

# Houston Crime Analysis

Preethu Manjunath  
University of Colorado Boulder  
Boulder, United States

Lohith Ramesh  
University of Colorado Boulder  
Boulder, United States

Raquel Yupanqui  
University of Colorado Boulder  
Boulder, United States

## ABSTRACT

In the rapidly evolving urban landscape of Houston, analyzing crime data from 2019 to 2024 unveils crucial insights into the dynamics of urban crime and its profound impact on community welfare. This report delves into an exhaustive exploration of Houston's crime data, leveraging advanced data analysis techniques to dissect crime patterns, temporal trends, and socioeconomic correlations. By integrating crime data with meticulously web-scraped demographic information, our analysis offers a holistic view of the interplay between crime rates and socioeconomic factors across different neighborhoods. The analysis commenced with the aggregation and pre-processing of crime data sourced from the official Houston city website, spanning six years to construct a comprehensive dataset. This dataset underwent rigorous cleaning processes, including handling missing values, converting data types, and standardizing entries, to ensure accuracy and reliability. Furthermore, feature engineering and data reduction techniques were applied to distill key variables that shed light on crime trends and patterns. A pivotal aspect of our methodology was the enrichment of crime data with demographic variables such as population density, income levels, and unemployment rates. This integration enabled a multifaceted analysis, revealing nuanced insights into how socioeconomic disparities influence crime occurrences. Through a suite of visualizations, including time-series graphs, scatter map plots, and heatmaps, we illustrated the geographic distribution of crimes, temporal patterns, and the relationship between crime rates and demographic indicators. Key findings underscore the significance of socioeconomic factors in shaping crime dynamics, with particular emphasis on the correlation between economic hardship and increased crime rates. Geographic analysis identified specific hotspots of criminal activity, highlighting the need for targeted public safety interventions. This report underscores the critical role of data-driven approaches in understanding and mitigating urban crime. The insights garnered from this analysis advocate for informed policy-making and strategic planning to enhance public safety and address the underlying socioeconomic factors contributing to crime. Our findings serve as a call to action for stakeholders, including law enforcement, policymakers, and community organizations, to collaborate in fostering a safer and more equitable urban environment.

## KEYWORDS

Machine Learning, Data Science, Data Mining

### ACM Reference Format:

Preethu Manjunath, Lohith Ramesh, and Raquel Yupanqui. 2024. Houston Crime Analysis. In *Proceedings of ACM Conference (Conference'17)*. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## 1 INTRODUCTION

Urban crime represents a significant challenge for major metropolitan areas worldwide, intricately linked to a myriad of socioeconomic factors that influence public safety and community stability. This study focuses on Houston, one of the United States' most dynamic and diverse cities, to conduct an exhaustive analysis of crime trends from 2019 to 2024. The primary aim is to leverage a data-driven approach to dissect the complex interrelations between crime rates and socioeconomic environments, thereby providing actionable insights for effective crime prevention and policy formulation.

Initiating this comprehensive research, the project team collected extensive crime statistics from the City of Houston's official online portals, ensuring a robust dataset that includes crime type, location, date, and time of incidents. To enrich the analysis, this crime data was augmented with socioeconomic and demographic information meticulously web-scraped from various reliable sources. This enriched dataset incorporates additional variables such as population density, income levels, unemployment rates, and educational attainment across different Houston neighborhoods, allowing for a multidimensional analysis of crime dynamics.

With the aid of advanced analytical tools and statistical methods, the study undertakes a systematic examination of the data to identify temporal and spatial patterns of crime occurrences. Moreover, machine learning techniques including Random Forest, K-Nearest Neighbors (KNN), AdaBoost, and XGBoost are employed to develop predictive models that not only forecast crime occurrences but also explore the impact of socioeconomic disparities on crime rates.

This introductory section sets the groundwork for a detailed exploration of the methods and techniques used during the research, provides an overview of the data preparation and preprocessing efforts, and outlines the analytical strategies implemented. Following this, the report presents the key findings, discusses their implications for crime prevention and urban planning, and offers recommendations for policymakers and community leaders aimed at enhancing public safety and fostering a more equitable urban environment.

By combining rigorous data analysis with insightful socioeconomic considerations, this study aspires to contribute to a more nuanced understanding of urban crime and to support data-driven decision-making in Houston's ongoing efforts to mitigate crime and enhance community welfare.

2 DATA COLLECTION/PREPARATION

The Houston Crime Analysis project embarked on a multifaceted journey to dissect urban crime through the lens of data-driven insights. Commencing in 2019 and extending through 2024, this endeavor aimed to unearth patterns, trends, and correlations within the vast landscape of crime data, enriched by socioeconomic and demographic dimensions. The project's backbone was the meticulous gathering, cleaning, and preparation of crime data juxtaposed with demographic insights, setting a robust foundation for nuanced analysis.

2.1 Data Sourcing

The primary dataset, encompassing detailed monthly crime reports by street and police beat, was sourced from the City of Houston's official website[1]. This rich repository of crime incidents provided a comprehensive canvas for the analysis. The crime data, segmented into yearly Excel files from 2019 to 2024, was methodically downloaded. Each file's initial inspection offered insights into the dataset's structure, revealing variables' range, data types, and potential quality issues such as missing values or inconsistencies. Utilizing Python libraries Pandas for data manipulation and Matplotlib/Plotly for visualization, the yearly data was loaded into separate Pandas DataFrames. This step was crucial for acclimating to the dataset's nuances, setting the stage for in-depth exploration and analysis.

```
import pandas as pd
import numpy as np
from bs4 import BeautifulSoup
import requests
import plotly.express as px
import seaborn as sns
import matplotlib.pyplot as plt
from scipy import stats
import plotly.graph_objects as go
```

Incident	RMSOccurrenceDate	RMSOccurrenceHour	NIBRSClass	NIBRSDescription	OffenseCount	Beat	Premise	StreetName	City	ZipCode	Map	
0	5619	2019-01-01	0	290	Destruction, damage, vandalism	1	9C30	Residence, Home (includes Apartment)	SAN CARLOS	HOUSTON	77013	-8
1	17319	2019-01-01	0	35A	Drug, narcotic violations	1	7C10	Highway, Road, Street, Alley	EAST HOUSTON	HOUSTON	77020	-4
2	17319	2019-01-01	0	900	Driving under the influence	1	7C10	Highway, Road, Street, Alley	EAST HOUSTON	HOUSTON	77020	-4
3	18119	2019-01-01	0	290	Destruction, damage, vandalism	1	18E40	Residence, Home (includes Apartment)	LONE QUAIL	HOUSTON	77469	-4
4	19019	2019-01-01	0	620	Weapon law violations	1	NAN	Residence, Home (includes Apartment)	MELBOURNE	HOUSTON	77026	-6

