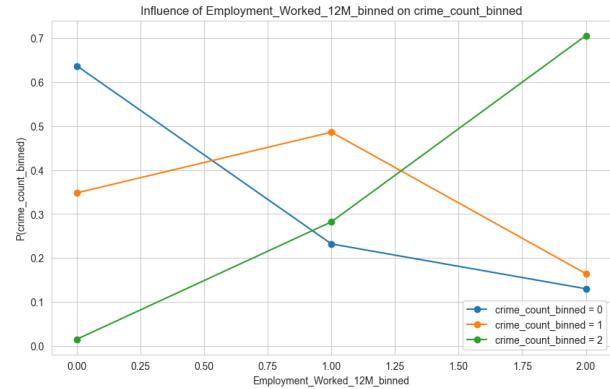


Fig. 17. Bayesian Network - Influence of "Employment_Worked_12M_binned" on Crime Counts

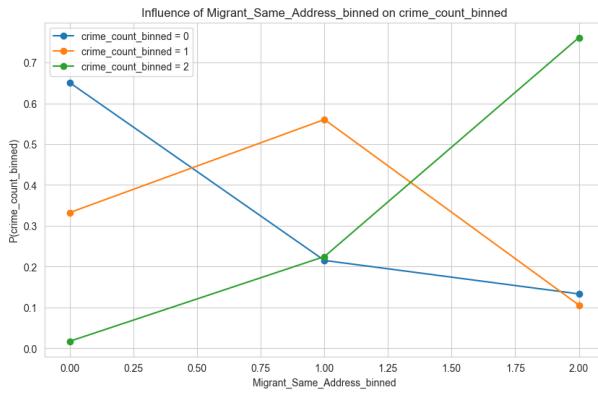


Fig. 15. Bayesian Network - Influence of "Migrant_Same_Address_binned" on Crime Counts

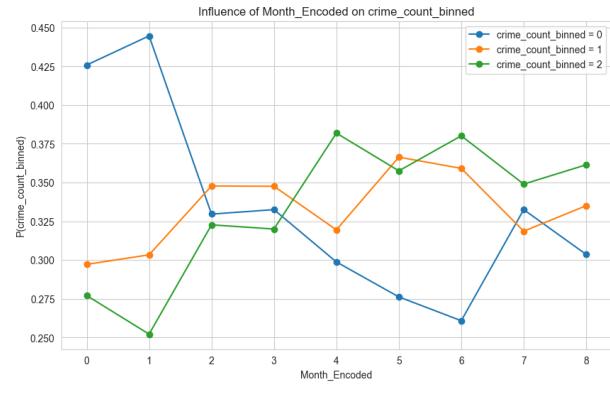


Fig. 18. Bayesian Network - Influence of "Month_Encoded" on Crime Counts

4) *Influence of "Ethnic_White_British_binned" on Crime:* It was discovered that there was a positive correlation between the percentage of Ethnic White British people living in a given area and the crime rate (Figure 16).

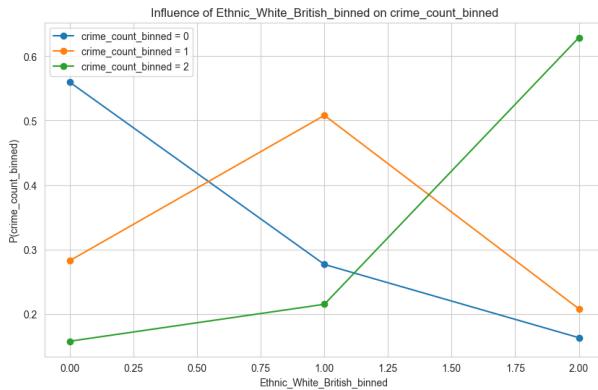


Fig. 16. Bayesian Network - Influence of "Ethnic_White_British_binned" on Crime Counts

5) *Influence of "Employment_Worked_12M_binned" on Crime:* As seen in Figure ??, higher employment levels surprisingly correlate with an increase in high crime counts. This contradicts traditional expectations that higher employment reduces crime.

6) *Influence of "Month_Encoded" on Crime:* Figure 18 reveals seasonal trends in crime, with notable fluctuations throughout the year.

7) *Discussion: Socio-Economic Factors Driving Crime:* The investigation of Bayesian networks shed light on the intricate connections between socioeconomic variables and criminal activity. Policies targeted at lessening financial hardship as a method of preventing crime are supported by the high correlation between unemployment and crime rates.

8) *Policy Implications:* When developing crime prevention methods, the Bayesian Network analysis emphasizes how crucial it is to take socioeconomic issues like unemployment, education, and stable migration into account. Given the higher likelihood of rising crime in places with high unemployment and erratic job patterns, policymakers ought to concentrate their efforts in these

areas.

F. Ethical Considerations

By coming up with trustworthy protection measures over data, adherence to regulations such as GDPR ensures the personal information of a person regarding their privacy and security. Equality can only be guaranteed through identification of the errors in data and in the models. Trust of the people can only be maintained by ensuring that predictive models are built with much accuracy and their usage is being given proper attention. Interventions should not disproportionately hurt the vulnerable populations.

V. CONCLUSION AND FUTURE WORKS

This study successfully integrated LSTM models, Explainable AI (XAI) techniques, and causal inference to analyze UK crime time series data. The LSTM model showed good performance in crime trend forecasting but faced generalization challenges during external disruptions like the COVID-19 pandemic. SHAP & LIME were employed to improve model interpretability. SHAP offering comprehensive global insights & LIME providing quick instance-specific explanations. Random Forest outperformed Linear Regression in predicting crime counts, highlighting unemployment, education, and employment instability as significant socio-economic factors influencing crime rates. Bayesian Network analysis confirmed causal relationships, identifying unemployment and education as key drivers of crime. For future work, incorporating external factors, which may responsible for crimes, such as whether data or public events or CCTV positions [41] [42], into the LSTM model could enhance robustness. Further exploration of hybrid models and advanced causal inference techniques, such as Structural Equation Models (SEM) [43], along with more granular socio-economic data, will improve prediction accuracy and provide deeper insights into crime dynamics.

REFERENCES

- [1] O. Z. Apene, N. V. Blamah, and G. I. Aimufua, "Advancements in crime prevention and detection: From traditional approaches to artificial intelligence solutions," *European Journal of Applied Science, Engineering and Technology*, vol. 2, no. 2, pp. 285–297, 2024.
- [2] Y. Zhou, F. Wang, and S. Zhou, "The spatial patterns of the crime rate in london and its socio-economic influence factors," *Social Sciences*, vol. 12, no. 6, p. 340, 2023.
- [3] E. Cesario, P. Lindia, and A. Vinci, "Multi-density crime predictor: An approach to forecast criminal activities in multi-density crime hotspots," *Journal of Big Data*, vol. 11, no. 1, 2024.
- [4] M. T. Ribeiro, S. Singh, and C. Guestrin, ""why should i trust you?"," *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2016.
- [5] J. Pearl, "Causal inference in statistics: An overview," *Statistics Surveys*, vol. 3, 2009.
- [6] A. I. Weinberg *et al.*, "Causality from bottom to top: A survey," 2024. [Online]. Available: <https://arxiv.org/html/2403.11219v1>
- [7] X. Zhang *et al.*, "Interpretable machine learning models for crime prediction," *Computers, Environment and Urban Systems*, vol. 94, p. 101789, 2022.
- [8] J. Zeng, B. Ustun, and C. Rudin, "Interpretable classification models for recidivism prediction," *Journal of the Royal Statistical Society Series A: Statistics in Society*, vol. 180, no. 3, p. 689–722, 2016.
- [9] E. Parkar *et al.*, "Comparative study of deep learning explainability and causal ai for fraud detection," *International Journal on Smart Sensing and Intelligent Systems*, vol. 17, no. 1, 2024.
- [10] D. Wang *et al.*, "Understanding the spatial distribution of crime based on its related variables using geospatial discriminative patterns," *Computers, Environment and Urban Systems*, vol. 39, p. 93–106, 2013.
- [11] Office for Statistics Regulation, "The quality of police recorded crime statistics for england and wales," 2024. [Online]. Available: <https://osr.statisticsauthority.gov.uk/publication/the-quality-of-police-recorded-crime-statistics-for-england-and-wales>
- [12] P. Linardatos, V. Papastefanopoulos, and S. Kotsiantis, "Explainable ai: A review of machine learning interpretability methods," *MDPI*, 2020. [Online]. Available: <https://www.mdpi.com/1099-4300/23/1/18>
- [13] O. Kounadi *et al.*, "A systematic review on spatial crime forecasting," *Crime Science*, vol. 9, no. 1, 2020.
- [14] J. Wu, S. M. Abrar *et al.*, "Enhancing short-term crime prediction with human mobility flows and deep learning architectures," *EPJ Data Science*, vol. 11, no. 1, 2022.
- [15] S. Chainey, L. Tompson, and S. Uhlig, "The utility of hotspot mapping for predicting spatial patterns of crime," *Security Journal*, vol. 21, no. 1-2, pp. 4–28, 2008.
- [16] H.-W. Kang and H.-B. Kang, "Prediction of crime occurrence from multi-modal data using deep learning," *PLOS ONE*, vol. 12, no. 4, 2017.
- [17] A. Casolaro *et al.*, "Deep learning for time series forecasting: Advances and open problems," *Information*, vol. 14, no. 11, p. 598, 2023.
- [18] V. Mandalapu *et al.*, "Crime prediction using machine learning and deep learning: A systematic review and future directions," *IEEE Access*, vol. 11, p. 60153–60170, 2023.
- [19] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, p. 1735–1780, 1997.
- [20] D. Paul and D. Yeasin, "Tslstm: Long short term memory (lstm) model for time series forecasting," *CRAN: Contributed Packages*, 2022.
- [21] U. M. Ramirez-Alcocer, E. Tello-Leal, and J. A. Mata-Torres, "Smart 2019: The eighth international conference on smart cities, systems, devices and technologies," 2019. [Online]. Available: https://personales.upv.es/thinkmind/SMART/SMART_2019/smart_2019_2_30_40010.html
- [22] M. Solís and L.-A. Calvo-Valverde, "Deep learning for crime forecasting of multiple regions, considering spatial-temporal correlations between regions," in *ITISE 2024*, 2024.
- [23] Y. Rayhan and T. Hashem, "Aist: An interpretable attention-based deep learning model for crime prediction," *ACM Transactions on Spatial Algorithms and Systems*, vol. 9, no. 2, p. 1–31, 2023.
- [24] A. Chaddad *et al.*, "Survey of explainable ai techniques in healthcare," *Sensors*, vol. 23, no. 2, p. 634, 2023.
- [25] C. O. Retzlaff *et al.*, "Post-hoc vs ante-hoc explanations: Xai design guidelines for data scientists," *Cognitive Systems Research*, vol. 86, p. 101243, 2024.
- [26] A. Holzinger *et al.*, "Causability and explainability of artificial intelligence in medicine," *WIREs Data Mining and Knowledge Discovery*, vol. 9, no. 4, 2019.
- [27] A. Gramegna and P. Giudici, "Shap and lime: An evaluation of discriminative power in credit risk," *Frontiers in Artificial Intelligence*, vol. 4, 2021.
- [28] M. Louhichi *et al.*, "Shapley values for explaining the black box nature of machine learning model clustering," *Procedia Computer Science*, vol. 220, pp. 806–811, 2023.

