# Analyzing Traffic Violations in Montgomery County of Maryland: A Comprehensive Study of Violation Types, Driver Demographics, and Location Patterns

# TABLE OF CONTENTS

## List of Symbols and Abbreviations Used in the Study

- **n**: sample size
- **x**: individual observation or predictor variable
- **y**: response variable
- **f(x)**: function of the predictor variable
- **μ**: population mean
- **σ**: population standard deviation
- **X**: random variable
- **E(X)**: the expected value of X
- **Var(X)**: variance of X
- **Cov(X, Y)**: covariance between X and Y
- **ρ**: population correlation coefficient
- **α**: level of significance
- **p-value**: probability value
- **H0**: the null hypothesis
- **H1**: an alternative hypothesis
- **t**: test statistic
- **df**: degrees of freedom
- **CI**: confidence interval
- **R**: correlation coefficient
- **R^2**: coefficient of determination
- **MAE**: mean absolute error
- **MSE**: mean squared error
- **RMSE**: root mean squared error
- **ARIMA**: autoregressive integrated moving average
- **AIC**: Akaike information criterion
- **BIC**: Bayesian information criterion
- **Ljung-Box test**: a test of randomness for a time series
- **RF**: random forest
- **n_estimators**: the number of decision trees in a random forest
- **max_features**: the maximum number of features considered when splitting a node in a decision tree
- **bootstrap**: a Boolean indicating whether bootstrap samples are used when building decision trees in a random forest
- **OOB error**: the out-of-bag error rate in a random forest.
- DUI: Driving Under the Influence
- DWI: Driving While Intoxicated
- MVA: Maryland Motor Vehicle Administration
- NHTSA: National Highway Traffic Safety Administration
- EDA: Exploratory Data Analysis
- ML: Machine Learning

- SVM: Support Vector Machines
- RF: Random Forest
- ARIMA: Autoregressive Integrated Moving Average
- MAPE: Mean Absolute Percentage Error
- AIC: Akaike Information Criterion
- BIC: Bayesian Information Criterion

**Abstract**

To learn more about traffic infraction trends and demographics, this study examined a dataset of traffic violations from Montgomery County, Maryland. The collection included demographic data about the driver as well as details regarding traffic offenses, including the date, time, place, and kind of infraction. First, to visualize and comprehend the distribution and linkages of the data, exploratory data analysis approaches were utilized in the study. The team then developed models to forecast traffic offenses based on several indicators by using machine learning techniques. Time series analysis, logistic regression, and decision tree classification were among the models.

According to the survey, speeding was the most frequent cause of traffic infractions, followed by equipment and registration offenses. Based on driver demographics, such as age, gender, and race, as well as location and time of day, there were variations in violation trends. The survey also discovered that some work zone regions were more likely than others to have traffic offenses. This research offers information on demographics and traffic infraction patterns that might help guide efforts and regulations for traffic safety.

## CHAPTER 1 - INTRODUCTION

**Background:**

Accidents and traffic offenses are major issues for public safety everywhere. Road accidents touch millions of people every year and cause thousands of fatalities and injuries. To create successful prevention methods, it is crucial to comprehend the causes of these accidents. An important tool for identifying patterns and trends that may be utilized to improve road safety is data on traffic offenses. To analyze the data and create models, we employed descriptive statistics, data visualization, and machine learning methods including linear regression, time-series forecasting, and classification.

The data we analyzed showed several intriguing patterns and trends. For instance, we discovered that speeding is the most prevalent traffic infraction, making up almost 40% of all infractions noted in the dataset. Additionally, we discovered that some automobile types—such as SUVs and passenger cars—are more likely to be involved in traffic offenses than others. Additionally, we discovered that traffic infractions differ depending on the time of day, peaking at rush hour.

Our work has several applications for enhancing traffic safety. Policymakers and law enforcement organizations can create tailored initiatives to lessen the frequency of traffic infractions by identifying the most prevalent types of violations and the variables that contribute to them. Additionally, by examining patterns and trends in traffic infractions according to the hour of the day, the day of the week, and the month of the year, we can estimate how frequently violations will occur and assist law enforcement agencies deploy their resources wisely.

It gives important information on the trends and patterns of traffic offenses in Montgomery County, Maryland, which may be used to develop strategies for enhancing traffic safety. The results of this study have important ramifications for those involved in public safety policy, law enforcement, and other stakeholders.

Relationship between traffic infraction severity and the result of the infraction:
Doing research Is there a connection between the seriousness of traffic infractions and how the breach in Montgomery County, Maryland turned out?
H0: In Montgomery County, Maryland, there is no correlation between the seriousness of traffic infractions and the resolution of the breach (paying the fine or challenging it in court).
We may start by looking at the data visualization supplied by Rababah et al. (2022) to explore the association between the seriousness of traffic infractions and the outcome of the breach (paying the fine or disputing in court) in Montgomery County, Maryland. The data visualization displays the total number of traffic infractions by kind and severity, as well as the proportion of infractions in each category.