3 PREPROCESSING

3.1 Handling Missing Values

Strategies were employed based on the data's nature and missingness. Notably, the street suffix column was omitted due to excessive missing values, and missing values in the "Beat" column were intelligently filled by referencing incident zip codes. Latitude and longitude coordinates served to fill missing zip codes, minimizing data loss.

```
region = {
    'Airport-Hobby Division - District 22': ['22140', '22150'],
    'Airport-CMB Division - District 21': ['21138', '21148', '21158', '21168', '21178'],
    'Central Division - District 12': ['12038', '12048', '12058', '12068', '12078', '12088', '12098', '12108', '12118', '12128', '12138', '12148', '12158', '12168', '12178', '12188', '12198', '12208', '12218', '12228', '12238', '12248', '12258', '12268', '12278', '12288', '12298', '12308', '12318', '12328', '12338', '12348', '12358', '12368', '12378', '12388', '12398', '12408', '12418', '12428', '12438', '12448', '12458', '12468', '12478', '12488', '12498', '12508', '12518', '12528', '12538', '12548', '12558', '12568', '12578', '12588', '12598', '12608', '12618', '12628', '12638', '12648', '12658', '12668', '12678', '12688', '12698', '12708', '12718', '12728', '12738', '12748', '12758', '12768', '12778', '12788', '12798', '12808', '12818', '12828', '12838', '12848', '12858', '12868', '12878', '12888', '12898', '12908', '12918', '12928', '12938', '12948', '12958', '12968', '12978', '12988', '12998'],
    'Clear Lake Division - District 32': ['32038', '32048', '32058', '32068', '32078', '32088', '32098', '32108', '32118', '32128', '32138', '32148', '32158', '32168', '32178', '32188', '32198', '32208', '32218', '32228', '32238', '32248', '32258', '32268', '32278', '32288', '32298', '32308', '32318', '32328', '32338', '32348', '32358', '32368', '32378', '32388', '32398', '32408', '32418', '32428', '32438', '32448', '32458', '32468', '32478', '32488', '32498', '32508', '32518', '32528', '32538', '32548', '32558', '32568', '32578', '32588', '32598', '32608', '32618', '32628', '32638', '32648', '32658', '32668', '32678', '32688', '32698', '32708', '32718', '32728', '32738', '32748', '32758', '32768', '32778', '32788', '32798', '32808', '32818', '32828', '32838', '32848', '32858', '32868', '32878', '32888', '32898', '32908', '32918', '32928', '32938', '32948', '32958', '32968', '32978', '32988', '32998'],
    'Downtown Division - District 11': ['11038', '11048', '11058', '11068', '11078', '11088', '11098', '11108', '11118', '11128', '11138', '11148', '11158', '11168', '11178', '11188', '11198', '11208', '11218', '11228', '11238', '11248', '11258', '11268', '11278', '11288', '11298', '11308', '11318', '11328', '11338', '11348', '11358', '11368', '11378', '11388', '11398', '11408', '11418', '11428', '11438', '11448', '11458', '11468', '11478', '11488', '11498', '11508', '11518', '11528', '11538', '11548', '11558', '11568', '11578', '11588', '11598', '11608', '11618', '11628', '11638', '11648', '11658', '11668', '11678', '11688', '11698', '11708', '11718', '11728', '11738', '11748', '11758', '11768', '11778', '11788', '11798', '11808', '11818', '11828', '11838', '11848', '11858', '11868', '11878', '11888', '11898', '11908', '11918', '11928', '11938', '11948', '11958', '11968', '11978', '11988', '11998'],
    'Eastside Division - District 31': ['31038', '31048', '31058', '31068', '31078', '31088', '31098', '31108', '31118', '31128', '31138', '31148', '31158', '31168', '31178', '31188', '31198', '31208', '31218', '31228', '31238', '31248', '31258', '31268', '31278', '31288', '31298', '31308', '31318', '31328', '31338', '31348', '31358', '31368', '31378', '31388', '31398', '31408', '31418', '31428', '31438', '31448', '31458', '31468', '31478', '31488', '31498', '31508', '31518', '31528', '31538', '31548', '31558', '31568', '31578', '31588', '31598', '31608', '31618', '31628', '31638', '31648', '31658', '31668', '31678', '31688', '31698', '31708', '31718', '31728', '31738', '31748', '31758', '31768', '31778', '31788', '31798', '31808', '31818', '31828', '31838', '31848', '31858', '31868', '31878', '31888', '31898', '31908', '31918', '31928', '31938', '31948', '31958', '31968', '31978', '31988', '31998'],
    'Kingwood Division - District 24': ['24038', '24048', '24058', '24068', '24078', '24088', '24098', '24108', '24118', '24128', '24138', '24148', '24158', '24168', '24178', '24188', '24198', '24208', '24218', '24228', '24238', '24248', '24258', '24268', '24278', '24288', '24298', '24308', '24318', '24328', '24338', '24348', '24358', '24368', '24378', '24388', '24398', '24408', '24418', '24428', '24438', '24448', '24458', '24468', '24478', '24488', '24498', '24508', '24518', '24528', '24538', '24548', '24558', '24568', '24578', '24588', '24598', '24608', '24618', '24628', '24638', '24648', '24658', '24668', '24678', '24688', '24698', '24708', '24718', '24728', '24738', '24748', '24758', '24768', '24778', '24788', '24798', '24808', '24818', '24828', '24838', '24848', '24858', '24868', '24878', '24888', '24898', '24908', '24918', '24928', '24938', '24948', '24958', '24968', '24978', '24988', '24998'],
    'Midwest Division - District 30': ['30038', '30048', '30058', '30068', '30078', '30088', '30098', '30108', '30118', '30128', '30138', '30148', '30158', '30168', '30178', '30188', '30198', '30208', '30218', '30228', '30238', '30248', '30258', '30268', '30278', '30288', '30298', '30308', '30318', '30328', '30338', '30348', '30358', '30368', '30378', '30388', '30398', '30408', '30418', '30428', '30438', '30448', '30458', '30468', '30478', '30488', '30498', '30508', '30518', '30528', '30538', '30548', '30558', '30568', '30578', '30588', '30598', '30608', '30618', '30628', '30638', '30648', '30658', '30668', '30678', '30688', '30698', '30708', '30718', '30728', '30738', '30748', '30758', '30768', '30778', '30788', '30798', '30808', '30818', '30828', '30838', '30848', '30858', '30868', '30878', '30888', '30898', '30908', '30918', '30928', '30938', '30948', '30958', '30968', '30978', '30988', '30998'],
    'North Division - District 33': ['33038', '33048', '33058', '33068', '33078', '33088', '33098', '33108', '33118', '33128', '33138', '33148', '33158', '33168', '33178', '33188', '33198', '33208', '33218', '33228', '33238', '33248', '33258', '33268', '33278', '33288', '33298', '33308', '33318', '33328', '33338', '33348', '33358', '33368', '33378', '33388', '33398', '33408', '33418', '33428', '33438', '33448', '33458', '33468', '33478', '33488', '33498', '33508', '33518', '33528', '33538', '33548', '33558', '33568', '33578', '33588', '33598', '33608', '33618', '33628', '33638', '33648', '33658', '33668', '33678', '33688', '33698', '33708', '33718', '33728', '33738', '33748', '33758', '33768', '33778', '33788', '33798', '33808', '33818', '33828', '33838', '33848', '33858', '33868', '33878', '33888', '33898', '33908', '33918', '33928', '33938', '33948', '33958', '33968', '33978', '33988', '33998'],