- [29] K. Baumgartner, S. Ferrari, and G. Palermo, “Constructing bayesian networks for criminal profiling from limited data,” *Knowledge-Based Systems*, vol. 21, no. 7, pp. 563–572, 2008.
- [30] T. Hu *et al.*, “Urban crime prediction based on spatio-temporal bayesian model,” *PLOS ONE*, vol. 13, no. 10, 2018.
- [31] Q. Wen *et al.*, “Robuststl: A robust seasonal-trend decomposition algorithm for long time series,” *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, p. 5409–5416, 2019.
- [32] A. Sohrabbeig, O. Ardakanian, and P. Musilek, “Decompose and conquer: Time series forecasting with multisessional trend decomposition using loess,” *Forecasting*, vol. 5, no. 4, p. 684–696, 2023.
- [33] C. Huang *et al.*, “Deepcrime,” in *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*, 2018.
- [34] O. Bespalova, “Gdp forecasts: Informational asymmetry of the spf and fomc minutes,” *International Journal of Forecasting*, vol. 36, no. 4, pp. 1531–1540, 2020.
- [35] D. R. Setiadi *et al.*, “Digital image steganography survey and investigation,” *Signal Processing*, vol. 206, p. 108908, 2023.
- [36] Q. Wen *et al.*, “Time series data augmentation for deep learning: A survey,” *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence*, p. 4653–4660, 2021.
- [37] N. Srivastava *et al.*, “Dropout: A simple way to prevent neural networks from overfitting,” *Journal of Machine Learning Research*, 1970. [Online]. Available: <http://jmlr.org/papers/v15/srivastava14a.html>
- [38] N. Jiang, K. Miao, Y. Chai, D. Lu, and J. Wu, “Spatio-temporal prediction of crime based on data mining and lstm network,” *2023 IEEE 6th Information Technology, Networking, Electronic and Automation Control Conference (ITNEC)*, vol. 6, pp. 672–676, 2023.
- [39] X. Zhang, L. Liu, L. Xiao, and J. Ji, “Comparison of machine learning algorithms for predicting crime hotspots,” *IEEE Access*, vol. 8, pp. 181 302–181 310, 2020.
- [40] T. P. Lowe *et al.*, “Violent crime and covid-19 in england and wales,” *Studies of Organized Crime*, pp. 279–294, 2022.
- [41] M. Ranson, “Crime, weather, and climate change,” *Journal of Environmental Economics and Management*, vol. 67, no. 3, pp. 274–302, 2014.
- [42] B. C. Welsh and D. P. Farrington, “Public area cctv and crime prevention: An updated systematic review and meta-analysis,” *Justice Quarterly*, vol. 26, no. 4, pp. 716–745, 2009.
- [43] A. S. Albahri *et al.*, “Hybrid artificial neural network and structural equation modeling techniques: A survey,” *Complex & Intelligent Systems*, vol. 8, no. 2, pp. 1781–1801, 2021.

APPENDIX A DATASET IMAGES

	education_lower_tier_local_authorities_code	education_lower_tier_local_authorities	education_does_not_apply	economy_lowest_tier	economy_lowest_tier
1	E04000001	Hertfordshire	17652	3749	56937
2	E04000002	Middlesex	29986	56434	57594
3	E04000003	Nottinghamshire	24189	55421	56721
4	E04000004	Nottingham-on-Tees	38833	56413	73107
5	E04000005	Nottingham	35695	47998	58818
6	E04000006	Notton	24530	58786	45242
7	E04000007	Macclesfield	38256	183191	49527
8	E04000008	Macclesfield with Dron	32332	55977	59911
9	E04000009	Blackpool	24497	57562	56977
10	E04000010	Kingston upon Hull	53137	117789	96888
11	E04000011	East Riding of Yorkshire	54075	156229	158220
12	E04000012	North East Lincolnshire	26442	48849	58475
13	E04000013	North Lincolnshire	30213	76233	63244
14	E04000014	York	30738	94195	75865
15	E04000015	Bury	53151	118371	93999
16	E04000016	Leicester	77557	154594	176088
17	E04000017	Orkney	6751	18783	15515
18	E04000018	Notland	59310	128668	135405
19	E04000019	Nottinghamshire	29986	56411	56722
20	E04000020	Telford and Wrekin	36867	85497	53177
21	E04000021	Stoke-on-Trent	52469	112581	93534
22	E04000022	Bath and North East Somerset	31791	89812	73056
23	E04000023	Wiltshire	32791	59391	112149
24	E04000024	North Somerset	38161	102792	75868
25	E04000025	South Gloucestershire	53058	147652	87754
26	E04000026	Plymouth	46071	121618	97686
27	E04000027	Torbay	22929	54955	58977
28	E04000028	Sandown	46259	117866	49285
29	E04000029	Isle of Wight	46166	100848	64495

Fig. 19. Economy status

	education_level_data_transformed.csv	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10	C11	C12	C13	C14	C15	C16	C17	C18	C19	C20	C21	C22	C23	C24	C25	C26	C27	C28	C29	C30	C31	C32	C33	C34	C35	C36	C37	C38	C39	C40	C41	C42	C43	C44	C45	C46	C47	C48	C49	C50	C51	C52	C53	C54	C55	C56	C57	C58	C59	C60	C61	C62	C63	C64	C65	C66	C67	C68	C69	C70	C71	C72	C73	C74	C75	C76	C77	C78	C79	C80	C81	C82	C83	C84	C85	C86	C87	C88	C89	C90	C91	C92	C93	C94	C95	C96	C97	C98	C99	C100	C101	C102	C103	C104	C105	C106	C107	C108	C109	C110	C111	C112	C113	C114	C115	C116	C117	C118	C119	C120	C121	C122	C123	C124	C125	C126	C127	C128	C129	C130	C131	C132	C133	C134	C135	C136	C137	C138	C139	C140	C141	C142	C143	C144	C145	C146	C147	C148	C149	C150	C151	C152	C153	C154	C155	C156	C157	C158	C159	C160	C161	C162	C163	C164	C165	C166	C167	C168	C169	C170	C171	C172	C173	C174	C175	C176	C177	C178	C179	C180	C181	C182	C183	C184	C185	C186	C187	C188	C189	C190	C191	C192	C193	C194	C195	C196	C197	C198	C199	C200	C201	C202	C203	C204	C205	C206	C207	C208	C209	C210	C211	C212	C213	C214	C215	C216	C217	C218	C219	C220	C221	C222	C223	C224	C225	C226	C227	C228	C229	C230	C231	C232	C233	C234	C235	C236	C237	C238	C239	C240	C241	C242	C243	C244	C245	C246	C247	C248	C249	C250	C251	C252	C253	C254	C255	C256	C257	C258	C259	C260	C261	C262	C263	C264	C265	C266	C267	C268	C269	C270	C271	C272	C273	C274	C275	C276	C277	C278	C279	C280	C281	C282	C283	C284	C285	C286	C287	C288	C289	C290	C291	C292	C293	C294	C295	C296	C297	C298	C299	C300	C301	C302	C303	C304	C305	C306	C307	C308	C309	C310	C311	C312	C313	C314	C315	C316	C317	C318	C319	C320	C321	C322	C323	C324	C325	C326	C327	C328	C329	C330	C331	C332	C333	C334	C335	C336	C337	C338	C339	C340	C341	C342	C343	C344	C345	C346	C347	C348	C349	C350	C351	C352	C353	C354	C355	C356	C357	C358	C359	C360	C361	C362	C363	C364	C365	C366	C367	C368	C369	C370	C371	C372	C373	C374	C375	C376	C377	C378	C379	C380	C381	C382	C383	C384	C385	C386	C387	C388	C389	C390	C391	C392	C393	C394	C395	C396	C397	C398	C399	C400	C401	C402	C403	C404	C405	C406	C407	C408	C409	C410	C411	C412	C413	C414	C415	C416	C417	C418	C419	C420	C421	C422	C423	C424	C425	C426	C427	C428	C429	C430	C431	C432	C433	C434	C435	C436	C437	C438	C439	C440	C441	C442	C443	C444	C445	C446	C447	C448	C449	C450	C451	C452	C453	C454	C455	C456	C457	C458	C459	C460	C461	C462	C463	C464	C465	C466	C467	C468	C469	C470	C471	C472	C473	C474	C475	C476	C477	C478	C479	C480	C481	C482	C483	C484	C485	C486	C487	C488	C489	C490	C491	C492	C493	C494	C495	C496	C497	C498	C499	C500	C501	C502	C503	C504	C505	C506	C507	C508	C509	C510	C511	C512	C513	C514	C515	C516	C517	C518	C519	C520	C521	C522	C523	C524	C525	C526	C527	C528	C529	C530	C531	C532	C533	C534	C535	C536	C537	C538	C539	C540	C541	C542	C543	C544	C545	C546	C547	C548	C549	C550	C551	C552	C553	C554	C555	C556	C557	C558	C559	C560	C561	C562	C563	C564	C565	C566	C567	C568	C569	C570	C571	C572	C573	C574	C575	C576	C577	C578	C579	C580	C581	C582	C583	C584	C585	C586	C587	C588	C589	C590	C591	C592	C593	C594	C595	C596	C597	C598	C599	C600	C601	C602	C603	C604	C605	C606	C607	C608	C609	C610	C611	C612	C613	C614	C615	C616	C617	C618	C619	C620	C621	C622	C623	C624	C625	C626	C627	C628	C629	C630	C631	C632	C633	C634	C635	C636	C637	C638	C639	C640	C641	C642	C643	C644	C645	C646	C647	C648	C649	C650	C651	C652	C653	C654	C655	C656	C657	C658	C659	C660	C661	C662	C663	C664	C665	C666	C667	C668	C669	C670	C671	C672	C673	C674	C675	C676	C677	C678	C679	C680	C681	C682	C683	C684	C685	C686	C687	C688	C689	C690	C691	C692	C693	C694	C695	C696	C697	C698	C699	C700	C701	C702	C703	C704	C705	C706	C707	C708	C709	C710	C711	C712	C713	C714	C715	C716	C717	C718	C719	C720	C721	C722	C723	C724	C725	C726	C727	C728	C729	C730	C731	C732	C733	C734	C735	C736	C737	C738	C739	C740	C741	C742	C743	C744	C745	C746	C747	C748	C749	C750	C751	C752	C753	C754	C755	C756	C757	C758	C759	C760	C761	C762	C763	C764	C765	C766	C767	C768	C769	C770	C771	C772	C773	C774	C775	C776	C777	C778	C779	C770	C771	C772	C773	C774	C775	C776	C777	C778	C779	C780	C781	C782	C783	C784	C785	C786	C787	C788	C789	C780	C781	C782	C783	C784	C785	C786	C787	C788	C789	C790	C791	C792	C793	C794	C795	C796	C797	C798	C799	C790	C791	C792	C793	C794	C795	C796	C797	C798	C799	C800	C801	C802	C803	C804	C805	C806	C807	C808	C809	C800	C801	C802	C803	C804	C805	C806	C807	C808	C809	C810	C811	C812	C813	C814	C815	C816	C817	C818	C819	C810	C811	C812	C813	C814	C815	C816	C817	C818	C819	C820	C821	C822	C823	C824	C825	C826	C827	C828	C829	C820	C821	C822	C823	C824	C825	C826	C827	C828	C829	C830	C831	C832	C833	C834	C835	C836	C837	C838	C839	C830	C831	C832	C833	C834	C835	C836	C837	C838	C839	C840	C841	C842	C843	C844	C845	C846	C847	C848	C849	C840	C841	C842	C843	C844	C845	C846	C847	C848	C849	C850	C851	C852	C853	C854	C855	C856	C857	C858	C859	C850	C851	C852	C853	C854	C855	C856	C857	C858	C859	C860	C861	C862	C863	C864	C865	C866	C867	C868	C869	C860	C861	C862	C863	C864	C865	C866	C867	C868	C869	C870	C871	C872	C873	C874	C875	C876	C877	C878	C879	C870	C871	C872	C873	C874	C875	C876	C877	C878	C879	C880	C881	C882	C883	C884	C885	C886	C887	C888	C889	C880	C881	C882	C883	C884	C885	C886	C887	C888	C8