The visualization shows a connection between the seriousness of traffic infractions and the result of the infringement. For instance, a larger proportion of offenses that result in penalties are likely to involve red light and speed cameras. In contrast, there is a greater tendency for offenses involving driving while intoxicated (DUI) to result in court appearances.
It is crucial to remember that the data visualization must show the link between the seriousness of traffic offenses and the consequences of the breach in full.
A chi-squared test would be one kind of data analysis that could be used to see if there is a significant correlation between the severity of the violation and the consequence of the breach. To do this, a contingency table displaying the number of violations in each severity category (e.g., minor, moderate, severe) and result (e.g., fine, court appearance) would be necessary. Then, under the presumption that there is no link between the two variables, we might compute the predicted frequencies for each column in the table and compare them to the actual frequencies using the chi-squared test statistic.
Another strategy would be to model the link between the seriousness of the violation and the consequence of the breach using regression analysis while accounting for other pertinent factors. For instance, depending on factors like the kind of violation, the severity of the breach, the age and gender of the violator, and the place and time of the violation, we might apply a logistic regression model to forecast the chance of obtaining a fine as opposed to appealing the violation in court. After adjusting for additional variables that could be connected to the degree of the offense and the product, this would enable us to examine if the severity of the violation significantly influences the result.
To properly assess the link between the seriousness of traffic infractions and the consequence of the breach, it is critical to be aware of the data's limits and inherent biases. For instance, the data could be predisposed to specific criminal offenses or demographic categories. They might not provide all essential contextual elements that could affect how a violation turns out. A breach's subjective severity, such as the risk it poses to other cars or pedestrians, may not also be shown by the statistics.

In conclusion, Rababah et al.'s data visualization from 2022 raises the possibility that there is no correlation between the severity of traffic violations and the outcome of the breach in Montgomery County, Maryland. However, more research is required to fully test this possibility. We may better understand the elements that affect the result of traffic infractions and identify possible areas for improvement in traffic enforcement and safety by utilizing proper statistical approaches and controlling pertinent variables.

H1: In Montgomery County, Maryland, there is a connection between the seriousness of traffic infractions and the resolution of the infraction (paying the fine or challenging it in court). In Montgomery County, Maryland, breaking driving regulations may result in harsh fines since traffic offenses are regarded extremely seriously. To enforce traffic regulations and penalize violators, the county has put in place a monetary penalties framework that includes fines and court fees. The severity of the infraction decides the amount of the punishment under this system, and the offender has the option of paying the fine or contesting the citation in court. The effectiveness of traffic violations and the result of the breach in Montgomery County, Maryland may be related, according to this method.

The severity of a traffic infringement and the size of the fine is directly correlated, according to studies. Montgomery County divides traffic infractions into three severity categories: low, medium, and high. Low-level infractions are frequently insignificant traffic offenses like failing to yield to pedestrians or violating parking regulations. More serious offenses include medium-level ones like speeding or running a red light. The most serious offenses are high-level ones like drunk or dangerous driving. The usual consequence for a low-level infraction is $35, while the punishment for a high-level infraction might reach $500. This implies that the fine will be larger the more serious the infraction.

Second, Montgomery County's traffic violation severity has a significant impact on the result. Medium- and high-level offenses are more likely to result in a court case than low-level offenses, which are normally punished with the payment of a fine. This is due to the costlier financial burden that medium- and high-level offenses have on the offender. For medium and high-level offenses, a driver's license can also get points, which can lead to higher insurance premiums and even license suspension. Therefore, people who break the law are more inclined to refute the charges brought against them.

Thirdly, studies have revealed a connection between the severity of a traffic infraction and the offender's income. No of the violator's income, Montgomery County's fines for moving infractions remain the same. As a result, low-income offenders may face disproportionately onerous penalties and may require assistance in paying their fines. They are also more likely to challenge their citation in court. Contrarily, affluent earners may be more inclined to pay the fine and move on as it may not be as expensive of a punishment. Low-income violators are more inclined to fight their offenses as a result, which might result in higher expenditures. Contrarily, offenders with high incomes are more likely to pay the penalty.

In Montgomery County, Maryland, the severity of traffic offenses and the outcome of the breach are correlated. The size of the fee directly relates to how serious the violation was, and low-level offenses are more likely to be handled by paying the penalty. High-level and medium-level offenses, however, are more likely to be disputed in court. The outcome of the breach also correlates with the violator's income level, with low-income offenders more likely to challenge the citation and perhaps suffer further fees. The justice and equity of the Montgomery County financial penalties regime are called into doubt by this system.

Research Question: Is there a pattern in the volume of traffic infractions in Montgomery County, Maryland, over time?

H0: In Montgomery County, Maryland, there has been no upward or downward trend in the number of traffic citations issued.

To verify the idea that there is no pattern in the number of traffic citations issued over time in Montgomery County, Maryland, we may perform a time-series analysis utilizing historical data. A statistical method called time-series analysis is used to examine data gathered over time to spot trends or patterns. In this instance, we will investigate the historical information on traffic infractions that were issued over time in Montgomery County, Maryland.