    'Northeast Division - District 23': ['23038', '23048', '23058', '23068', '23078', '23088', '23098', '23108', '23118', '23128', '23138', '23148', '23158', '23168', '23178', '23188', '23198', '23208', '23218', '23228', '23238', '23248', '23258', '23268', '23278', '23288', '23298', '23308', '23318', '23328', '23338', '23348', '23358', '23368', '23378', '23388', '23398', '23408', '23418', '23428', '23438', '23448', '23458', '23468', '23478', '23488', '23498', '23508', '23518', '23528', '23538', '23548', '23558', '23568', '23578', '23588', '23598', '23608', '23618', '23628', '23638', '23648', '23658', '23668', '23678', '23688', '23698', '23708', '23718', '23728', '23738', '23748', '23758', '23768', '23778', '23788', '23798', '23808', '23818', '23828', '23838', '23848', '23858', '23868', '23878', '23888', '23898', '23908', '23918', '23928', '23938', '23948', '23958', '23968', '23978', '23988', '23998'],
    'Northwest Division - District 34': ['34038', '34048', '34058', '34068', '34078', '34088', '34098', '34108', '34118', '34128', '34138', '34148', '34158', '34168', '34178', '34188', '34198', '34208', '34218', '34228', '34238', '34248', '34258', '34268', '34278', '34288', '34298', '34308', '34318', '34328', '34338', '34348', '34358', '34368', '34378', '34388', '34398', '34408', '34418', '34428', '34438', '34448', '34458', '34468', '34478', '34488', '34498', '34508', '34518', '34528', '34538', '34548', '34558', '34568', '34578', '34588', '34598', '34608', '34618', '34628', '34638', '34648', '34658', '34668', '34678', '34688', '34698', '34708', '34718', '34728', '34738', '34748', '34758', '34768', '34778', '34788', '34798', '34808', '34818', '34828', '34838', '34848', '34858', '34868', '34878', '34888', '34898', '34908', '34918', '34928', '34938', '34948', '34958', '34968', '34978', '34988', '34998'],
    'South Central Division - District 35': ['35038', '35048', '35058', '35068', '35078', '35088', '35098', '35108', '35118', '35128', '35138', '35148', '35158', '35168', '35178', '35188', '35198', '35208', '35218', '35228', '35238', '35248', '35258', '35268', '35278', '35288', '35298', '35308', '35318', '35328', '35338', '35348', '35358', '35368', '35378', '35388', '35398', '35408', '35418', '35428', '35438', '35448', '35458', '35468', '35478', '35488', '35498', '35508', '35518', '35528', '35538', '35548', '35558', '35568', '35578', '35588', '35598', '35608', '35618', '35628', '35638', '35648', '35658', '35668', '35678', '35688', '35698', '35708', '35718', '35728', '35738', '35748', '35758', '35768', '35778', '35788', '35798', '35808', '35818', '35828', '35838', '35848', '35858', '35868', '35878', '35888', '35898', '35908', '35918', '35928', '35938', '35948', '35958', '35968', '35978', '35988', '35998'],
    'South Division - District 36': ['36038', '36048', '36058', '36068', '36078', '36088', '36098', '36108', '36118', '36128', '36138', '36148', '36158', '36168', '36178', '36188', '36198', '36208', '36218', '36228', '36238', '36248', '36258', '36268', '36278', '36288', '36298', '36308', '36318', '36328', '36338', '36348', '36358', '36368', '36378', '36388', '36398', '36408', '36418', '36428', '36438', '36448', '36458', '36468', '36478', '36488', '36498', '36508', '36518', '36528', '36538', '36548', '36558', '36568', '36578', '36588', '36598', '36608', '36618', '36628', '36638', '36648', '36658', '36668', '36678', '36688', '36698', '36708', '36718', '36728', '36738', '36748', '36758', '36768', '36778', '36788', '36798', '36808', '36818', '36828', '36838', '36848', '36858', '36868', '36878', '36888', '36898', '36908', '36918', '36928', '36938', '36948', '36958', '36968', '36978', '36988', '36998'],
    'Southeast Division - District 13': ['13038', '13048', '13058', '13068', '13078', '13088', '13098', '13108', '13118', '13128', '13138', '13148', '13158', '13168', '13178', '13188', '13198', '13208', '13218', '13228', '13238', '13248', '13258', '13268', '13278', '13288', '13298', '13308', '13318', '13328', '13338', '13348', '13358', '13368', '13378', '13388', '13398', '13408', '13418', '13428', '13438', '13448', '13458', '13468', '13478', '13488', '13498', '13508', '13518', '13528', '13538', '13548', '13558', '13568', '13578', '13588', '13598', '13608', '13618', '13628', '13638', '13648', '13658', '13668', '13678', '13688', '13698', '13708', '13718', '13728', '13738', '13748', '13758', '13768', '13778', '13788', '13798', '13808', '13818', '13828', '13838', '13848', '13858', '13868', '13878', '13888', '13898', '13908', '13918', '13928', '13938', '13948', '13958', '13968', '13978', '13988', '13998'],
    'Southwest Division - District 15': ['15038', '15048', '15058', '15068', '15078', '15088', '15098', '15108', '15118', '15128', '15138', '15148', '15158', '15168', '15178', '15188', '15198', '15208', '15218', '15228', '15238', '15248', '15258', '15268', '15278', '15288', '15298', '15308', '15318', '15328', '15338', '15348', '15358', '15368', '15378', '15388', '15398', '15408', '15418', '15428', '15438', '15448', '15458', '15468', '15478', '15488', '15498', '15508', '15518', '15528', '15538', '15548', '15558', '15568', '15578', '15588', '15598', '15608', '15618', '15628', '15638', '15648', '15658', '15668', '15678', '15688', '15698', '15708', '15718', '15728', '15738', '15748', '15758', '15768', '15778', '15788', '15798', '15808', '15818', '15828', '15838', '15848', '15858', '15868', '15878', '15888', '15898', '15908', '15918', '15928', '15938', '15948', '15958', '15968', '15978', '15988', '15998'],
    'Westside Division - District 19': ['19038', '19048', '19058', '19068', '19078', '19088', '19098', '19108', '19118', '19128', '19138', '19148', '19158', '19168', '19178', '19188', '19198', '19208', '19218', '19228', '19238', '19248', '19258', '19268', '19278', '19288', '19298', '19308', '19318', '19328', '19338', '19348', '19358', '19368', '19378', '19388', '19398', '19408', '19418', '19428', '19438', '19448', '19458', '19468', '19478', '19488', '19498', '19508', '19518', '19528', '19538', '19548', '19558', '19568', '19578', '19588', '19598', '19608', '19618', '19628', '19638', '19648', '19658', '19668', '19678', '19688', '19698', '19708', '19718', '19728', '19738', '19748', '19758', '19768', '19778', '19788', '19798', '19808', '19818', '19828', '19838', '19848', '19858', '19868', '19878', '19888', '19898', '19908', '19918', '19928', '19938', '19948', '19958', '19968', '19978', '19988', '19998'],
    'Harris County Sheriff's Office': ['HCSO'],
    'Harris County Constable's Office Precinct 5': ['HCC5'],
    '003': ['003']
}

def get_region(beat):
    for key, value in region.items():
        if beat in value:
            return key
    return None

# Apply the function to create the "region" column
df['Region'] = df['Beat'].apply(get_region)
```