	C1	C2	C3	C4
1	ethnic_lower_tier_local_authorities_code	ethnic_lower_tier_local_authorities	ethnic_asian_asian_british_or_asian_welsh_banglades...	ethnic_asian_asian_british_or_asian_welsh_banglades...
2	E04000001	Hartlepool	278	237
3	E04000002	Middlesbrough	595	669
4	E04000003	Milton Keynes	120	108
5	E04000004	Stockton-on-Tees	236	690
6	E04000005	Darlington	759	308
7	E04000006	Wolverhampton	58	418
8	E04000007	Warrington	127	124
9	E04000008	Blackburn with Darwen	1473	400
10	E04000009	Bury	442	542
11	E04000010	Kingston upon Hull	279	252
12	E04000011	East Riding of Yorkshire	179	766
13	E04000012	North East Lincolnshire	318	438
14	E04000013	North Lincolnshire	1348	451
15	E04000014	Tyne and Wear	423	1099
16	E04000015	Derby	877	1456
17	E04000016	Liverpool	7065	2481
18	E04000017	Bolton	19	149
19	E04000018	Newcastle upon Tyne	2232	1045
20	E04000019	Herefordshire	55	339
21	E04000020	Telford and Wrekin	207	818
22	E04000021	Stoke-on-Trent	1177	778
23	E04000022	Redcar and Cleveland	359	289
24	E04000023	Bristol	2616	5466
25	E04000024	North Somerset	435	725
26	E04000025	South Gloucestershire	100	2113

Access is granted from Event Dispatch Thread (EDT) only see <https://qa.qspjplatform-threading.readthedocs.io/en/latest/>

Current thread: ThreadApplicationImpl posted thread 2795.main@471684762 (EventQueue\$1DispatchThread!=false)

SystemEventQueueThread:ThreadImpl\$EventQueue\$1DispatchThread@471684762

Fig. 22. Ethnic Data

	C1	C2	C3	C4
1	Crime ID	Month	Reported by	Falls within
2				Longitude
3				Latitude
4				Location
5				LSDA code
6				LSDA name
7				Crime
8				
9				
10				
11				
12				
13				
14				
15				
16				
17				
18				
19				
20				
21				
22				
23				
24				
25				
26				
27				
28				
29				
30				
31				
32				
33				
34				
35				
36				
37				
38				
39				
40				
41				
42				
43				
44				
45				
46				

Fig. 23. UK Crime Data

APPENDIX B ANNOTATED PAPERS

You can access the folder containing the papers used in the study at the following link:

[Click here to view the papers](#)

APPENDIX C LSTM MODEL CODE

```

1 #%% md
2 # Crime Data Analysis and Forecasting using
3 # LSTM and Explainability Techniques
4
5 This notebook demonstrates the process of
6 analyzing crime data, forecasting future
7 trends using an LSTM model, and
8 interpreting the model's predictions using
9 LIME and SHAP.
10 #%% md
11 # 1. Imports
12
13 Let's begin by collecting all the import
14 statements together at the beginning. This
15 ensures that all dependencies are loaded
16 before any logic is executed.
17
18 #%%
19 import numpy as np
20 import pandas as pd
21 import matplotlib.pyplot as plt
22 from sklearn.preprocessing import MinMaxScaler
23 from sklearn.metrics import mean_absolute_error
24 , mean_squared_error
25 from tensorflow.keras.models import Sequential
26 from tensorflow.keras.layers import Dense, LSTM
27 , Dropout
28 from statsmodels.tsa.seasonal import
29 seasonal_decompose
30 from lime import lime_tabular
31 import shap
32 import os # Added to handle file paths more
33 robustly
34
35 #%% md
36 # 2. Loading and Preparing the Dataset
37
38 Now, I'll load the dataset and perform some
39 initial data cleaning and transformation:
40
41 #%%
42 # Define the file path more robustly using os.
43 path
44 file_path = os.path.join('/Users', 'arkamandol',
45 , 'DataspellProjects', ,
46 'Desertation_arka_23023023', 'data_files', ,
47 'UK_Police_Street_Crime_2018-10-01
48 -to_2021_09_31.csv')
49
50 # Load the data from the specified file path
51 try:
52     df = pd.read_csv(file_path)
53 except FileNotFoundError:
54     print(f"File not found: {file_path}")
55     raise
56 except Exception as e:
57     print(f"An error occurred while loading the
58 file: {e}")
59     raise
60
61 # Display the initial rows of the dataset to
62 understand its structure
63 df.head()
64
65 # 1. Remove any rows where 'Longitude' or
66 'Latitude' values are missing
67 df_cleaned = df.dropna(subset=['Longitude', ,
68 'Latitude'])

```

```

47
48 # 2. Focus the dataset on "Violence and sexual
49 # offences" crimes
50 df_cleaned = df_cleaned[df_cleaned['Crime type',
51 ] == "Violence and sexual offences"]
52
53 # 3. Keep only the necessary columns for
54 # analysis
55 columns_to_keep = ['Month', 'Longitude', 'Latitude']
56 df_cleaned = df_cleaned[columns_to_keep]
57
58 #%% md
59 **Explanation**:
60 - **File Path Handling**: Improved the file
   path handling with `os.path.join()` for
   better cross-platform compatibility.
61 - **Error Handling**: Added exception handling
   during file loading to catch and report
   issues like file not found or other IO
   errors.
62 - **Data Cleaning**: The data is filtered to
   remove rows with missing geographical
   information and focus on specific crime
   types.
63
64 #%% md
65 # 3. Time Series Transformation
66
67 Here I transform the data into a time series
   format suitable for analysis.
68
69 # Convert the 'Month' column to datetime format
   for time series analysis
70 data = df_cleaned.copy()
71 data['Month'] = pd.to_datetime(data['Month'])
72
73 # Set the 'Month' column as the index
74 data.set_index('Month', inplace=True)
75
76 # Aggregate the data to get the number of
   crimes per month
77 monthly_crime_count = data.resample('M').size()
78
79 # Display the first few entries of the
   aggregated data
80 monthly_crime_count.head()
81
82 #%% md
83 **Explanation**:
84 - **Datetime Conversion**: Converts the 'Month'
   column into a datetime format, which is
   essential for time series operations.
85 - **Indexing**: The 'Month' column is set as
   the index to facilitate resampling and time
   series analysis.
86 - **Aggregation**: The number of crimes per
   month is calculated.
87
88 #%% md
89 # 4. Visualization of Time Series Data
90
91 Next, visualize the monthly crime counts over
   time.
92
93 #%%
94 # Generate a line plot to visualize monthly
   crime counts over time
95 plt.figure(figsize=(14, 7))
96 monthly_crime_count.plot()
97 plt.title('Monthly Crime Counts')
98 plt.xlabel('Month')
99 plt.ylabel('Number of Crimes')
100 plt.grid(True)
101 plt.show()
102
103 #%% md
104 **Explanation**:
105 - **Visualization**: A line plot is generated
   to visualize the crime trend over time,
   which helps in understanding the data
   distribution.
106
107 #%% md
108 # 5. Seasonal Decomposition of Time Series
109
110 Decompose the time series to analyze the
   underlying components:
111
112 #%%
113 # Decompose the time series into trend,
   seasonal, and residual components
114 decomposition = seasonal_decompose(
   monthly_crime_count, model='additive')
115
116 # Plot the decomposed components to better
   understand the underlying patterns
117 fig = decomposition.plot()
118 fig.set_size_inches(14, 7)
119 plt.show()
120
121 #%% md
122 **Explanation**:
123 - **Decomposition**: The time series is
   decomposed into trend, seasonal, and
   residual components to identify underlying
   patterns in the data.
124
125 #%% md
126 # 6. LSTM Model for Time Series Forecasting
127
128 Now need to build and train an LSTM model to
   predict future crime counts:
129
130 #%%
131 # Convert the series to a supervised learning
   dataset
132 def create_dataset(data, time_step=1):
133     X, Y = [], []
134     for i in range(len(data) - time_step):
135         X.append(data[i:(i + time_step), 0])
136         Y.append(data[i + time_step, 0])
137     return np.array(X), np.array(Y)
138
139 # Reshape the data into a format suitable for
   LSTM
140 scaler = MinMaxScaler(feature_range=(0, 1))
141 scaled_data = scaler.fit_transform(
   monthly_crime_count.values.reshape(-1, 1))
142
143 # Create sequences of 12 months (1 year)
144 time_step = 12
145 X, Y = create_dataset(scaled_data, time_step)
146
147 # Reshape input to be [samples, time steps,
   features]
148 X = X.reshape(X.shape[0], X.shape[1], 1)
149
150 # Split data into training and testing sets

```

```

151 train_size = int(len(X) * 0.7)
152 test_size = len(X) - train_size
153 X_train, X_test = X[:train_size], X[train_size:]
154 Y_train, Y_test = Y[:train_size], Y[train_size:]
155
156 #%% md
157 **Explanation**:
158 - **Supervised Dataset Creation**: The 'create_dataset' function converts the time series data into sequences for supervised learning.
159 - **Normalization**: The data is scaled between 0 and 1 using MinMaxScaler, which helps in stabilizing the LSTM model training.
160 - **Train-Test Split**: The dataset is split into training and testing sets, maintaining 70% of the data for training.
161
162 #%% md
163 # 7. Building the LSTM Model
164
165 Let's define, compile, and train the LSTM model :
166
167 #%%
168 # Define the LSTM model with dropout
169 model = Sequential()
170 model.add(LSTM(50, return_sequences=True,
171               input_shape=(time_step, 1)))
171 model.add(Dropout(0.3)) # Adding dropout after the first LSTM layer
172 model.add(LSTM(50, return_sequences=False))
173 model.add(Dropout(0.3)) # Adding dropout after the second LSTM layer
174 model.add(Dense(25))
175 model.add(Dense(1))
176
177 # Compile the model
178 model.compile(optimizer='adam', loss='mean_squared_error')
179
180 # Train the model
181 model.fit(X_train, Y_train, batch_size=1,
182             epochs=500)
183
184 #%% md
185 **Explanation**:
186 - **Model Architecture**: A Sequential LSTM model is defined with two LSTM layers, each followed by a Dropout layer to prevent overfitting. The model ends with Dense layers to produce the final output.
187 - **Compilation**: The model is compiled with the Adam optimizer and mean squared error loss function.
188 - **Training**: The model is trained with a small batch size and a large number of epochs for better learning.
189
190 #%% md
191 # 8. Making Predictions and Evaluating the Model
192 After training, predict and evaluate the model:
193
194 #%%
195 # Predict using the trained model
196 train_predict = model.predict(X_train)

197 test_predict = model.predict(X_test)
198
199 # Inverse transform the predictions to original scale
200 train_predict = scaler.inverse_transform(
201   train_predict)
202 test_predict = scaler.inverse_transform(
203   test_predict)
204 Y_train = scaler.inverse_transform([Y_train])
205 Y_test = scaler.inverse_transform([Y_test])
206
207 # Calculate MSE, RMSE and MAE for training and testing data
208 train_mse = mean_squared_error(Y_train.T,
209   train_predict)
210 train_rmse = np.sqrt(train_mse)
211 train_mae = mean_absolute_error(Y_train.T,
212   train_predict)
213 test_mse = mean_squared_error(Y_test.T,
214   test_predict)
215 test_rmse = np.sqrt(test_mse)
216 test_mae = mean_absolute_error(Y_test.T,
217   test_predict)
218
219 print(f'Train MSE: {train_mse}')
220 print(f'Train RMSE: {train_rmse}')
221 print(f'Train MAE: {train_mae}')
222 print(f'Test MSE: {test_mse}')
223 print(f'Test RMSE: {test_rmse}')
224 print(f'Test MAE: {test_mae}')
225
226 #%% md
227 **Explanation**:
228 - **Predictions**: The model is used to predict both the training and testing data.
229 - **Inverse Scaling**: The predictions are transformed back to the original scale.
230 - **Evaluation**: Mean Squared Error (MSE), Root Mean Squared Error (RMSE), and Mean Absolute Error (MAE) are calculated to evaluate model performance.
231
232 #%% md
233 # 9. Future Predictions
234
235 Predicting crime counts for the next 6 months:
236
237 # Predict the next 6 months (October 2021 to March 2022)
238 future_steps = 6
239 future_predictions = []
240 last_sequence = X[-1] # The last sequence of the test set
241
242 for _ in range(future_steps):
243   next_pred = model.predict(last_sequence.
244     reshape(1, time_step, 1))
245   future_predictions.append(next_pred[0, 0])
246   last_sequence = np.append(last_sequence[1:], next_pred, axis=0)
247
248 # Inverse transform future predictions
249 future_predictions = scaler.inverse_transform(
250   np.array(future_predictions).reshape(-1, 1))
251
252 # Generate date range for predictions
253 last_date = monthly_crime_count.index[-1]
254 future_dates = pd.date_range(last_date + pd.
255   DateOffset(months=1), periods=future_steps,