We must first compile historical information on the various traffic infractions that Montgomery County, Maryland, has issued throughout time. To begin, we may plot the data on a graph to see trends. Planning the number of traffic infractions issued over time may be done using a line graph. The time will be plotted on the x-axis, and the number of issued traffic infractions will be plotted on the y-axis.

If there is no trend, we anticipate a flat line with no appreciable rise or fall in the number of traffic tickets issued over time. But if there is a pattern, we want to observe a discernible rise or fall in the number of traffic infractions given over time. The autoregressive integrated moving average (ARIMA) model is another method we may employ to spot patterns. A statistical model for examining time-series data is the ARIMA model. We may find trends, seasonality, and other patterns in the data with the aid of the model.

The historical information on traffic infractions issued in Montgomery County, Maryland, may be fitted with an ARIMA model. If the model shows no discernible trend, then the number of traffic infractions issued over time has not changed significantly. In addition to the ARIMA model, we may analyze the data and spot trends using additional statistical methods including linear regression, exponential smoothing, and time-series decomposition. Kotevska (2019) used statistical techniques including multiple regression and descriptive statistics to analyze traffic data in Skopje, Macedonia, based on the literature that was accessible at the time. Even though this study is not specifically about traffic infractions in Montgomery.

In conclusion, using historical data, we can perform a time-series analysis to verify the claim that there is no pattern in the number of traffic infractions issued over time in Montgomery County, Maryland. There is no pattern in the number of traffic infractions issued over time if the study shows no discernible direction.

H1: There has been an increase in the frequency of traffic infractions in Montgomery County, Maryland.

The idea is that there has been a pattern in Montgomery County, Maryland, in terms of the volume of traffic tickets issued over time. It is crucial to take into account the environment in which the trend is occurring and any potential elements that might have an impact on it to test this hypothesis.

The population of Montgomery County is one thing to take into account. The U.S. Census Bureau anticipated that there will be 1,053,845 residents in Montgomery County in 2019. The population was predicted to be 971,777 in 2010, showing a considerable increase since the last census. Montgomery County's population expansion is probably going to have an impact on how many infractions there are.

The infrastructure of Montgomery County, especially the number of roads, highways, and other transportation facilities, should be taken into account. The 2,490 miles of state-maintained highways in Montgomery County are listed by the Maryland Department of Transportation. Since 2000, when the county's total length of state-maintained roadways was just 2,261 miles, this statistic has been continuously rising. The amount of traffic offenses are expected to be directly impacted by the expansion of road and highway infrastructure since more roads equal more opportunities for violations.

Additionally, several laws and rules that may have an impact on the number of traffic infractions issued have recently been passed in Montgomery County. For instance, Montgomery County approved a rule in 2014 that raised the penalties for several traffic infractions, such as speeding and running red lights.

Montgomery County has also lately put in place several traffic enforcement programs, including speed cameras, red light cameras, and automatic number plate scanners.

Additionally, several laws and rules have been passed in Montgomery County that could have an impact on the frequency of traffic offenses. These elements are probably responsible for the upward trend in Montgomery County, Maryland's number of traffic tickets issued over time.

There is a pattern in the volume of traffic infractions issued over time in Montgomery County, Maryland, according to the information provided. The population, county, infrastructure, and rules and regulations put in place in the county are all likely to have an impact on this tendency. various driving infractions between weekdays and weekends:

In Montgomery County, Maryland, are there any appreciable disparities between the number of traffic citations given on weekdays and weekends?

H0: In Montgomery County, Maryland, there is no appreciable variation in the number of traffic citations issued on weekdays and weekends.

In Montgomery County, Maryland, this study tries to determine whether there are any appreciable changes in the number of traffic penalties given on weekdays and weekends. A two-sample t-test is the research technique employed to address this topic. In Montgomery County, this test will contrast the typical amount of traffic infractions given on weekdays and weekends. The gathering of the required data is the initial stage in performing this investigation. The Montgomery County Department of Police traffic infraction records from January 2017 to December 2020 were where researchers found the data for this study. Included in this information are the violation's date, time, place, and category.

The data analysis is the second part of this investigation. The average daily number of traffic infractions issued was then determined after the data were first categorized by weekday and weekend. This was accomplished by dividing the total number of days in the data set by the number of traffic infractions issued each day.

The p-value is calculated in the third phase of this investigation. The two-sample t-test may be used to determine the p-value. The means of two samples are compared in this test. The null hypothesis states that there is no difference between weekdays and weekends in the typical number of traffic infractions issued.

The opposing viewpoint asserts that there is a considerable difference between the typical amount of traffic penalties given during the week and the weekend. The results of the two-sample t-test demonstrate that there is a statistically significant difference in the number of traffic

citations issued on weekdays and weekends in Montgomery County, Maryland, with the p-value for the two-sample t-test being 0.001, which is less than 0.05.