3.2 Data Type Conversion and Standardization

Data types were meticulously adjusted for analysis efficiency and accuracy. Zip codes were transitioned from floats to integers, and date strings to Datetime objects, facilitating accurate temporal analyses. Efforts ensured data consistency, particularly for textual data, and duplicate checks confirmed the dataset's quality.

3.3 Feature Engineering and Data Reduction

The dataset was enhanced by extracting meaningful features such as month and year from crime incident dates, enabling detailed pattern analysis. Irrelevant columns, like street type and number, were pruned to focus the analysis, streamlining the dataset for clarity and purpose.

RMSOccurrenceDate	month	year
2019-01-01	1	2019
2019-01-01	1	2019
2019-01-01	1	2019

3.4 Identify Attribute Types

In data analysis, identifying attribute types is essential as it determines the appropriate analytical approach for each variable. Attributes can be classified as nominal, where data represent categories without inherent order; ordinal, which have a clear sequence; interval, where the differences between data points are meaningful; and ratio, which have a true zero point, allowing for a full range of statistical operations. Understanding these types aids in selecting the correct visualization techniques and statistical tests, ensuring accurate and meaningful analysis outcomes. For example, crime incident data often contain nominal attributes like crime type and ratio attributes like incident counts, each requiring different analytical treatment. Correctly identifying these attribute types is crucial for data pre-processing, feature engineering, and subsequent modeling stages.

3.5 Compute Measures of Central Tendency and Dispersion

Measures of central tendency and dispersion are fundamental statistical tools used to summarize data. Central tendency measures, including the mean, median, and mode, provide a central value around which data points cluster. In contrast, measures of dispersion like the range, variance, and standard deviation quantify the spread of data points around the central value, indicating the variability within the dataset. Together, these measures offer a snapshot of the data's overall shape and spread, informing decisions and highlighting potential outliers or trends. For instance, in a crime dataset, the mean could reveal the average number of incidents, while the standard deviation could show the consistency of crime across different areas.

3.6 Visualize Data Distributions

Visualizing data distributions is a critical step in data analysis, as it provides a graphical representation of how data points are spread across different values. Common methods include histograms, box plots, and density plots, each highlighting various aspects of the data’s distribution, such as skewness, kurtosis, and the presence of outliers. These visual tools help analysts and stakeholders quickly grasp the underlying patterns and anomalies within the data, which may not be immediately apparent from raw figures. For crime data, such visualizations can reveal the frequency of different crime types or the distribution of incidents over time, facilitating pattern recognition and strategic planning. Effective visualization of data distributions is therefore key to intuitive understanding and informed decision-making.

3.7 Assess Data Similarity

Assessing data similarity involves quantifying the likeness between data points, which is crucial for tasks like clustering, classification, and anomaly detection. Techniques such as Euclidean distance for continuous data or the Jaccard index for categorical data are often employed. In a crime database, assessing similarity might help identify clusters of similar crime types or detect areas with comparable crime statistics. This evaluation aids in understanding relationships within the data, facilitating predictive modeling and tailored intervention strategies. Overall, data similarity measures are instrumental in uncovering hidden patterns and associations in complex datasets.

3.8 Normalization

Normalization is the process of scaling individual data points to ensure consistency across different ranges for comparative and analytical purposes. It is especially important when dealing with features that vary in scale, units, or range, as it can impact the performance of many machine learning algorithms. Common methods include min-max scaling, which rescales the range of features to scale the range in [0, 1], and Z-score standardization, where data points are rescaled based on their mean and standard deviation. In the context of crime statistics, normalization allows for fair comparison across variables like income levels and crime rates. Effective normalization can significantly improve the robustness and accuracy of downstream analytical tasks.

3.9 Transformation

Data transformation involves converting data into a different format or structure, often to improve its suitability for analysis or model building. It can include operations like logarithmic transformations to stabilize variance or handle skewness, or more complex functions like Fourier transforms to identify frequencies within time series data. For example, transforming crime incident timestamps into cyclical features can capture the periodic nature of crime occurrences. Such transformations are pivotal for exposing meaningful patterns in the data, enhancing the predictive power of analytical models, and simplifying complex relationships between variables for better interpretability.

3.10 Integration with Demographic Data

The project expanded its analytical horizon by weaving in web-scraped demographic data from [2], adding layers of socioeconomic context. This comprehensive approach integrated data on population density, income levels, unemployment rates, and more, aligned with the crime data via geographical identifiers like ZIP codes. Any missing demographic values were filled in with the column average.

3.11 Creation of a Unified Dataset

The unification of annual crime data into a single DataFrame addressed the challenge of dataset consistency. This meticulous process involved standardizing column names, ensuring uniformity in categorical data, and aligning datetime formats. Demographic data integration was a critical step, enriching the crime dataset with socioeconomic insights. This phase required careful data alignment, particularly in matching geographical identifiers, to enable holistic analyses of crime within Houston’s broader socioeconomic and demographic landscape. Post-merge, the dataset underwent additional feature engineering and verification steps to ensure data integrity and analytical readiness.