```

```

248 freq='M')
249 #%% md
250 **Explanation**:
251 - **Future Predictions**: The model is used to
predict the next 6 months of crime counts.
252 - **Date Range Generation**: A corresponding
date range is generated for the predicted
values.
253 #%% md
254 # 10. Visualization of Predictions
255
256 Plotting actual vs predicted crime counts:
257
258 #%%
259 # Plot the results
260 plt.figure(figsize=(14, 7))
261 plt.plot(monthly_crime_count.index,
monthly_crime_count.values, label='Original
Data', color='blue')
262 plt.plot(monthly_crime_count.index[time_step:
time_step + len(train_predict)],
train_predict,
label='Training Predictions', color=
orange')
263 plt.plot(monthly_crime_count.index[time_step +
len(train_predict):time_step + len(
train_predict) + len(test_predict)],
test_predict, label='Testing
Predictions', color='green')
264 plt.plot(future_dates, future_predictions,
label='Future Predictions (Oct 2021 – Mar
2022)', color='red')
265
266 plt.title('LSTM Model Predictions vs. Actual
Data')
267 plt.xlabel('Time')
268 plt.ylabel('Crime Count')
269 plt.xticks(rotation=45)
270 plt.grid(True, which='both', axis='x',
linestyle='--', linewidth=0.7)
271 plt.legend()
272
273 # Set the x-ticks to be all months in the range
274 plt.gca().set_xticks(pd.date_range(start=
monthly_crime_count.index[0], end=
future_dates[-1], freq='M'))
275
276 # Format the x-ticks to show month and year
277 plt.gca().set_xticklabels([date.strftime('%Y-%m'
) for date in
pd.date_range(start=
monthly_crime_count.index[0], end=
future_dates[-1], freq='M')])
278
279 plt.show()
280 #%% md
281 **Explanation**:
282 - **Visualization**: The plot compares the
original data with both the training/
testing predictions and future predictions,
showing how well the model fits and
forecasts.
283 #%% md
284 # 11. LIME for Model Interpretability
285
286 Using LIME to explain the predictions of the
LSTM model:
287
288
289
290
291
292
293 #%%
294 # Define the prediction function for LIME
295 def lstm_predict(input_data):
296     """ Wrapper function to adapt input data
for the LSTM model's expected input shape
"""
297     reshaped_data = input_data.reshape(
input_data.shape[0], time_step, 1) #
Reshape to [samples, time steps, features]
298     return model.predict(reshaped_data)
299
300 # Reshape X_train for LIME
301 X_train_2d = X_train.reshape(X_train.shape[0],
X_train.shape[1])
302
303 # Initialize the LIME Explainer
304 explainer = lime_tabular.LimeTabularExplainer(
305     training_data=X_train_2d,
306     feature_names=['t-' + str(i) for i in range(
307         X_train.shape[1])], # Create feature
308     names for each time step
309     mode='regression')
310
311 # Loop through each instance in X_test
312 lime_explanations = []
313
314 for i in range(X_test.shape[0]):
315     # Reshape test instance to 2D for LIME
316     test_instance = X_test[i].reshape(1, -1)[0]
317
318     # Generate explanation for the current
319     # instance
320     exp = explainer.explain_instance(
321         data_row=test_instance,
322         predict_fn=lstm_predict,
323         num_features=12 # Set this to ensure
324         # all features are considered
325     )
326
327     lime_explanations.append(exp)
328
329     # Optionally, print or plot each
330     # explanation
331     print(f"Explanation for instance {i}:")
332     exp.show_in_notebook(show_table=True,
333     show_all=False)
334
335     # Plotting each explanation
336     features, weights = zip(*[(i[0], i[1]) for
337         i in exp.local_exp[0]])
338     features = [f"t-{X_train.shape[1] - idx}" for
339         idx in features]
340
341     plt.figure(figsize=(10, 5))
342     plt.bar(features, weights, color='skyblue')
343     plt.xlabel('Features (Time Steps)')
344     plt.ylabel('Weights')
345     plt.title(f'LIME Feature Contributions for
346     Instance {i}')
347     plt.xticks(rotation=45)
348     plt.tight_layout()
349     plt.show()
350
351 #%% md
352 **Explanation**:
353 - **LIME Explanation**: The code uses LIME to
interpret the LSTM models predictions,
identifying the most influential features
for each prediction.

```

```

347 #%% md
348 # 12. SHAP for Model Interpretability
349
350 Lastly, SHAP is used to interpret the model:
351
352 #%%
353 # Define a custom prediction function that
354 # reshapes the input before predicting
355 def predict_reshaped(data):
356     reshaped_data = data.reshape(data.shape[0],
357                                   -1, 1)
358     return model.predict(reshaped_data)
359
360 # Use a small subset of training data as
361 # background for KernelExplainer
362 background = X_train.reshape(X_train.shape[0],
363                             -1) # Flatten to 2D for KernelExplainer
364
365 # Initialize KernelExplainer with the custom
366 # predict function
367 explainer = shap.KernelExplainer(
368     predict_reshaped, background)
369
370 # Compute SHAP values for test set (reshaped to
371 # 2D)
372 X_test_2d = X_test.reshape(X_test.shape[0], -1)
373 # Flatten to 2D
374 shap_values = explainer.shap_values(X_test_2d,
375                                     nsamples=100)
376
377 # Reshape SHAP values from (8, 12, 1) to (8,
378 # 12)
379 shap_values_reshaped = np.array(shap_values).reshape(X_test.shape[0], X_test.shape[1])
380
381 # Now use SHAP's summary plot, passing in the
382 # original X_test without the last dimension
383 X_test_reshaped = X_test.reshape(X_test.shape[0], X_test.shape[1]) # Remove the last
384 dimension
385 shap.summary_plot(shap_values_reshaped,
386                   X_test_reshaped)

```

Listing 1. LSTM Model Implementation

APPENDIX D BAYESIAN NETWORK CODE

```

1 #%% md
2
3 #%% md
4 # Crime Data Analysis and Bayesian Network
4 # Modeling
5
6 #%% md
7 ## 1. Imports and Setup
8
9 #%% md
10 ### Imports and Setup
11
12 - **pandas**: For data manipulation and
12 analysis.
13 - **numpy**: For numerical operations.
14 - **matplotlib and seaborn**: For visualization
14 .
15 - **sklearn**: For machine learning models and
15 preprocessing.
16 - **pgmpy**: For building and analyzing
16 Bayesian networks.
17 - **shap and lime**: For model explainability (
17 although not used explicitly in the code).
18
19 will add the imports at the beginning of the
19 script to ensure all dependencies are
19 loaded before executing any logic.
20
21 #%%
22 # Import necessary libraries
23 import pandas as pd
24 import numpy as np
25 import matplotlib.pyplot as plt
26 import seaborn as sns
27 from pgmpy.estimators import HillClimbSearch,
27 BicScore, BayesianEstimator, K2Score
28 from pgmpy.models import BayesianNetwork
29 from pgmpy.inference import VariableElimination
30 import networkx as nx
31
32 #%%
33 ## 2. Loading and Preprocessing Crime Data
34 #%% md
35 ### Loading and Preprocessing Crime Data
36
37 - **Loading Crime Data**: load the dataset
37 containing crime data from a specified CSV
37 file.
38 - **Date Conversion**: The 'Month' column is
38 converted to datetime format for easier
38 filtering and manipulation.
39 - **Filtering Data**: The data is filtered
39 to retain only records from the year 2021
39 and those related to the "Violence and
39 sexual offences" crime type.
40
41 #%%
42 # Load the dataset containing crime data
43 crime_data_path = "/Users/arkamadol/
43 DataspellProjects/Desertation_arka_23023023
43 /data_files/UK_Police_Street_Crime_2018
43 -10-01_to_2021_09_31.csv"
44 crime_df = pd.read_csv(crime_data_path)
45
46 # Convert 'Month' to datetime and filter for
46 2021 and specific crime type
47 crime_df['Month'] = pd.to_datetime(crime_df['
47 Month'])

```

```

48 filtered_crime_df = crime_df[(crime_df['Month']
49     ].dt.year == 2021) &
50         (crime_df['Crime
51             type'] == "Violence and sexual offences")]
52 ## 3. Data Cleaning and Aggregation
53
54 ## Data Cleaning and Aggregation
55
56 - **Modifying 'LSOA name'**: modify the 'LSOA
57     name' column to remove the last 5
58     characters , which might represent redundant
59     information .
60
61 # Modify 'LSOA name' to remove the last 5
62 # characters
63 filtered_crime_df['LSOA name'] =
64     filtered_crime_df['LSOA name'].str[:-5]
65 # Rename 'LSOA name' to '
66 #     lower_tier_local_authorities '
67 filtered_crime_df.rename(columns={'LSOA name':
68     'lower_tier_local_authorities'}, inplace=
69     True)
70
71 # Aggregate the data by 'Month' and '
72 #     lower_tier_local_authorities' and count
73 # occurrences
74 aggregated_crime_data = filtered_crime_df.
75     groupby(['Month',
76             'lower_tier_local_authorities']).size().
77     reset_index(name='crime_count')
78
79 ## 4. Loading and Merging Socioeconomic Data
80
81 ## Loading and Merging Socioeconomic Data
82
83 - **Loading Socioeconomic Datasets**: load
84     multiple CSV files containing different
85     socioeconomic data such as ethnic groups ,
86     economy status , employment history ,
87     education level , and migrant data .
88
89 - **Renaming Columns**: In each dataset , the
90     locality code columns are renamed to '
91         area_code' for consistency .
92
93 - **Merging Datasets**: All the
94     socioeconomic datasets are merged into a
95     single DataFrame based on the 'area_code',
96     column .
97
98 ## Load datasets
99 # Load datasets
100 ethnic_data = pd.read_csv("/Users/arkamadol/
101     /DataspellProjects/Desertation_arka_23023023
102     /data_files/Ethnic_Group_Data_Transformed.
103     csv")
104
105 economy_status_data = pd.read_csv("/Users/
106     arkamadol/DataspellProjects/
107
108     Desertation_arka_23023023 / data_files /
109     Economy_Status_Data_Transformed.csv")
110
111 employment_history_data = pd.read_csv("/Users/
112     arkamadol/DataspellProjects/
113     Desertation_arka_23023023 / data_files /
114     Employment_History_Data_Transformed.csv")
115
116 education_level_data = pd.read_csv("/Users/
117     arkamadol/DataspellProjects/
118     Desertation_arka_23023023 / data_files /
119     Education_Level_Data_Transformed.csv")
120
121 migrant_data = pd.read_csv("/Users/arkamadol/
122     /DataspellProjects/Desertation_arka_23023023
123     /data_files/Migrant_Data_Transformed.csv")
124
125 # Rename the locality code columns to '
126 #     area_code'
127 ethnic_data.rename(columns={''
128         ethnic_lower_tier_local_authorities_code':
129             'area_code'}, inplace=True)
130
131 economy_status_data.rename(columns={''
132         economy_lower_tier_local_authorities_code':
133             'area_code'}, inplace=True)
134
135 employment_history_data.rename(columns={''
136         employment_lower_tier_local_authorities_code':
137             'area_code'}, inplace=True)
138
139 education_level_data.rename(columns={''
140         education_lower_tier_local_authorities_code':
141             'area_code'}, inplace=True)
142
143 migrant_data.rename(columns={''
144         migrant_lower_tier_local_authorities_code':
145             'area_code'}, inplace=True)
146
147 # Merge all datasets on the 'area_code' column
148 merged_data = ethnic_data.merge(
149     economy_status_data, on='area_code', how='
150         outer')
151
152 merged_data = merged_data.merge(
153     employment_history_data, on='area_code', how='
154         outer')
155
156 merged_data = merged_data.merge(education_level_data,
157     on='area_code', how='outer')
158
159 merged_data = merged_data.merge(migrant_data,
160     on='area_code', how='outer')
161
162 ## 5. Cleaning and Merging Final Dataset
163
164 ## Cleaning and Merging Final Dataset
165
166 - **Dropping Redundant Columns**: identify and
167     drop any redundant ,
168     'lower_tier_local_authorities' , columns from
169     the merged dataset .
170
171 - **Final Merging**: The cleaned socioeconomic
172     data is merged with the aggregated crime
173     data using the ,
174     'lower_tier_local_authorities' , column .
175
176 - **Initial Data Inspection**: inspect the
177     first few rows and the structure of the
178     final merged dataset to ensure correctness .
179
180 ##%
181 # Load the additional dataset
182 additional_df = merged_data.copy()
183
184 # Identify and drop redundant ,
185 #     'lower_tier_local_authorities' , columns
186 columns_to_drop = [col for col in additional_df
187     .columns if 'lower_tier_local_authorities' ]