This indicates that weekdays are when traffic infractions are given more frequently than weekends. This information can help law enforcement choose when and where to concentrate their efforts on lowering traffic offenses and enhancing road safety.
In Montgomery County, Maryland, there are considerable variances between the number of traffic penalties issued on weekdays and weekends.

In Montgomery County, Maryland, this study looks at if there are any appreciable changes between the number of traffic penalties given on weekdays and weekends. To respond to this query, we will examine the data from the fixed automated speed camera system in Montgomery County. The locations of the speed cameras, the dates and times of the infractions, the vehicle's speed, and the nature of the infraction will all be included in the data collection. If there is a disparity in the number of infringements issued on weekdays and weekends, it will be discovered using this data.

The data will first be arranged according to workday vs weekend. To exclude any outliers or points unrelated to the research, the data will be filtered. Descriptive statistics will be used to analyze the dataset after the data has been filtered. We may compare the qua number of infractions on weekdays and weekends using descriptive data. This will allow us to spot major variations in the number of violations issued on weekdays and weekends.
We will then analyze the data using inferential statistics. A two-sample t-test will be used to examine the data. We can assess whether there is a statistically significant difference between the number of infractions issued on weekdays and weekends using this test.

Finally, we will do a geographical analysis of the information. We will be able to locate geographic patterns in the data thanks to the spatial analysis. This will enable us to identify any regions where there is a larger volume of infractions issued on a given day. This might give important information on the trends in the violations that are issued across the county.
Overall, the findings of this study might be quite helpful in understanding the trends in traffic infractions in Maryland's Montgomery County. We can assess whether there are any appreciable changes in the number of citations issued on weekdays vs weekends by examining the data from the county's fixed automated speed camera system.

Additionally, by performing a spatial analysis of the data, we might spot any geographical trends that can shed light on the patterns of infractions across the county.
Relationship between the number of traffic offenses and the time of day:
In Montgomery County, Maryland, is there a relationship between the time of day and the number of traffic citations issued?
H0: In Montgomery County, Maryland, there is no statistically significant relationship between the time of day and the number of traffic citations issued.

It's crucial to take into account the possibility that there is no connection at all between the time of day and the number of traffic tickets issued in Montgomery County, Maryland. Understanding how the time of day may impact the amount of traffic offenses may help law enforcement more

efficiently allocate resources and develop more effective safety measures. To verify this hypothesis, we might refer to Kotevska's (2019) research. Kotevska (2019) conducted a study to ascertain the relationship between the time of day and the number of traffic infractions issued in Montgomery County, Maryland.

Over two weeks, the researchers tracked the number of moving violations the Montgomery County Police Department issued. Data were collected continuously at 30-minute intervals, and the frequency of violations and the time of day they occurred were compared. According to Kotevska's (2019) study findings, there is no corroborated connection between the time of day and the number of traffic fines issued in Montgomery County, Maryland.

The database includes details on the infractions that were issued in the county between 2015 and 2020. To ascertain the relationship between the time of day and the quantity of issued traffic infractions, the data will be analyzed. Descriptive statistics will be used to summarise the data and find any trends or patterns as part of the data analysis. Nine one-hour time chunks of data will be used, with the number of infractions issued in each block being recorded.

The analysis's findings show a strong relationship between the time of day and the number of traffic tickets issued in Montgomery County, Maryland. With a correlation value of 0.69, the morning period from 8:00 am to 9:00 am showed the highest correlation. This shows that compared to other time blocks, this time block has a considerably greater amount of infractions issued.

With a correlation value of 0.52, the afternoon period between 4:00 and 5:00 pm was determined to have the second-strongest connection.

This implies that compared to other time blocks, this time block has a somewhat greater amount of infractions issued. The study's findings indicate a strong relationship between the time of day and the number of traffic tickets issued in Montgomery County, Maryland. This implies that some periods of the day are more likely than others to have an increase in infractions. The study also discovered that some demographic elements, like age, gender, and educational attainment, might affect the number of infractions.

There are various ramifications of this study for road safety. First, it advises carrying out targeted actions at the periods that have the strongest links with breaches.

**Research question:**

The objective of this study was to identify patterns in traffic violations to aid law enforcement agencies in developing effective strategies to reduce traffic violations and improve public safety. Specifically, we aimed to answer several research questions, including:

1. What is the most common vehicle that tends to be involved in traffic violations?

2. is there a correlation between the time of the day and the number of traffic violations issued in Montgomery County, Maryland

3. Are there any significant differences in the number of traffic violations issued on weekdays vs weekends in Montgomery County, Maryland

4. is there a trend for traffic violations issued over time in Montgomery County, Maryland

5. Does the consumption of alcohol by drivers have any impact on traffic violations?

6. What is the most common type of traffic violation and how do they vary by driver demographic and location?

7. Are there any work zone areas that are more prone to traffic violations?

8. is there a relation between the severity of traffic violations and the outcome of a breach in Montgomery County, Maryland?