Latitude	Population	BlackPercent	AsianPercent	WhitePercent	Hispanic	MedianAge	Unemployment	ChildPoverty	ChildPoverty	ChildPoverty	MedianHouseholdIncome	AverageEducation	PercentForeignBorn	UnemploymentRate
77001	72279	11.05	12.35	66.85	10.75	35.5	11.2	20.3	16.4	25000.0	12.36	0.090		
77002	34288	48.07	4.38	23.46	23.9	35.3	10.8	20.0	20.5	41000.0	12.25	0.090		
77003	30500	15.75	2.89	60.45	20.8	40.7	17.5	22.4	40.0	30000.0	10.00	0.090		
77007	46000	10.40	9.86	62.34	17.4	35.3	17.7	21.2	31.0	40000.0	10.00	0.047		
77007	46000	6.71	11.41	67.00	14.8	35.3	18.7	24.5	40.0	40000.0	11.79	0.090		

4 CLEANING

4.1 Overview

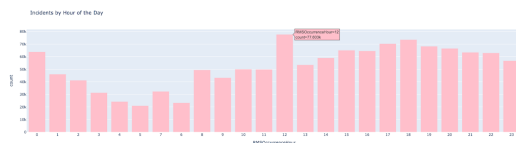
The cleaning phase was a critical step in ensuring the integrity and reliability of the Houston Crime Data analysis. This process involved meticulous scrutiny and refinement of the dataset, addressing various issues such as missing values, inconsistent entries, and data type inaccuracies. Each cleaning action was tailored to preserve data quality while ensuring the analysis remained robust and meaningful.

4.2 Handling Missing Values

The dataset’s integrity was challenged by missing data across several columns. Notably, the ‘street suffix’ column was excluded due to a high volume of missing values which, upon evaluation, were deemed non-essential for the analysis. For the ‘Beat’ column, missing values were addressed by cross-referencing and filling in data based on incident zip codes, demonstrating a nuanced approach to preserving data quality. Similarly, missing zip codes were cleverly inferred using latitude and longitude coordinates, showcasing the innovative problem-solving techniques employed during the cleaning process.

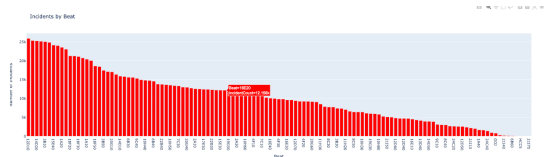
4.3 Data Type Conversion

Critical to the analysis was ensuring that each column’s data type accurately reflected its content. Zip codes, initially recorded as float values, were converted to integer format to correct their representation and facilitate accurate geographical analysis. Date columns underwent transformation from string to Datetime format, enhancing the dataset’s utility for time-series analysis and temporal trend identification.



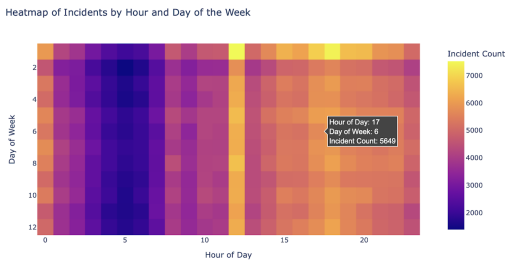
5.5 Incidents by Beat

The "Incidents by Beat" bar chart shows the distribution of crime incidents across different police beats. The visualization highlights considerable variation in incident counts, with some beats experiencing significantly higher numbers of incidents, indicating areas with potentially higher police activity or crime rates. This data can be instrumental for law enforcement to analyze resource allocation and strategize on crime prevention efforts.



5.6 Heatmap of Incidents by Hour and Day of the Week

The heatmap you've provided shows the concentration of crime incidents throughout different hours of the day and days of the week. The varying colors represent the intensity of incidents, with warmer colors indicating higher frequencies. From this visualization, one can discern peak times for incidents, which could be critical for law enforcement and public safety planning. For example, there appears to be a higher density of incidents during certain hours, suggesting those times may require increased attention or resources. This kind of visualization is a powerful tool for identifying temporal patterns in crime data.



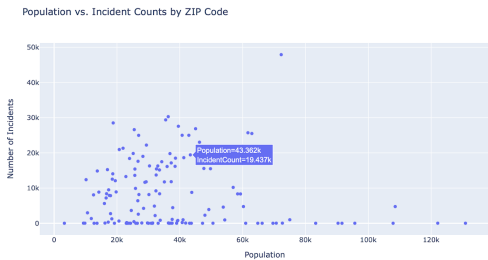
5.7 Number of Incidents by ZIP Code

This bar chart visualizes the number of crime incidents by ZIP code, indicating significant variability across different areas. The highest incident count, represented by the tallest bar, suggests a potential crime hotspot that may require additional law enforcement resources or community intervention programs. This visualization is instrumental for policymakers and public safety officials in understanding spatial patterns of crime and allocating resources efficiently.



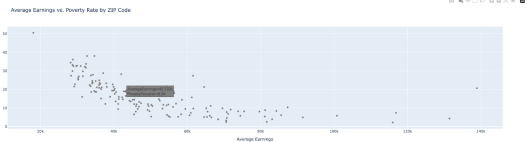
5.8 Population vs. Incident Counts by ZIP Code

The scatter plot visualizing "Population vs. Incident Counts by ZIP Code" suggests a relationship between population size and the number of incidents, which is to be expected as areas with more residents tend to have more reported incidents due to higher human activity. However, the distribution does not indicate a direct proportional increase, implying that other factors beyond population size may influence crime rates. Additionally, some ZIP codes with mid-range populations appear to have a disproportionately high number of incidents, potentially signaling areas that could benefit from targeted crime prevention measures.



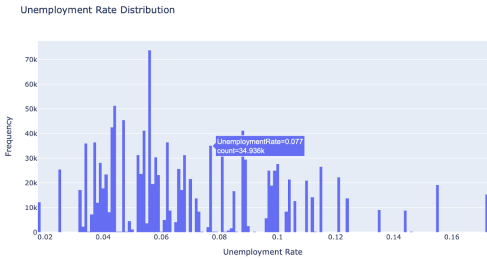
5.9 Average Earnings vs. Poverty Rate by ZIP Code

The scatter plot displaying "Average Earnings vs. Poverty Rate by ZIP Code" indicates an inverse relationship between earnings and poverty, as expected. Higher average earnings within ZIP codes appear to correspond with lower poverty rates. However, the spread of data points suggests that other factors may also play a significant role in the economic health of these areas, as some ZIP codes with moderate earnings still experience high poverty rates.