```

```

118     in col]
119 columns_to_drop.remove('
120     ethnic_lower_tier_local_authorities') # 157
121     Keep one relevant column
122
123 # Drop the identified columns
124 df_cleaned = additional_df.drop(columns=
125     columns_to_drop) 158
126
127 # Rename the relevant column to '
128     lower_tier_local_authorities'
129 df_cleaned.rename(columns={' 160
130     ethnic_lower_tier_local_authorities': ' 161
131     lower_tier_local_authorities'}, inplace=
132     True) 162
133
134 # Merge the crime data with the cleaned dataset
135     on 'lower_tier_local_authorities' 163
136 merged_data = aggregated_crime_data.merge(
137     df_cleaned, on=' 164
138     lower_tier_local_authorities', how='inner')
139
140 # Display basic information and initial rows of
141     the merged dataset 165
142 print(merged_data.head())
143 print(merged_data.info())
144
145 #%% md
146 ## 6. Data Preparation for Analysis
147
148 #%% md
149 ### Data Preparation for Analysis
150
151 - **Dropping Unnecessary Columns**: drop
152     columns that are not needed for further
153     analysis, like '
154     lower_tier_local_authorities' and '
155     area_code'.
156
157 - **Renaming Columns**: rename certain columns
158     to shorter, more intuitive names for easier
159     reference. 174
160
161 #%
162 # Create a copy for processing
163 data = merged_data.copy() 175
164
165 # Drop columns not needed for analysis
166 data.drop(columns=[',
167     'lower_tier_local_authorities', 'area_code'],
168     inplace=True) 176
169
170 # Mapping new column names for easier reference
171 new_column_names = { 177
172     'ethnic_asian_asian_british_or_asian_welsh_bang':
173         'Ethnic_Asian_Bangladeshi', 177
174     'ethnic_asian_asian_british_or_asian_welsh_chin':
175         'Ethnic_Asian_Chinese', 178
176     'ethnic_asian_asian_british_or_asian_welsh_indi':
177         'Ethnic_Asian_Indian', 179
178     'ethnic_asian_asian_british_or_asian_welsh_othe':
179         'Ethnic_Asian_Other', 180
180     'ethnic_asian_asian_british_or_asian_welsh_paki':
181         'Ethnic_Asian_Pakistani', 181
181     'ethnic_black_black_british_black_welsh_caribbe':
182         'Ethnic_Black_African', 182
183
184     'ethnic_black_black_british_black_welsh_caribbean_or_african':
185         'Ethnic_Black_Caribbean', 185
186     'ethnic_black_black_british_black_welsh_caribbean_or_african':
187         'Ethnic_Black_Other', 187
188     'ethnic_mixed_or_multiple_ethnic_groups_other_mixed_or_mult':
189         'Ethnic_Mixed_Other', 189
190
191     'ethnic_mixed_or_multiple_ethnic_groups_white_and_asian':
192         'Ethnic_Mixed_White_Asian', 192
193
194     'ethnic_mixed_or_multiple_ethnic_groups_white_and_black_afri':
195         'Ethnic_Mixed_White_Black_African', 195
196
197     'ethnic_mixed_or_multiple_ethnic_groups_white_and_black_cari':
198         'Ethnic_Mixed_White_Black_Caribbean', 198
199
200     'ethnic_other_ethnic_group_any_other_ethnic_group':
201         'Ethnic_Other', 201
202     'ethnic_other_ethnic_group_arab':
203         'Ethnic_Arab', 203
204
205     'ethnic_white_english_welsh_scottish_northern_irish_or_briti':
206         'Ethnic_White_British', 206
207     'ethnic_white_gypsy_or_irish_traveller':
208         'Ethnic_White_Gypsy_Traveller', 208
209     'ethnic_white_irish':
210         'Ethnic_White_Irish', 210
211     'ethnic_white_other_white':
212         'Ethnic_White_Other', 212
213     'ethnic_white_roma':
214         'Ethnic_White_Roma', 214
215     'economy-employed':
216         'Economy_Employed', 216
217     'economy-not-employed':
218         'Economy_Not_Employed', 218
219     'employment-not-in-employment-never-worked':
220         'Employment_Never_Worked', 220
221
222     'employment-not-in-employment-not-worked-in-the-last-12-months':
223         'Employment_Not_Worked_12M', 223
224
225     'employment-not-in-employment-worked-in-the-last-12-months':
226         'Employment_Worked_12M', 226
227
228     'education_level_1_and_entry_level_qualifications_1_to_4_gcses':
229         'Education_Level_1', 229
230
231     'education_level_2_qualifications_5_or_more_gcses_a':
232         'Education_Level_2', 232
233
234     'education_level_3_qualifications_2_or_more_a_levels_or_vces':
235         'Education_Level_3', 235
236
237     'advanced_gnvq_city_and_guilds_advanced_craft_onc_ond_btce':
238         'Education_Level_4+', 238
239     'education_no_qualifications':
240         'Education_No_Qualifications', 240
241
242     'education_other_apprenticeships_vocational_or_work':
243
244     'related_qualifications_other_qualifications_achieved_in_englis':
245         'Education_Other_Qualifications', 245
246
247     'migrant_address_one_year_ago_is_student_term-time_or_boarding_school_address_in_the_uk':
248         'Education_Other_Qualifications', 248

```

```

182     : 'Migrant_Student_UK',
183     migrant_address_one_year_ago_is_the_same_as_th
184     ': 'Migrant_Same_Address',
185     'migrant_does_not_apply': 'Migrant_NA',
186     migrant_migrant_from_outside_the_uk_address_on
187     ': 'Migrant_Outside_UK',
188     migrant_migrant_from_within_the_uk_address_one_for i, var in enumerate(variables, 1):
189     ': 'Migrant_Within_UK' 231 # Plotting the distributions
190   }
191   # Rename columns according to the mapping 232 plt.figure(figsize=(15, 10))
192   data.rename(columns=new_column_names, inplace= 233 for i, var in enumerate(variables, 1):
193       True) 234 plt.subplot(3, 2, i)
194   ## md 235 sns.histplot(data[var], kde=True)
195   ## 7. Correlation Analysis and Visualization 236 plt.title(f'Distribution of {var}')
196   ## md 237 plt.xlabel(var)
197   ### Correlation Analysis and Visualization 238 plt.ylabel('Frequency')
198   - **Correlation Matrix Calculation**: calculate 239
199     the correlation matrix for the dataset to 240
200     identify relationships between variables. 241 plt.tight_layout()
201   - **Heatmap Visualization**: A heatmap is 242 plt.show()
202     plotted to visualize the correlation matrix 243
203     , making it easier to spot strong 244 ## Logarithmic Transformation
204     correlations. 245
205   ## md 246
206   # Calculate the correlation matrix 247 ## Logarithmic Transformation
207   correlation_matrix = data.corr() 248
208   # Set the size of the plot 249 - **Log Transformation**: apply a logarithmic
209   plt.figure(figsize=(10, 8)) # Increase the 250 transformation to the selected variables to
210   figure size for better readability 251
211   # Create the heatmap with annotations 252 ## Applying logarithmic transformation, adding 1
212   heatmap = sns.heatmap(correlation_matrix, annot= 253
213     =True, fmt=".1f", annot_kws={'size': 6}, 254 log_transformed_data = data.copy()
214     cmap='coolwarm', linewidths=.5) 255 log_transformed_data['crime_count_log'] = np.
215   # Rotate the x and y labels for better 256 log_transformed_data['Education_Level_3_log'] =
216   readability 257 np.log1p(log_transformed_data['Education_Level_3'])
217   plt.xticks(rotation=90) # Rotate the x-axis 258 log_transformed_data['Migrant_Same_Address_log'] =
218   labels for better readability 259 np.log1p(log_transformed_data['Migrant_Same_Address'])
219   plt.yticks(rotation=0) # y-axis labels can be 260 log_transformed_data['Employment_Worked_12M_log'] =
220   left horizontal or adjusted as needed 261 np.log1p(log_transformed_data['Employment_Worked_12M'])
221   # Show the plot 262 log_transformed_data['Economy_Not_Employed_log'] =
222   plt.show() 263 np.log1p(log_transformed_data['Economy_Not_Employed'])
223   ## md 264 log_transformed_data['Ethnic_White_British_log'] =
224   ## 8. Distribution Analysis 265 np.log1p(log_transformed_data['Ethnic_White_British'])
225   ## md 266
226   ### Distribution Analysis 267 # Plotting the log-transformed distributions
227   - **Variable Selection**: select key variables 268 plt.figure(figsize=(15, 10))
228     related to crime, education, migration, 269 for i, var in enumerate(['crime_count_log', ,
229     employment, economy, and ethnicity for 270   'Education_Level_3_log', ,
230     further analysis. 271   'Migrant_Same_Address_log', ,
231     - **Distribution Plots**: plot the 272   'Employment_Worked_12M_log', ,
232     distribution of these variables using 273   'Economy_Not_Employed_log', ,
233     histograms with kernel density estimation ( 274   'Ethnic_White_British_log'], 1):
234     KDE) to understand their underlying 275   plt.subplot(3, 2, i)
235     distribution. 276   sns.histplot(log_transformed_data[var], kde
236   ## md 277   =True)
238   ## Log-Transformed Distribution Plots**: plot 278   plt.title(f'Log-Transformed Distribution of
239     the distributions of the log-transformed 279   var)
```

```

269     {var}'))  
270     plt.xlabel(var)  
271     plt.ylabel('Frequency')  
272 plt.tight_layout()  
273 plt.show()  
274  
275 #%% md  
276 ## 10. Quantile Binning  
277  
278 #%% md  
279 ### Quantile Binning  
280  
281 - **Quantile Binning**: The log-transformed variables are further processed using quantile binning to categorize them into discrete categories (e.g., 'Low', 'Medium', 'High').  
282 - **Binned Data**: The binned data provides a simplified view of the variables, making it easier to analyze and interpret their relationships.  
283  
284 #%%  
285 # Applying quantile binning on log-transformed variables  
286 log_binned_data = log_transformed_data.copy()  
287  
288 # Defining the number of bins (e.g., 3 bins for low, medium, high)  
289 n_bins = 3  
290  
291 # Applying quantile binning  
292 log_binned_data['crime_count_binned'] = pd.qcut(log_binned_data['crime_count_log'], q=n_bins, labels=['Low', 'Medium', 'High'])  
293 log_binned_data['Education_Level_3_binned'] = pd.qcut(log_binned_data['Education_Level_3_log'], q=n_bins, labels=['Low', 'Medium', 'High'])  
294 log_binned_data['Migrant_Same_Address_binned'] = pd.qcut(log_binned_data['Migrant_Same_Address_log'], q=n_bins, labels=['Low', 'Medium', 'High'])  
295 log_binned_data['Employment_Worked_12M_binned'] = pd.qcut(log_binned_data['Employment_Worked_12M_log'], q=n_bins, labels=['Low', 'Medium', 'High'])  
296 log_binned_data['Economy_Not_Employed_binned'] = pd.qcut(log_binned_data['Economy_Not_Employed_log'], q=n_bins, labels=['Low', 'Medium', 'High'])  
297 log_binned_data['Ethnic_White_British_binned'] = pd.qcut(log_binned_data['Ethnic_White_British_log'], q=n_bins, labels=['Low', 'Medium', 'High'])  
298  
299 # Display the first few rows of the binned data  
300 print(log_binned_data[['crime_count_binned', 'Education_Level_3_binned', 'Migrant_Same_Address_binned', 'Employment_Worked_12M_binned', 'Economy_Not_Employed_binned', 'Ethnic_White_British_binned']].head())  
301  
302 #%% md  
303 ## 11. Categorical Encoding of Date Variables  
304  
305 #%% md  
306 ### Categorical Encoding of Date Variables  
307  
308 - **Month Encoding**: The 'Month' column, which is initially in datetime format, is converted into a categorical variable representing just the month. This is useful for modeling purposes where month-based seasonality might be relevant.  
309  
310 #%%  
311 # Convert 'Month' to a datetime object if it isn't already  
312 log_binned_data['Month'] = pd.to_datetime(log_binned_data['Month'])  
313  
314 # Extract the month part and encode it as a categorical variable  
315 log_binned_data['MonthEncoded'] = log_binned_data['Month'].dt.month.astype('category')  
316  
317 # Display the first few rows to verify the changes  
318 print(log_binned_data[['Month', 'MonthEncoded']].head())  
319  
320 #%% md  
321 ## 12. Bayesian Network Structure Learning  
322  
323 #%% md  
324 ### Bayesian Network Structure Learning  
325  
326 - **Data Preparation for Bayesian Network**: select relevant binned and encoded variables to build a Bayesian Network model. Categorical variables are converted into numerical codes required for Bayesian modeling.  
327 - **Structure Learning**: The structure of the Bayesian Network is learned using Hill-Climb Search with BIC (Bayesian Information Criterion) and K2 scoring methods. The learned structure provides insights into the relationships between different variables.  
328  
329 #%%  
330 # Selecting the relevant binned and encoded variables for the Bayesian network  
331 bayesian_data = log_binned_data[['MonthEncoded', 'crime_count_binned', 'Education_Level_3_binned', 'Migrant_Same_Address_binned', 'Employment_Worked_12M_binned', 'Economy_Not_Employed_binned', 'Ethnic_White_British_binned']]  
332  
333 # Encoding categorical variables (MonthEncoded is already categorical)  
334 for col in bayesian_data.columns:  
335     if bayesian_data[col].dtype.name == 'category':  
336         bayesian_data[col] = bayesian_data[col].cat.codes  
337  
338 # Performing structure learning using Hill-Climb Search and BIC scoring  
339 hc = HillClimbSearch(bayesian_data)  
340 best_model = hc.estimate(scoring_method=K2Score(bayesian_data))