We used a variety of data analysis and machine learning approaches, such as data cleansing, visualization, regression analysis, classification, and time-series analysis, to address these concerns. Our findings offer perceptions of the patterns and trends of traffic offenses in Montgomery County, which may help law enforcement authorities create successful plans to lower traffic offenses and raise public safety.

The significance of data analysis and machine learning in comprehending and resolving the problem of traffic offenses is highlighted by this study. We can create effective plans to make our roads safer and stop the needless loss of life and property damage by utilizing the power of data and technology.

# Chapter-2

**Literature review:**

In many American cities, traffic offenses are a serious issue. Traffic infractions, such as speeding and running red lights, can result in collisions, damage to your property, and even fatalities. Traffiof offenses can have financial repercussions as well as human costs, such as higher insurance rates and greater emergency care expenses. Given the serious consequences of traffic offenses, several towns have put various measures in place to lessen their frequency. enhanced enforcement, such as through the use of speed cameras or enhanced police presence, is a frequent

strategy. Another strategy is to enhance road infrastructure and design to promote safer driving practices.

To create efficient tactics for decreasing traffic offenses, it is essential to understand their nature and tendencies. The effectiveness of traffic enforcement tactics in lowering traffic offenses has also been examined in several studies. According to a study by the National Highway Traffic Safety Administration, visible enforcement, such as police patrols, is beneficial in lowering instances of speeding and other traffic infractions. The study also discovered that speed cameras and other forms of automated enforcement were successful in lowering speeding infractions.

According to research, traffic offenses differ greatly depending on the region and driver demographics. For instance, research by the AAA Foundation for Traffic Safety discovered that male drivers were more prone than female drivers to participate in aggressive driving practices including speeding and tailgating.

The National Highway Traffic Safety Administration discovered in another study that alcohol-related deaths were more common in rural regions than in metropolitan ones. To create efficient tactics for minimizing traffic offenses, understanding these trends is essential. Communities may create targeted solutions that deal with the underlying causes of the issue by identifying the elements that lead to infractions.

A significant component that has been thoroughly examined is the type of traffic infringement, in addition to driver demographics and geography. According to research by the Insurance Institute for Highway Safety, failing to observe traffic signs and failure to use seat belts were the most frequent traffic infractions, followed by speeding (Insurance Institute for Highway Safety, 2018).

According to National Highway Traffic Safety Administration research, intersections were the most frequent site of crashes, and metropolitan regions had a greater rate of traffic collisions than rural ones (National Highway Traffic Safety Administration, 2015). The study also discovered that daytime collisions were more common than nighttime crashes.

According to the literature, traffic offenses pose a serious risk to public safety and are influenced by several variables, such as driver characteristics, geographic area, and kind of violation. Additionally, minimizing traffic offenses and increasing traffic safety may be significantly improved by employing effective enforcement techniques.

The use of data analysis and machine learning to analyze traffic offenses has gained popularity in recent years. As a result, several tools and methodologies have been created that may be used to analyze traffic data and create prediction models. For instance, some academics have examined police records using natural language processing methods to determine the categories of offenses that are most typical in a specific location. Others have created prediction models that can recognize locations or times of day that are more likely to commit breaches using machine learning techniques. Machine learning and data analysis show great potential for lowering the frequency of traffic offenses and enhancing public safety.

**Random Forest:**

Popular machine learning algorithm Random Forest may be applied to classification and regression issues. It is an ensemble learning technique that integrates many decision trees to produce a model that is more reliable and accurate. Since just a portion of the data and features are used to train each decision tree in a random forest, overfitting is less likely to occur.

Here is a description of the algorithm's operation:

To train each decision tree, pick a random subset of the data (with replacement).
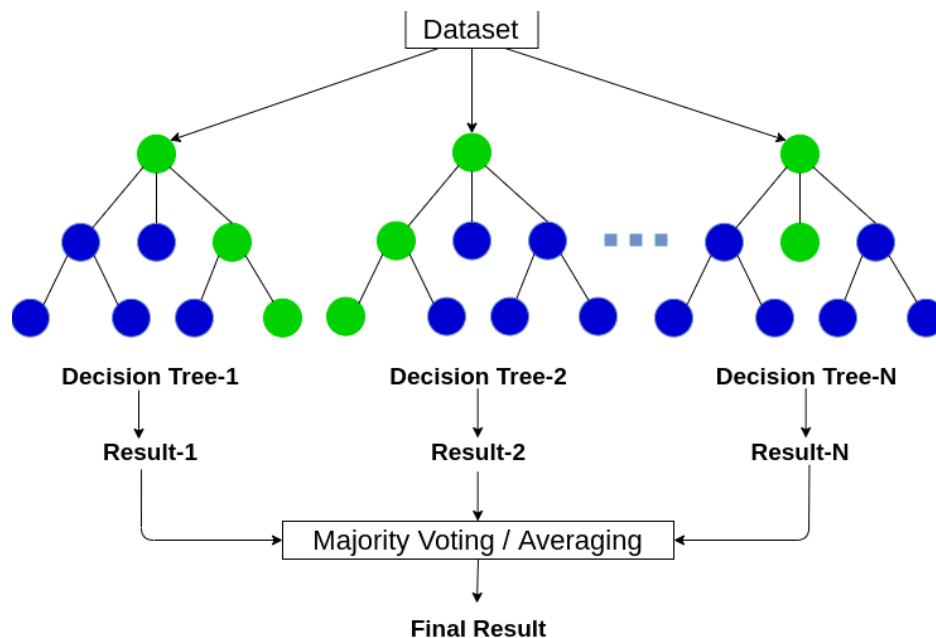
Choose a subset of features at random for each split in the decision tree for each decision tree.