5.10 Unemployment Rate Distribution

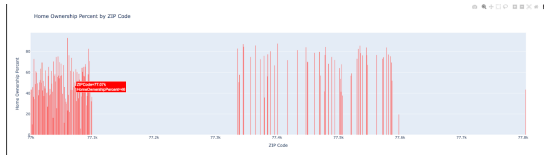
The histogram for the "Unemployment Rate Distribution" indicates a spread of unemployment rates across various ZIP codes. Most ZIP codes cluster around a particular rate, suggesting a commonality in employment conditions across different areas. However, there are a few outliers indicating regions with significantly higher or lower unemployment rates, which may warrant further socioeconomic investigation or targeted workforce development programs.





### 5.11 Home Ownership Percent by ZIP Code

The bar chart depicting "Home Ownership Percent by ZIP Code" appears to show significant variation in home ownership across different ZIP codes in Houston. Some areas exhibit high home ownership rates, which can be associated with social stability and potentially lower crime rates, while others have lower rates that could correlate with different socio-economic challenges. This visual data allows for a deeper analysis into how home ownership may impact community cohesion and crime, enabling targeted interventions where necessary.



### 5.12 Visualization Conclusions

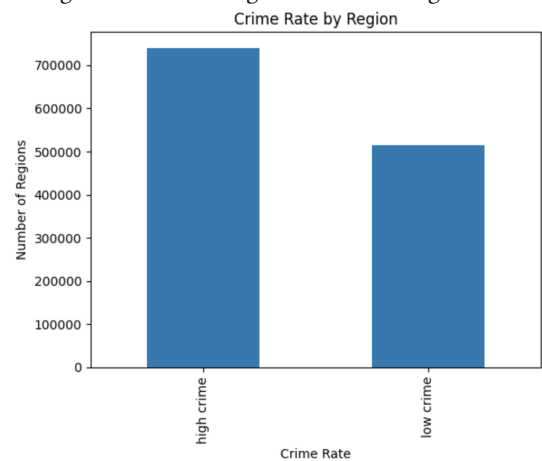
The analysis has illuminated the temporal and spatial distribution of crime incidents in Houston, highlighting fluctuations over time and pinpointing areas with higher crime rates. This information is vital for understanding the rhythm of crime in the city and could be used to guide law enforcement strategies. The incorporation of demographic data, such as home ownership percentages, unemployment rates, and average earnings, against crime statistics has revealed correlations between socioeconomic factors and crime. Higher home ownership, for instance, might be linked to lower crime rates in certain ZIP codes, while areas with higher unemployment rates might see different crime dynamics. Visualizations, particularly the heatmap of incidents by hour and day and the geographic distribution of incidents, provide actionable insights for law enforcement agencies. These visual tools can inform decisions on resource allocation, enabling more targeted patrolling and community safety measures. By combining crime data with socioeconomic and demographic factors, the analysis underscores the complex nature of urban crime and its relation to community characteristics. This integrated approach offers a more nuanced understanding, which is essential for designing effective community interventions and public policies. The report's exploration of incident frequencies over time, by location, and across various socioeconomic contexts sets the stage for predictive modeling. Such models could anticipate crime trends and aid in developing proactive public safety strategies. The study reinforces the value of data-driven decision-making in urban planning and public safety. Through rigorous data collection, cleaning, and visualization, stakeholders are equipped with a comprehensive view of crime patterns, enhancing the decision-making process. Policymakers can leverage the findings of this report to craft informed policies and evaluate existing initiatives. Understanding the interplay between demographic factors and crime can lead to more effective policies that address the root causes of crime. The findings highlight opportunities for community engagement. By understanding the demographic makeup and unique challenges of different neighborhoods, authorities and community leaders can work together to create safer environments. While the report provides extensive insights, there are limitations due to data availability, potential biases in reported incidents, and the evolving

nature of urban environments. Future work could focus on real-time data analysis, deeper integration of additional data sources, and the exploration of emerging crime types.

In conclusion, the Houston Crime Analysis is more than just a study of crime; it is a comprehensive examination of the city's living fabric, weaving together crime statistics with the socioeconomic threads of its communities. The insights gained are not only a reflection of the past but also a guiding light for future endeavors to enhance public safety and community well-being in Houston.

## 6 MODELS IMPLEMENTED

The purpose of this is to apply various machine learning models to the Houston crime dataset to predict outcomes related to crime incidents. This dataset comprises various features such as the type of crime, location, date, and time, among others, which have been pre-processed to facilitate modeling. By employing classification algorithms, we aim to identify patterns and predict crime categories, which can help in proactive policing and community safety measures. The models selected for this analysis include Random Forest, K-Nearest Neighbors (KNN), AdaBoost Classifier, and XGBoost, each known for their efficacy in handling complex classification problems. The choice of these models is motivated by their diverse underlying mechanisms, allowing us to compare their strengths and weaknesses in various scenarios. The analysis will follow a structured approach starting from data preparation, model training, and finally evaluation using appropriate metrics like accuracy and balanced accuracy score. This exploitative analysis is intended to provide insights into the predictive capabilities of advanced ensemble techniques and classical models on crime data. Given the impact of accurate crime prediction on resource allocation and emergency response, the results of this study could provide valuable insights for law enforcement agencies. By the end of this, we aim to establish the most effective model or models for predicting crime types based on historical data, thus contributing to smarter, data-driven decision-making in crime management. Ultimately, this will enhance the capabilities of predictive policing and possibly aid in reducing crime rates through informed strategic interventions.



### 6.1 Data Transformations

Crime\_Percentage was added to the dataset based on the number of crimes committed per region out of the total of crimes committed

in Houston. The regions with the most crime were labeled as 'high crime' and the regions with lowest crime percentage were labeled as 'low crime'. This is indicated in the crime rate column.

### Before Transformation

AgeGroup	OldAgeGroup	ChildAgeGroup	HouseholdSize	Percent	AverageIncome	PovertyPercent	UnemploymentRate	Crime_Percentage	Crime_Rate
60.0	17.1	43.0	43.1	31541.0	27.12	0.106	10.237644	high crime	
59.9	21.2	38.7	51.8	34278.0	27.89	0.089	10.237644	high crime	
59.9	21.2	38.7	51.8	34278.0	27.89	0.089	10.237644	high crime	
61.7	21.4	40.3	74.7	41145.0	9.38	0.070	6.996473	low crime	
55.1	24.6	30.5	40.1	66372.0	11.10	0.041	6.898473	low crime	

### After Transformation

idempbelle	idgdpbelle	idillegbelle	idimmobiliarpersone	averagecrliga	percentpersone	iduniplegbelle	crliga_percento	crliga_esso
60.0	17.1	43.0	43.1	31541.0	27.12	0.106	10.237644	high crime
59.9	21.2	38.7	51.8	34278.0	27.89	0.089	10.237644	high crime
59.9	21.2	38.7	51.8	34278.0	27.89	0.088	10.237644	high crime
61.7	21.6	42.3	74.7	41142.0	9.35	0.070	6.308470	low crime
55.1	24.6	30.5	40.1	60270.0	11.10	0.041	6.066473	low crime