```

```

344 # Display the structure of the learned Bayesian
345      network
346 print(best_model.edges())
347
348 %% md
349 ## 13. Fitting and Evaluating Bayesian Network
350
351 %% md
352 ### Fitting and Evaluating Bayesian Network
353
354 - **Model Fitting**: The learned Bayesian
      Network structure is fitted to the data
      using the Bayesian Estimator.
355 - **Model Evaluation**: The Bayesian Network is
      evaluated using the BIC and K2 scores to
      assess its performance and adequacy in
      representing the relationships in the data.
356
357 %% md
358 # Fit the Bayesian network
359 model = BayesianNetwork(best_model.edges())
360 model.fit(bayesian_data, estimator=
            BayesianEstimator)
361
362 # Evaluate the model using BIC score
363 bic_score = BicScore(bayesian_data).score(model)
364 print(f"BIC Score: {bic_score}")
365
366 # Evaluate the model using K2 score
367 k2_score = K2Score(bayesian_data).score(
            best_model)
368 print(f"K2 Score: {k2_score}")
369
370 %% md
371 ## 14. Inference and Analysis
372
373 %% md
374 ### Inference and Analysis
375
376 - **Variable Elimination**: perform inference
      using the Variable Elimination method to
      determine the most probable causes of high
      crime rates (categorized as 'crime_count_binned').
377 - **Crime Cause Analysis**: The results of the
      inference provide insights into the likely
      causes of high crime rates based on the
      relationships captured in the Bayesian
      Network.
378
379 %% md
380 # Perform inference to determine the most
      probable causes of high crime_count_binned
381 inference = VariableElimination(model)
382 crime_cause = inference.map_query(variables=['crime_count_binned'])
383 print(f"Most likely cause of crime: {crime_cause}")
384
385 %% md
386 ## 15. Visualization of Bayesian Network
387
388 %% md
389 ### Visualization of Bayesian Network
390
391 - **Network Structure Visualization**: use the
      'networkx' library to visualize the
      structure of the Bayesian Network. The
      circular layout helps in clearly displaying
      the relationships between variables.

392
393 %%%
394 # Extracting the structure of the Bayesian
      Network
395 edges = model.edges()
396
397 # Create a directed graph from the Bayesian
      network structure
398 G = nx.DiGraph()
399 G.add_edges_from(edges)
400
401 # Set node positions using a circular layout
      for better organization
402 pos = nx.circular_layout(G) # Circular layout
      for a more organized appearance
403
404 # Draw the network
405 plt.figure(figsize=(15, 5)) # Increase the
      figure size for a more spacious view
406 nx.draw(
407     G,
408     pos,
409     with_labels=True,
410     node_size=3000, # Even larger nodes
411     node_color='lightgreen', # Softer, calming
412     color for nodes
413     font_size=8, # Slightly larger font for
414     labels
415     font_weight='bold',
416     edge_color='gray', # Softer edge color
417     linewidths=2, # Thicker node borders for
418     emphasis
419     arrowsize=20, # Larger arrows for better
420     visibility
421     arrowstyle='fancy', # Different arrow
422     style for a unique look
423 )
424
425 plt.title("Bayesian Network Structure -"
426           "Circular Layout", fontsize=16) # Larger
427           title font size
428 plt.show()

429 %% md
430 ## 16. Analysis of Variable Influence on Crime
431
432 %% md
433 ### Analysis of Variable Influence on Crime
434
435 - **Influence Plotting**: A custom function is
      created to plot the influence of different
      variables on the probability distribution
      of the target variable ('crime_count_binned'').
      This helps in understanding how changes
      in one variable affect the likelihood of
      different crime levels.

436
437 %%%
438 import matplotlib.pyplot as plt
439 from pgmpy.inference import VariableElimination
440
441 def plot_variable_influence(model, data,
        target_variable='crime_count_binned',
        variables_to_plot=None):
442     """
443     Plots the influence of different variables
        on the target variable's (crime_count_binned) probability
        distribution.
444
445     :param model: The fitted Bayesian network
446     """

```

```

439 model.
440 :param data: The DataFrame used to fit the
441 model.
442 :param target_variable: The target variable
443 to analyze (default is 'crime_count_binned'
444').
445 :param variables_to_plot: A list of
446 variables for which to plot influence on
447 the target variable.
448 """
449
450 if variables_to_plot is None:
451     variables_to_plot = data.columns.tolist()
452
453     variables_to_plot.remove(
454         target_variable)
455
456 # Initialize inference object
457 inference = VariableElimination(model)
458
459 for var in variables_to_plot:
460     # Get unique values of the variable to
461     # plot
462     unique_vals = sorted(data[var].unique())
463
464     probabilities = {bin_label: [] for
465 bin_label in data[target_variable].cat.
466 categories}
467
468     for val in unique_vals:
469         # Query the model for the
470         # probability distribution of the target
471         # variable given the current variable's value
472         query_result = inference.query(
473             variables=[target_variable], evidence={var:
474                 val})
475
476         for bin_label in probabilities:
477             probabilities[bin_label].append(
478                 query_result.values[bin_label])
479
480     # Plot the results
481     plt.figure(figsize=(10, 6))
482     for bin_label, prob_values in
483     probabilities.items():
484         plt.plot(unique_vals, prob_values,
485             marker='o', label=f'{target_variable} = {bin_label}')
486
487         plt.title(f'Influence of {var} on {target_variable}')
488         plt.xlabel(var)
489         plt.ylabel(f'P({target_variable})')
490         plt.legend()
491         plt.grid(True)
492         plt.show()
493
494 # Ensure the relevant columns are converted to
495 # categorical dtype
496 categorical_columns = ['Month_Encoded', ,
497 'crime_count_binned', ,
498 'Education_Level_3_binned', ,
499 'Migrant_Same_Address_binned', ,
500 'Employment_Worked_12M_binned', ,
501 'Economy_Not_Employed_binned', ,
502 'Ethnic_White_British_binned']
503
504 # Convert each column to a categorical dtype
505 for col in categorical_columns:
506     bayesian_data[col] = bayesian_data[col].
507     astype('category')