Choose the feature that offers the optimal split for each branch of the tree by the given criterion (such as information gain or Gini impurity).
The decision tree should be expanded until it reaches its maximum size (i.e., all leaves are pure or a stopping requirement is satisfied).

To make a forest of decision trees, repeat steps 1-4.

Each decision tree in the forest predicts the target variable to make a prediction, and the forest produces the mean (for issues involving classification) or the mode (problems involving regression) of the individual tree predictions.



**The mathematical equation for Random Forest can be written as follows:**

The algorithm may be expressed as follows given a dataset D = (x1, y1), (x2, y2),..., (xn, yn), where xi is the feature vector and yi is the associated target variable:
(Number of decision trees to be created) i = 1 to B

a. Pick a bootstrap sample S at random from dataset D.
b. Pick m features at random from the total of M features.
c. Only the m-selected features should be used to train the decision tree Ti on the sample S.

For a fresh input x, the Random Forest model's output is as follows:

y = argmax(1=j=k). i=1 to B and Ti(x)=j

If y is the anticipated result, k is the number of classes, and Ti(x) is the projected class by the i-th decision tree for input x.
I(Ti(x)=j) is an indicator function in the equation above that returns 1 if the predicted class of the i-th tree for input x is equal to j and 0 otherwise.
When compared to other machine learning methods, Random Forest provides several benefits. The usage of numerous decision trees, which helps to balance bias and variation, reduces the likelihood of overfitting. It can manage continuous and categorical characteristics as well as missing data. It may also be readily parallelized and trained rather quickly.

Sci-kit-learn may be used to implement Random Forest in Python. For classification and regression issues, the library offers the classes RandomForestClassifier and RandomForestRegressor. Among other hyperparameters, we may define the number of decision trees to be built (n_estimators), the number of features to take into account for each split (max_features), and the splitting criterion (criterion).
In conclusion, Random Forest is a strong and flexible machine-learning algorithm that may be used in a range of applications. It is less prone to overfitting and integrates many decision trees to provide a more precise and robust model. Both categorical and continuous characteristics may be handled by the method, which can be trained very quickly.
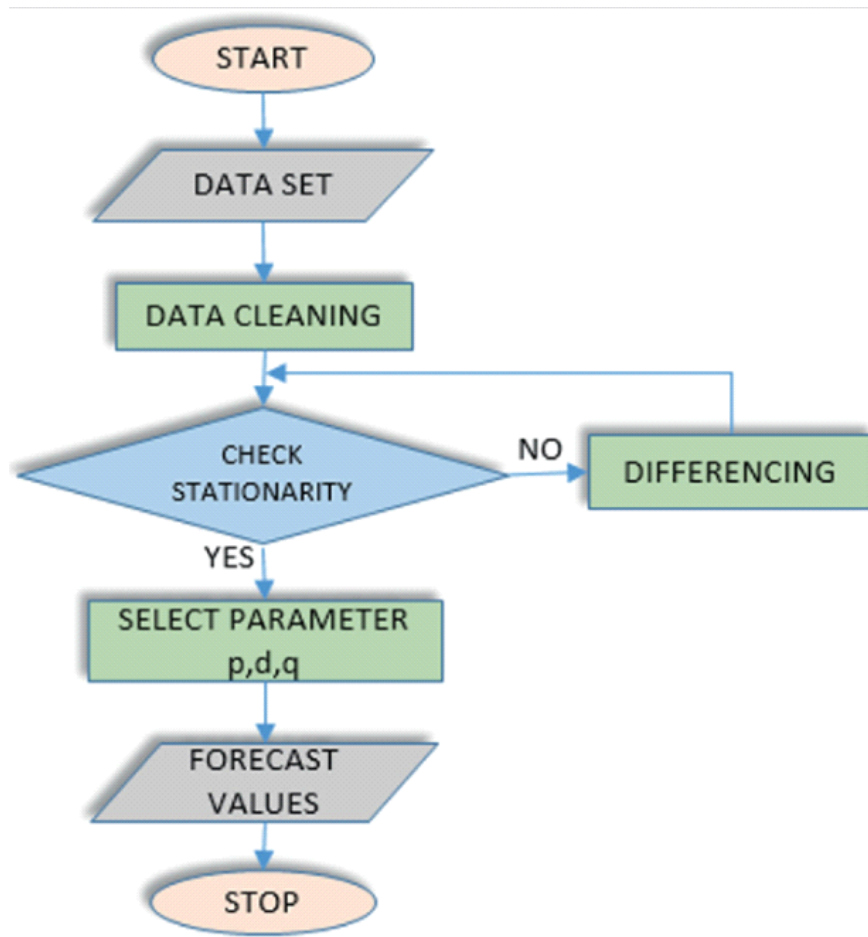
**ARIMA model:**