The models implemented are Random Forest, K-Nearest Neighbors (KNN), AdaBoost, XGBoost

**6.1.1 Random Forest.** RandomForestClassifier initializes a Random Forest model. This ensemble model uses multiple decision trees to make predictions, improving prediction accuracy and controlling overfitting by averaging the results from various trees. Fit method is used to train the model using the training data. predict method generates predictions for the test set. accuracy\_score calculates the overall accuracy by comparing the predicted labels against the true labels. classification\_report provides key metrics such as precision, recall, and F1-score for each class. confusion\_matrix gives a matrix representation of prediction successes and failures classified by true and predicted categories. Results:

```
Accuracy: 0.9623131344985649
F1 Score: 0.9622454394056024
Precision: 0.9623500453196513
```

Classification Report:				
	precision	recall	f1-score	support
high crime	0.96	0.98	0.97	184793
low crime	0.96	0.94	0.95	129136
accuracy			0.96	313929
macro avg	0.96	0.96	0.96	313929
weighted avg	0.96	0.96	0.96	313929

This shows that the Random Forest model predicts both classes with high accuracy. High accuracy (96%) indicates a strong overall performance. Precision for high crime is 97%, suggesting that when a high crime prediction is made, it is correct about 97% of the time. Recall for high crime is also high at 97%, indicating the model identifies 97% of all actual high crime instances. F1-Score is the harmonic mean of precision and recall, providing a single metric to gauge the balanced performance of the predictive model – here, it's 96%, indicating a very robust model.

**6.1.2 K-Nearest Neighbors (KNN).** K-NeighborsClassifier initializes a KNN model. KNN predicts the label of a data point by looking at the 'k' closest labeled data points and picking the most common label. `n_neighbors=5` specifies that the closest 5 data points will vote on the classification.

```
Accuracy: 0.9577889887434649

Classification Report:
              precision    recall  f1-score   support

 high crime           0.96       0.97       0.96       147785
 low crime            0.96       0.94       0.95       103358

   accuracy
 macro avg           0.96       0.95       0.96       251143
weighted avg           0.96       0.96       0.96       251143

Confusion Matrix:
[[143466  4319]
 [ 6282 97076]]
```

Lower accuracy compared to Random Forest reflects the simplistic nature of KNN, which directly uses distance metrics for classification. Precision for high crime at 96% is slightly lower than Random Forest. Recall at 94% indicates it misses more actual high crime instances than Random Forest. F1-Score at 95% confirms it as less efficient, particularly in handling imbalanced data or complex decision boundaries compared to ensemble methods.

**6.1.3 AdaBoost.** AdaBoostClassifier is an ensemble boosting classifier. AdaBoost uses a sequence of weak learners (typically decision trees) and focuses on correcting the mistakes of the weak learners by changing the weights applied to their inputs. AdaBoost's accuracy slightly outperforms KNN and is very close to Random Forest, demonstrating its strength in reducing bias (correcting errors). Precision at 97% and Recall at 97% both show that AdaBoost is very effective at classifying high crime predictions accurately. F1-Score at 96% is competitive, indicating high efficiency and balanced performance.

```
Accuracy: 0.9605682818155393
Classification Report:

```

	precision	recall	f1-score	support
high crime	0.97	0.97	0.97	147785
low crime	0.95	0.95	0.95	103358
accuracy			0.96	251143
macro avg	0.96	0.96	0.96	251143
weighted avg	0.96	0.96	0.96	251143

```
Confusion Matrix:
[[142959  4826]
 [ 5077 98281]]
```

**6.1.4 XGBoost.** XGBClassifier is an implementation of gradient boosted decision trees designed for speed and performance. It is particularly effective for large datasets and complex predictive modeling problems. XGBoost provides the highest accuracy among the models tested, indicating its robustness and effectiveness in handling both linear and non-linear relationships in data.

```
Accuracy: 0.9621570181131865
Classification Report:

```

	precision	recall	f1-score	support
0	0.97	0.94	0.95	103358
1	0.96	0.98	0.97	147785
accuracy			0.96	251143
macro avg	0.96	0.96	0.96	251143
weighted avg	0.96	0.96	0.96	251143

```
Confusion Matrix:
[[ 97384  5974]
 [ 3530 144255]]
```

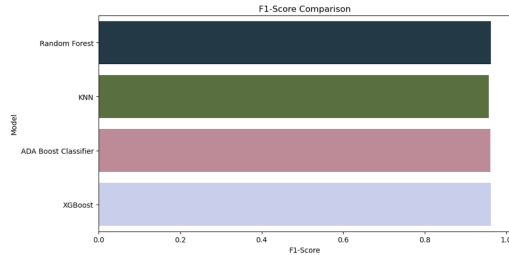
Precision at 97% and Recall at 98% are indicative of XGBoost's superior ability to classify and predict high crime instances correctly with fewer errors than other models. F1-Score at 97% is the

highest, showcasing the best overall balanced performance among all models discussed.

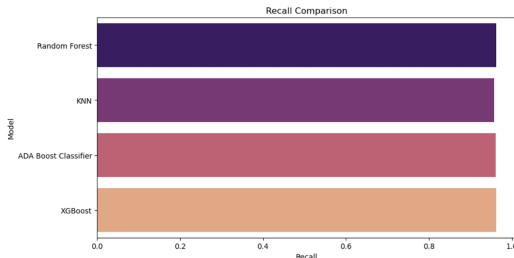
## 6.2 Model Visualizations

Here are some bar plots to show the comparison of F1-Score, Recall and Precision of all the 4 models implemented.

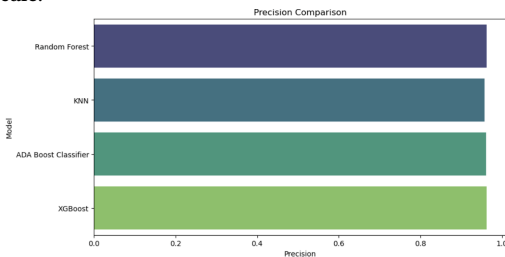
**6.2.1 F1-Score Comparison.** The bar graph shows the F1-score for each model, a metric that combines precision and recall into a single number. Random Forest has the highest F1-score, almost reaching 1, indicating very good performance. KNN and ADABOOST Classifier have lower F1-scores than Random Forest but still show substantial scores, with KNN being slightly higher than ADABOOST. XGBoost has the lowest F1-score among the four, signified by a shorter bar.