```

Listing 2. Bayesian Network Implementation

APPENDIX E

FEATURE IMPORTANCE ANALYSIS CODE

```

1 #%% md
2 # 1. Importing Libraries
3
4 **Explanation**:
5 - **pandas** and **numpy**: For data manipulation and numerical operations.
6 - **matplotlib** and **seaborn**: For data visualization.
7 - **sklearn.preprocessing.StandardScaler**: To scale features before feeding them into the model.
8 - **sklearn.model_selection.train_test_split**: To split the dataset into training and testing subsets.
9 - **sklearn.ensemble.RandomForestRegressor**: For training a Random Forest model.
10 - **sklearn.linear_model.LinearRegression**: For training a Linear Regression model.
11 - **sklearn.metrics.mean_squared_error**, r2_score**: To evaluate the models using metrics.
12 - **shap** and **lime.lime_tabular**: For model interpretability.
13
14 #%%
15 import pandas as pd
16 import numpy as np
17 import matplotlib.pyplot as plt
18 import seaborn as sns
19 from sklearn.preprocessing import StandardScaler
20 from sklearn.model_selection import train_test_split
21 from sklearn.ensemble import RandomForestRegressor
22 from sklearn.linear_model import LinearRegression
23 from sklearn.metrics import mean_squared_error, r2_score
24 import shap
25 import lime.lime_tabular
26
27 #%% md
28 # 2. Loading and Preparing the Crime Data
29
30 **Explanation**:
31 - **crime_df**: Load the CSV file containing crime data.
32 - **'Month' to datetime**: Convert the 'Month' column to a datetime object for easier filtering.
33 - **Filter the DataFrame**: Keep only the records from 2021 related to "Violence and sexual offences".
34 - **Modify 'LSOA name'**: Truncate the last 5 characters from 'LSOA name' to generalize it.
35 - **Rename column**: Change the column name to 'lower_tier_local_authorities'.
36 - **Aggregate Data**: Group by 'Month' and 'lower_tier_local_authorities', counting the number of crimes.
37
38 #%%
39 # Load the dataset containing crime data
40 crime_data_path = "/Users/arkamandol/DataspellProjects/Desertation_arka_23023023/data_files/UK_Police_Street_Crime_2018
41
42 -10-01_to_2021_09_31.csv"
43 crime_df = pd.read_csv(crime_data_path)
44
45 # Convert 'Month' to datetime and filter for 2021 and specific crime type
46 crime_df['Month'] = pd.to_datetime(crime_df['Month'])
47 filtered_crime_df = crime_df[(crime_df['Month'].dt.year == 2021) & (crime_df['Crime type'] == "Violence and sexual offences")]
48
49 # Modify 'LSOA name' to remove the last 5 characters
50 filtered_crime_df['LSOA name'] = filtered_crime_df['LSOA name'].str[:-5]
51
52 # Rename 'LSOA name' to 'lower_tier_local_authorities'
53 filtered_crime_df.rename(columns={'LSOA name': 'lower_tier_local_authorities'}, inplace=True)
54
55 # Aggregate the data by 'Month' and 'lower_tier_local_authorities' and count occurrences
56 aggregated_crime_data = filtered_crime_df.groupby(['Month', 'lower_tier_local_authorities']).size().reset_index(name='crime_count')
57
58 #%% md
59 # 3. Loading and Cleaning the Additional Dataset
60
61 **Explanation**:
62 - **file_path**: Load another dataset containing demographic and socio-economic data.
63 - **columns_to_drop**: Identify columns containing redundant 'lower_tier_local_authorities' labels.
64 - **Drop columns**: Remove all but one 'lower_tier_local_authorities' column to avoid duplication.
65 - **Rename column**: Standardize the column name for consistency.
66
67 # Load datasets
68 ethnic_data = pd.read_csv("/Users/arkamandol/DataspellProjects/Desertation_arka_23023023/data_files/Ethnic_Group_Data_Transformed.csv")
69 economy_status_data = pd.read_csv("/Users/arkamandol/DataspellProjects/Desertation_arka_23023023/data_files/Economy_Status_Data_Transformed.csv")
70 employment_history_data = pd.read_csv("/Users/arkamandol/DataspellProjects/Desertation_arka_23023023/data_files/Employment_History_Data_Transformed.csv")
71 education_level_data = pd.read_csv("/Users/arkamandol/DataspellProjects/Desertation_arka_23023023/data_files/Education_Level_Data_Transformed.csv")
72 migrant_data = pd.read_csv("/Users/arkamandol/DataspellProjects/Desertation_arka_23023023/data_files/Migrant_Data_Transformed.csv")
73
74 # Rename the locality code columns to '

```

```

    area_code'
75 ethnic_data.rename(columns={'area_code': 'area_code'}, inplace=True)
76 economy_status_data.rename(columns={'area_code': 'area_code'}, inplace=True)
77 employment_history_data.rename(columns={'area_code': 'area_code'}, inplace=True)
78 education_level_data.rename(columns={'area_code': 'area_code'}, inplace=True)
79 migrant_data.rename(columns={'area_code': 'area_code'}, inplace=True)
80
81 # Merge all datasets on the 'area_code' column
82 merged_data = ethnic_data.merge(
83     economy_status_data, on='area_code', how='outer')
83 merged_data = merged_data.merge(
84     employment_history_data, on='area_code',
85 merged_data = merged_data.merge(
86     education_level_data, on='area_code', how='outer')
85 merged_data = merged_data.merge(migrant_data,
87 # Save the merged dataset to a CSV file (optional)
88 # merged_data.to_csv('path_to/merged_dataset.csv', index=False)
89 merged_data
90 #%%
91 # Load the additional dataset
92 # file_path = '/Users/arkamandal/
93     DataspellProjects/Desertation_arka_23023023
93     /data_files/Merged_Dataset.csv'
94 additional_df = pd.read_csv(file_path)
94 additional_df = merged_data.copy()
95 # Identify and drop redundant 'lower_tier_local_authorities' columns
96 columns_to_drop = [col for col in additional_df
96     .columns if 'lower_tier_local_authorities' in col]
97 columns_to_drop.remove('ethnic_lower_tier_local_authorities') # Keep one relevant column
98
99 # Drop the identified columns
100 df_cleaned = additional_df.drop(columns=columns_to_drop)
101
102 # Rename the relevant column to 'lower_tier_local_authorities'
103 df_cleaned.rename(columns={'ethnic_lower_tier_local_authorities': 'lower_tier_local_authorities'}, inplace=True)
104
105 #%%
106 # 4. Merging the Datasets
107
108 **Explanation**:
109 - **Merge DataFrames**: Combine the crime data and the cleaned demographic data based on 'lower_tier_local_authorities'.
110 - **Display Info**: Print the first few rows and the structure of the merged DataFrame
111
112 #%% to verify the merge.
113 # Merge the crime data with the cleaned dataset on 'lower_tier_local_authorities'
114 merged_data = aggregated_crime_data.merge(
114     df_cleaned, on='lower_tier_local_authorities', how='inner')
115
116 # Display basic information and initial rows of the merged dataset
117 print(merged_data.head())
118 print(merged_data.info())
119
120 #%% md
121 # 5. Preprocessing and Feature Engineering
122
123 **Explanation**:
124 - **Drop columns**: Remove irrelevant columns like 'lower_tier_local_authorities' and 'area_code'.
125 - **Column mapping**: Simplify column names for easier reference.
126 - **Select features**: Choose the most correlated features within each category (e.g., Migrant, Education) with the target variable, 'crime_count'.
127 - **Scale Features**: Standardize the selected features for modeling.
128 - **Convert 'Month'**: Transform the 'Month' column into a numeric format for use in models.
129
130 #%%
131 # Create a copy for processing
132 data = merged_data.copy()
133
134 # Drop columns not needed for analysis
135 data.drop(columns=['lower_tier_local_authorities', 'area_code'],
135     inplace=True)
136
137 # Mapping new column names for easier reference
138 new_column_names = {
138
139     'ethnic_asian_asian_british_or_asian_welsh_bangladeshi':
139         'Ethnic_Asian_Bangladeshi',
140
140     'ethnic_asian_asian_british_or_asian_welsh_chinese':
140         'Ethnic_Asian_Chinese',
141
141     'ethnic_asian_asian_british_or_asian_welsh_indian':
141         'Ethnic_Asian_Indian',
142
142     'ethnic_asian_asian_british_or_asian_welsh_other_asian':
142         'Ethnic_Asian_Other',
143
143     'ethnic_asian_asian_british_or_asian_welsh_pakistani':
143         'Ethnic_Asian_Pakistani',
144
144     'ethnic_black_black_british_black_welsh_caribbean_or_african':
144         'Ethnic_Black_African',
145
145     'ethnic_black_black_british_black_welsh_caribbean_or_african':
145         'Ethnic_Black_Caribbean',
146
146     'ethnic_black_black_british_black_welsh_caribbean_or_african':
146         'Ethnic_Black_Other',
147
147     'ethnic_mixed_or_multiple_ethnic_groups_other_mixed_or_mult':
147         'Ethnic_Mixed_Other',

```

```

148 ,
149 ethnic_mixed_or_multiple_ethnic_groups_white_a
150 ': 'Ethnic_Mixed_White_Asian', 174
151 , 175
152 ethnic_mixed_or_multiple_ethnic_groups_white_a# Rename columns according to the mapping
153 ': 'Ethnic_Mixed_White_Black_African', 177 data.rename(columns=new_column_names, inplace=True)
154 , 178
155 ethnic_mixed_or_multiple_ethnic_groups_white_a# Select columns for further analysis
156 ': 'Ethnic_Mixed_White_Black_Caribbean', 179 prefixes = ['Migrant', 'Education', 'Economy',
157 , 180 'Ethnic', 'Employment']
158 ethnic_other_ethnic_group_any_other_ethnic_gro selected_columns = []
159 ': 'Ethnic_Other', 181
160 'ethnic_other_ethnic_group_arab': ' 182
161 Ethnic_Arab', 183
162 , 184
163 ethnic_white_english_welsh_scottish_northern_
164 ': 'Ethnic_White_British', 185
165 'ethnic_white_gypsy_or_irish_traveller': ' 186
166 Ethnic_White_Gypsy_Traveller', 187
167 'ethnic_white_irish': 'Ethnic_White_Irish', 188
168 'ethnic_white_other_white': ' 189
169 Ethnic_White_Other', 190 # No need to reload the data, proceed with the
170 'ethnic_white_roma': 'Ethnic_White_Roma', 191 selected_columns_with_month = ['Month'] +
171 'economy_employed': 'Economy_Employed', 192 selected_columns + ['crime_count']
172 'economy_not_employed': ' 193 final_data = data[selected_columns_with_month]
173 Economy_Not_Employed', 194
174 'employment_not_in_employment_never_worked' 195 # Scale features
175 ': 'Employment_Never_Worked', 196
176 , 197
177 employment_not_in_employment_not_worked_in_the scaled_features = scaler.fit_transform(
178 ': 'Employment_Not_Worked_12M', 198 final_data.drop(columns=['Month', 'crime_count']))
179 , 199
180 employment_not_in_employment_worked_in_the_las scaled_features_df = pd.DataFrame(
181 ': 'Employment_Worked_12M', 200 scaled_features, columns=selected_columns)
182 , 201
183 education_level_1_and_entry_level_qualification# # Combine scaled features with 'Month' and 'crime_count'
184 *# final_data = pd.concat([final_data[['Month']], scaled_features_df, final_data['crime_count']], axis=1)
185 _to_c_any_gcse_at_other_grades_o_levels_or_les# gnvq_basic_or_
186 ': 'Education_Level_1', 202
187 , 203
188 education_level_2_qualifications_5_or_more_ges# ces_intermediate
189 *# final_data['Month'] = pd.to_datetime(
190 _to_c_or_9_to_4_o_levels_passes_gcse_grade_1_ 201 final_data['Month'].astype(int) // 10**9
191 ': 'Education_Level_2', 202
192 , 203
193 education_level_3_qualifications_2_or_more_a204 certificate_program
194 ; 204 # Extract year and month from the 'Month'
195 _advanced_gnvq_city_and_guilds_advanced_craft_ 205 # column and treat them as categorical
196 ': 'Education_Level_3', 206 # Extract month from the 'Month' column and
197 , 207 # treat it as categorical
198 education_level_4_qualifications_or_above_degre final_data['Month'] = pd.to_datetime(final_data
199 ': 'Education_Level_4+', 208 ['Month']).dt.month.astype('category')
200 , 209
201 education_no_qualifications': ' 209
202 Education_No_Qualifications', 210
203 , 211
204 education_other_apprenticeships_vocational_or_210
205 _related_qualifications_other_qualifications_achi# %% l_outside_england
206 ': 'Education_Other_Qualifications', 211
207 , 212
208 migrant_address_one_year_ago_is_student_term212 final_data
209 -time_or_boarding_school_address_in_the_uk' 213 #%% md
210 ': 'Migrant_Student_UK', 214 # 6. Correlation Matrix Visualization
211 , 215
212 migrant_address_one_year_ago_is_the_same_as216 **Explanation**:
213 ': 'Migrant_Same_Address', 217 - **Correlation Matrix**: Display the
214 'migrant_does_not_apply': 'Migrant_NA', 218 correlation matrix to understand the
215 , 219 relationships between features and the
216 migrant_migrant_from_outside_the_uk_address_on target variable.
217 ': 'Migrant_Outside_UK', 218 - **Heatmap**: Visual representation using
218 , 219 Seaborn's heatmap for easier interpretation