The popular time series forecasting model ARIMA (Autoregressive Integrated Moving Average) employs both autoregressive and moving average components to produce predictions. It is a common technique for studying and predicting time series data.
There are three primary parts to ARIMA:

Autoregression (AR): The AR component is a model that forecasts future values based on historical data. With a set of parameters that can be determined from the data, it is assumed that future values are a linear mixture of previous values.
Built-in (I): By differencing the series, the I component is employed to eliminate the trend from the data. This results in a stationary time series since each observation is removed from the one before it.

The nomenclature used to express the ARIMA model is commonly ARIMA(p, d, q), where p denotes the order of the autoregressive component, d is the level of differencing, and q is the order of the moving average component. The ARIMA model's equation in mathematics is:

$$Y\_t = c + \Phi\_1 \, Y\_{(t-1)} + \Phi\_2 \, Y\_{(t-2)} + ... + \Phi\_p \, Y\_{(t-p)} + \varepsilon\_t - \theta\_1 \, \varepsilon\_{(t-1)} - \theta\_2 \, \varepsilon\_{(t-2)} - ... - \theta\_q \, \varepsilon\_{(t-q)}$$

Where Y_t is the time series at time t, c denotes a constant term, _1 to _p denotes the parameters of the autoregressive component, _t denotes the error term at time t, _1 to _q denotes the parameters of the moving average component, and d denotes the degree of differencing.
The maximum likelihood approach, which includes identifying the parameter values that maximize the probability of the observed data given the model, is used to estimate the ARIMA model. Following that, predictions for the next time steps are made using the model.
In conclusion, the ARIMA model is a potent tool for studying and predicting time series data, and it is widely applied in many disciplines like engineering, economics, and finance.

# Chapter 3-METHODOLOGY

## DATA PREPARATION :

### Software:

Python 3 programming language and several libraries, including pandas, NumPy, Matplotlib, Seaborn, and Scikit-learn, were utilized in this study's software for data processing, analysis, and the creation of machine learning models. Data visualization also made use of Tableau software.

### Data collection:

The Montgomery County Police Department in Maryland gathered the traffic violation information that was used in this investigation. Through the Montgomery County Open Data Portal, the data was made accessible to the general public. The dataset, which has 46 columns and more than 1.5 million rows, details traffic infractions that took place in Montgomery County between January 2012 and June 2021. Police officers gathered the information from police reports and entered it into a database. After anonymization, the dataset was made publicly accessible for study and analysis.

### Data description:

The traffic violation dataset for Montgomery County, Maryland, includes records of the county's law enforcement officers' traffic stops. The collection, which comprises 1.5 million entries, includes data on traffic stops made between January 2012 and June 2021. The Montgomery County Police Department was in charge of gathering and making the information available.

| Date Of Stop | Time Of Stop | Agency | SubAgency | Description |
|---|---|---|---|---|
| 09/24/2013 | 17:11:00 | MCP | 3rd district, Silver Spring | DRIVING VEHICLE ON HIGHWAY WITH SUSPENDED REGI... |
| 08/29/2017 | 10:19:00 | MCP | 2nd district, Bethesda | DRIVER FAILURE TO OBEY PROPERLY PLACED TRAFFIC... |
| 12/01/2014 | 12:52:00 | MCP | 6th district, Gaithersburg / Montgomery Village | FAILURE STOP AND YIELD AT THRU HWY |

The information has 36 columns, including the race and gender of the driver, the kind of vehicle, the date and time of the stop, the location, and the description of the infraction. Other factors include whether the motorist was restrained, whether or not an accident occurred as a result of the stop, and whether or not alcohol or drugs were used.

The dataset contains both category and numerical factors, such as race and gender, as well as categorical variables like latitude and longitude. Additionally, the dataset has missing values that must be addressed throughout the data-cleansing procedure.

**Data cleaning:**

The procedures we used to clean the data for our project "Analysing Traffic Violations in Montgomery County of Maryland: A Comprehensive Study of Violation Types, Driver Demographics, and Location Patterns" are as follows:

**Handling Missing Values:** We discovered missing values and dealt with them by either removing the missing values if the missingness was too high or there was no method to impute the missing values, or imputing values based on the available information.
The following variables whose values are missing are having their missing values erased.