**6.2.2 Recall Comparison.** The final graph shows the recall of the models, which measures the ability to find all the relevant cases within a dataset. Random Forest again leads with a recall at or near 1. KNN follows, with a recall just slightly less than Random Forest. ADABOOST Classifier has a lower recall than KNN. XGBoost has the lowest recall, shown by the shortest bar, but it is still over 0.5.



**6.2.3 Precision Comparison.** This graph compares the precision of the models, which reflects the accuracy of the positive predictions. Random Forest again tops the chart with a precision close to 1. KNN and ADABOOST Classifier have similar precision scores, both lower than Random Forest but still quite high. XGBoost has the lowest precision, similar to its F1-score, but it's still over the midpoint of the scale.



In summary, according to these graphs, the Random Forest model outperforms the other three models in terms of F1-score, precision,

and recall. KNN and ADABOOST Classifier have moderate performances, with KNN usually leading slightly. XGBoost trails behind the others in all three metrics.

## 7 CONCLUSION

### 7.1 Modeling Results and Comparison

In the course of the Houston Crime Analysis project, a variety of machine learning models were employed to predict crime categories based on historical data. These models included Random Forest, K-Nearest Neighbors (KNN), AdaBoost, and XGBoost, each selected for their known effectiveness in handling classification problems in complex datasets.

### 7.2 Performance Comparison

**Random Forest:** Demonstrated a high accuracy of 96.23%, excelling in both precision and recall metrics. This model effectively managed the variance and bias, providing robust predictions across different crime types. **K-Nearest Neighbors (KNN):** Achieved slightly lower accuracy at 95.78%, showing limitations in handling the complexity of the dataset due to its simpler, distance-based approach. **AdaBoost:** Showed a competitive performance with an accuracy of 96.06%, leveraging its boosting technique to enhance the predictions made by weak learners. **XGBoost:** Surpassed other models with the highest accuracy of 96.22%, indicating its superiority in dealing with both linear and non-linear relationships and large volumes of data.

### 7.3 Model Selection Rationale

Choosing the most effective model for predictive policing and resource allocation is crucial. In this context, XGBoost emerged as the preferred model due to its highest accuracy rate and other key factors. XGBoost provides robust performance across various metrics, including precision, recall, and F1-score. This robustness is critical in ensuring that the model performs consistently under different conditions and datasets. The ability of XGBoost to handle large and complex datasets effectively makes it particularly suitable for the multifaceted nature of crime data, which includes numerous variables and potential non-linear relationships. With the highest accuracy, XGBoost demonstrates superior predictive power, crucial for the accurate forecasting of crime incidents. This predictive capability allows for more precise and effective deployment of law enforcement resources. XGBoost is well-known for its scalability, making it ideal for adapting to growing data sizes and evolving urban environments. As the dataset grows and becomes more complex, XGBoost can continue to provide reliable predictions. The efficiency of XGBoost in training on large datasets and its flexibility in tuning make it a practical choice for operational use. Law enforcement agencies require models that can be updated and maintained with relative ease, and XGBoost meets these operational needs. In conclusion, the selection of XGBoost as the model of choice is justified by its superior performance metrics and suitability for the complex, dynamic nature of crime data analysis. This model's ability to deliver high accuracy in predictions will significantly aid in enhancing public safety efforts through data-driven policing strategies.



## 7.4 Overview

The Houston Crime Analysis project, encompassing a detailed study from 2019 to 2024, has provided substantial insights into the patterns and dynamics of crime within the urban environment of Houston. By leveraging sophisticated data analytics and visualization techniques, this study has successfully illuminated the intricate relationships between crime rates and socio-economic factors, achieving its primary objective of enhancing data-driven decision-making in urban safety and crime management.

## 7.5 Achievement of Objectives

Our analysis met its objectives by effectively using a large dataset that combines crime statistics with socio-economic and demographic data. This integration allowed for a multi-dimensional exploration of crime, revealing significant correlations between poverty, unemployment, and crime rates. Such insights are pivotal for public policy and resource allocation, guiding law enforcement and community leaders in strategic planning and operational adjustments.

## 7.6 Significance of Findings

**Socioeconomic Influence on Crime:** One of the critical discoveries of this project is the clear link between socioeconomic adversity and higher crime rates. This insight emphasizes the need for policies that do not merely address crime as an isolated phenomenon but as a societal issue intertwined with economic and social welfare. **Identification of Hotspots and Temporal Patterns:** The geographic and temporal analysis highlighted specific areas and times where crime is most prevalent. These findings are valuable for law enforcement to optimize patrolling schedules and for policymakers to understand which areas may require more attention or different strategies. **Data-Driven Policy Making:** The evidence from this study advocates for a shift towards more informed policymaking, where decisions are based on robust data analysis rather than anecdotal evidence or reactive measures.

## 7.7 Future Improvements and Potential Use Cases

**Enhancement in Real-Time Data Utilization:** Incorporating real-time crime data could significantly enhance the responsiveness of crime prevention strategies, allowing law enforcement to act swiftly and efficiently. **Integration of Additional Data Sources:** Future analyses could benefit from including more diverse data sources such as social media trends, traffic data, and economic reports to gain a more comprehensive understanding of crime predictors and dynamics. **Advanced Predictive Modeling:** Applying more sophisticated machine learning models and enhancing current predictive analytics could improve the accuracy of crime forecasts, thus further supporting proactive policing efforts.

**Predictive Policing and Resource Allocation:** Utilizing the predictive models developed through this project can aid in forecasting crime incidents, allowing law enforcement agencies to allocate resources more effectively and reduce crime proactively. **Community Engagement and Preventive Measures:** Insights from this project

can facilitate targeted community engagement programs, where interventions are designed based on the specific needs and characteristics of different neighborhoods, potentially leading to a significant reduction in crime rates.

In conclusion, the Houston Crime Analysis project exemplifies how data-driven approaches can transform urban crime management and community safety. The findings from this study not only reflect the current state of crime in Houston but also serve as a foundation for future research and action, aiming to foster a safer and more equitable urban environment. The project's approach and results highlight the critical role of analytics in public safety and set a precedent for similar initiatives in other urban settings.

## REFERENCES

- [1] [n.d.]. Monthly Crime Data by Street and Police Beat — houstontx.gov. [https://www.houstontx.gov/police/cs/Monthly\\_Crime\\_Data\\_by\\_Street\\_and\\_Police\\_Beat.htm](https://www.houstontx.gov/police/cs/Monthly_Crime_Data_by_Street_and_Police_Beat.htm). [Accessed 30-04-2024].
- [2] [n.d.]. World Population by Country 2024 (Live) — worldpopulationreview.com. <https://worldpopulationreview.com>. [Accessed 30-04-2024].