```

```

219 .
220 #%%
221 import matplotlib.pyplot as plt
222 import seaborn as sns
223
224 # Calculate the correlation matrix
225 corr_matrix = final_data.corr()
226
227 # Set up the matplotlib figure
228 plt.figure(figsize=(10, 8))
229
230 # Create a heatmap with annotations and a
231 #   stylish color map
232 sns.heatmap(corr_matrix, annot=True, fmt=".2f",
233               cmap='magma', linewidths=0.5, linecolor='white',
234               cbar_kws={'shrink': 0.8, 'aspect': 20, 'pad': 0.02})
235
236 # Customize the plot with a title, better
237 #   layout, and larger font sizes
238 plt.title('Correlation Matrix of Features and
239 Crime Count', fontsize=18, pad=20)
240 plt.xticks(fontsize=12, rotation=45, ha='right')
241 plt.yticks(fontsize=12)
242
243 # Remove the top and right spines for a cleaner
244 #   look
245 sns.despine()
246
247 # Adjust layout for better spacing
248 plt.tight_layout()
249
250 # Display the heatmap
251 plt.show()
252
253 #%% md
254 # 7. Train-Test Split and Model Training
255
256 #**Explanation**:
257 - **Train-test split**: Split the data into 70%
258 #   training and 30% testing subsets.
259 - **Model Training**: Train both a Random
260 #   Forest model and a Linear Regression model
261 #   on the training data.
262
263 #%% md
264 # Split the data into training and testing sets
265 X = final_data.drop(columns=['crime_count'])
266 y = final_data['crime_count']
267 X_train, X_test, y_train, y_test =
268     train_test_split(X, y, test_size=0.3,
269     random_state=42)
270
271 # Train Random Forest model
272 rf_model = RandomForestRegressor(random_state
273     =42)
274 rf_model.fit(X_train, y_train)
275
276 # Train Linear Regression model
277 linear_model = LinearRegression()
278 linear_model.fit(X_train, y_train)
279
280 #%% md
281 # 8. Model Evaluation
282
283 #**Explanation**:
284 - **Predict**: Generate predictions for the
285 #   test set using both models.
286
287 - **Evaluate**: Calculate and print the Mean
288 #   Squared Error (MSE) and R-squared (R )
289 #   values for both models to assess
290 #   performance.
291
292 #%% md
293 # Predict and evaluate Random Forest
294 y_pred_rf = rf_model.predict(X_test)
295 mse_rf = mean_squared_error(y_test, y_pred_rf)
296 r2_rf = r2_score(y_test, y_pred_rf)
297
298 # Predict and evaluate Linear Regression
299 y_pred_linear = linear_model.predict(X_test)
300 mse_linear = mean_squared_error(y_test,
301     y_pred_linear)
302 r2_linear = r2_score(y_test, y_pred_linear)
303
304 # Print evaluation metrics
305 print("Random Forest Evaluation:")
306 print(f'R-squared (R ): {r2_rf:.4f}')
307 print(f"Mean Squared Error (MSE): {mse_rf:.4f}\n")
308
309 print("Linear Regression Evaluation:")
310 print(f'R-squared (R ): {r2_linear:.4f}')
311 print(f"Mean Squared Error (MSE): {mse_linear
312     :.4 f}\n")
313
314 #%% md
315 # 9. Visualization of Predictions vs Actual
316 #   Values
317
318 #**Explanation**:
319 - **Plot**: Create scatter plots comparing
320 #   actual vs predicted values for both models
321 #   to visually assess their performance.
322
323 #%% md
324 # Visualization: Predicted vs Actual for Random
325 #   Forest and Linear Regression
326 plt.figure(figsize=(14, 6))
327
328 plt.subplot(1, 2, 1)
329 plt.scatter(y_test, y_pred_rf, alpha=0.6, color
330     ='blue')
331 plt.plot([y_test.min(), y_test.max()], [y_test.
332     min(), y_test.max()], 'k--', lw=2)
333 plt.title('Random Forest: Actual vs Predicted')
334 plt.xlabel('Actual Crime Count')
335 plt.ylabel('Predicted Crime Count')
336
337 plt.subplot(1, 2, 2)
338 plt.scatter(y_test, y_pred_linear, alpha=0.6, color
339     ='red')
340 plt.plot([y_test.min(), y_test.max()], [y_test.
341     min(), y_test.max()], 'k--', lw=2)
342 plt.title('Linear Regression: Actual vs
343 Predicted')
344 plt.xlabel('Actual Crime Count')
345 plt.ylabel('Predicted Crime Count')
346
347 plt.tight_layout()
348 plt.show()
349
350 #%% md
351 # 10. Feature Importance Visualization
352
353 #**Explanation**:
354 - **Feature Importance**: Calculate and
355 #   visualize feature importance from the
356 #   Random Forest model and Linear Regression

```

```

model coefficients.

#%%
# Random Forest Feature Importance
importances_rf = rf_model.feature_importances_
importance_df_rf = pd.DataFrame({'Feature':
    X_train.columns, 'Importance':
    importances_rf}).sort_values(by='Importance',
        ascending=False)

# Linear Regression Coefficients
coefficients_lr = linear_model.coef_
importance_df_lr = pd.DataFrame({'Feature':
    X_train.columns, 'Importance':
    coefficients_lr}).sort_values(by='
        Importance', ascending=False)

# Visualization of Feature Importances
plt.figure(figsize=(12, 6))

plt.subplot(1, 2, 1)
plt.barh(importance_df_rf['Feature'],
    importance_df_rf['Importance'], color='
        skyblue')
plt.title('Random Forest Feature Importance')
plt.xlabel('Importance')

plt.subplot(1, 2, 2)
plt.barh(importance_df_lr['Feature'], abs(
    importance_df_lr['Importance']), color='
        lightcoral')
plt.title('Linear Regression Feature Importance
        ')
plt.xlabel('Absolute Coefficient Value')

plt.tight_layout()
plt.show()

#%%
# 11. SHAP Values for Interpretation

**Explanation**:
- **SHAP Values**: Use SHAP to interpret
    feature importance and impact on
    predictions for both the Random Forest and
    Linear Regression models.

#%%
# SHAP interpretation for Random Forest
explainer_rf = shap.TreeExplainer(rf_model)
shap_values_rf = explainer_rf.shap_values(
    X_test)
shap.summary_plot(shap_values_rf, X_test,
    plot_type="bar")

# SHAP interpretation for Linear Regression
explainer_lr = shap.LinearExplainer(
    linear_model, X_train)
shap_values_lr = explainer_lr.shap_values(
    X_test)
shap.summary_plot(shap_values_lr, X_test,
    plot_type="bar")

#%%
# 12. LIME for Interpretation

**Explanation**:
- **LIME**: Generate local interpretable model-
    agnostic explanations using LIME for both
models on a specific instance from the test
set.

#%%
# LIME interpretation for Random Forest
explainer_rf_lime = lime.lime_tabular.
    LimeTabularExplainer(
        training_data=X_train.values,
        feature_names=X_train.columns,
        mode='regression'
    )
instance_to_explain = X_test.iloc[0]
explanation_rf = explainer_rf_lime.
    explain_instance(
        data_row=instance_to_explain.values,
        predict_fn=rf_model.predict
    )
explanation_rf.show_in_notebook(show_table=True
    )

# LIME interpretation for Linear Regression
explainer_lr_lime = lime.lime_tabular.
    LimeTabularExplainer(
        training_data=X_train.values,
        feature_names=X_train.columns,
        mode='regression'
    )
explanation_lr = explainer_lr_lime.
    explain_instance(
        data_row=instance_to_explain.values,
        predict_fn=linear_model.predict
    )
explanation_lr.show_in_notebook(show_table=True
    )

# Save LIME interpretation to HTML
rf_html_file = 'lime_rf_explanation.html'
explanation_rf.save_to_file(rf_html_file)
# Save LIME interpretation to HTML
lr_html_file = 'lime_lr_explanation.html'
explanation_lr.save_to_file(lr_html_file)

```

Listing 3. Feature Importance Analysis

APPENDIX F LSTM HOTSPOT MAPPING

```

1 #%%
2 import pandas as pd
3 import numpy as np
4
5 # Path to your dataset
6 file_path = '/Users/arkamadol/
7             DataspellProjects/Desertation_arka_23023023
8             /data_files/uk_crime_lat_long.csv',
9 data = pd.read_csv(file_path)
10
11 #%%
12 data.head()
13
14 # Convert 'Month' to datetime format for better
15 # date handling
16 data['Month'] = pd.to_datetime(data['Month'])
17
18 # Sort the DataFrame by the 'Month' column
19 data.sort_values('Month', inplace=True)
20 data.head()
21
22 # Group data by Month, Latitude, and Longitude
23 # and count occurrences
24 monthly_data = data.groupby(['Month', 'Latitude',
25                             'Longitude']).size().reset_index(name='Count')
26 monthly_data.head()
27
28 from sklearn.preprocessing import MinMaxScaler
29
30 # Initialize the MinMaxScaler
31 scaler = MinMaxScaler(feature_range=(0, 1))
32
33 # Fit and transform the 'Count' data to scale
34 # it
35 monthly_data['Normalized_Count'] = scaler.
36     fit_transform(monthly_data[['Count']])
37
38 # Function to create input sequences for LSTM
39 def create_dataset(dataset, look_back=1):
40     dataX, dataY = [], []
41     for i in range(len(dataset) - look_back):
42         a = dataset[i:(i + look_back)]
43         dataX.append(a)
44         dataY.append(dataset[i + look_back])
45     return np.array(dataX), np.array(dataY)
46
47 # Define the number of past months data to
48 # consider for predicting the next month
49 look_back = 3
50 X, y = create_dataset(monthly_data[
51     'Normalized_Count'].values, look_back)
52
53 # Reshape input to be [samples, time steps,
54 # features] for LSTM
55 X = np.reshape(X, (X.shape[0], look_back, 1))
56
57 # Define the LSTM model
58 model = Sequential()
59 model.add(LSTM(50, input_shape=(look_back, 1)))
60 # 50 LSTM units
61
62 # Output layer that
63 # predicts the future value
64 model.compile(loss='mean_squared_error',
65                 optimizer='adam')
66
67 # Fit the model on the dataset
68 model.fit(X, y, epochs=10, batch_size=512,
69             verbose=1)
70
71 # Predict future values
72 predictions = model.predict(X)
73
74 # Inverse transform to get predictions in the
75 # original count scale
76 predictions = scaler.inverse_transform(
77     predictions)
78
79 # Add predictions to the monthly_data DataFrame
80 monthly_data['Predicted_Count'] = np.
81     concatenate([np.zeros(look_back),
82                 predictions.flatten()])
83
84 # Create a map centered around an average
85 # location
86 center_lat, center_lon = monthly_data['Latitude'.
87     mean(), monthly_data['Longitude'].mean()]
88 map = folium.Map(location=[center_lat,
89                   center_lon], zoom_start=6)
90
91 # Add a heatmap to the map using predicted
92 # crime counts
93 heat_data = [[row['Latitude'], row['Longitude'],
94               row['Predicted_Count']] for index, row
95   in monthly_data.iterrows()]
96 HeatMap(heat_data).add_to(map)
97
98 # Save or display the map
99 map.save('/Users/arkamadol/DataspellProjects/
100 Desertation_arka_23023023/data_files/
101 crime_hotspots.html')
102
103 # Calculate RMSE
104 rmse = sqrt(mean_squared_error(y, predictions))
105 print('Root Mean Squared Error:', rmse)
106
107 #%%
108
109 #%%

```

Listing 4. LSTM Hotspot Mapping Script