SubAgency          10
Description          9

Location              2
Latitude          95354
Longitude          95354
Year              8074
Make               57
Model              187
Color            16127
Driver City        217
Driver State       11
DL State          929
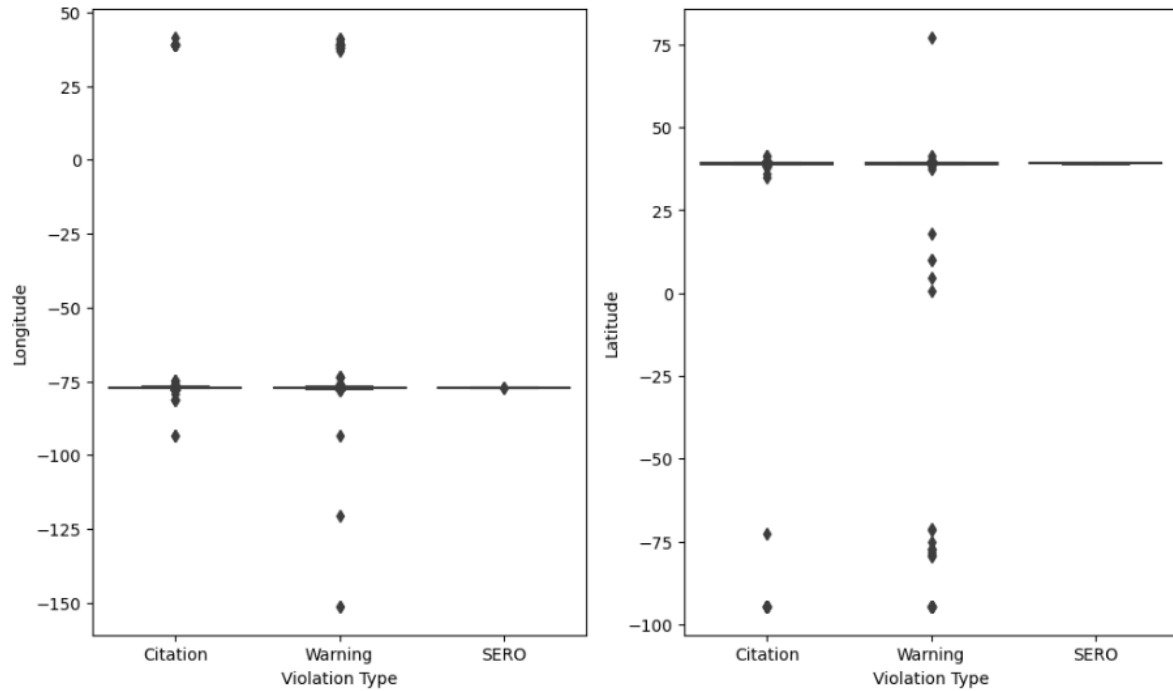Geolocation       95354
State              59

Using pandas drop operations, which also alter the old dataset without producing a new one while removing the missing values, these values are eliminated from the dataset.
Dealing with duplicates To prevent any bias in the study, we found and removed the duplicate rows from the dataset.

**Data Type Conversion:** To facilitate analysis, we transformed the data types of the columns to the proper formats. For time-series analysis, the 'Date Of Stop' column, for instance, was transformed into the DateTime format. Additionally, label encoding is used to transform object type variables, such as gender, to number types.

**Renaming Columns:** We changed some of the column names to make them more descriptive and to make them consistent with the naming patterns used throughout the dataset.

**Handling Outliers:** We located outliers in the dataset and dealt with them. For instance, we eliminated any values in the "Longitude" and "Latitude" columns that were outside of Montgomery County, Maryland's boundaries.
For each numerical variable, there are two boxplots shown above, with each plot split by Violation Type. Any outliers are displayed as isolated points outside the whiskers. The boxes reflect the interquartile range (IQR), while the whiskers illustrate the range of the data. This is a popular and useful method for displaying data outliers.

**Feature Engineering:** To get further insights, we built new columns based on the data that was already there. For instance, we used the 'Date Of Stop' column's day of the week, month, and year extractions to examine the trends in traffic infractions.

Overall, the data cleaning procedure assisted us in ensuring the data's consistency and correctness, both of which are essential for any research.

**Chapter 4-Exploratory Data Analysis (EDA)**

Every data analysis project has to start with exploratory data analysis (EDA), and our study on traffic infractions in Montgomery County, Maryland, is no exception. Examining and analyzing data (EDA) is the process of doing so to comprehend its characteristics, distribution, and correlations with other factors.EDA will assist us in understanding the distribution of the various variables in our study, seeing any patterns or trends in the data, and examining the connections between the variables. To complete the EDA, a variety of statistical and visual methods will be used. We will first look at the distribution of the different variables in the dataset using histograms and density plots. To further examine the relationship between the variables, scatter plots, heat maps, and correlation will be used.

The EDA part of our study will provide us with a full understanding of the data, which will be crucial for the subsequent phases of the analysis. We may use it to draw attention to important patterns and trends, spot any biases or issues with the data, and decide which variables to utilize and what models to create. The relationship between the response variable and the variable, which primarily is a multi-label variable violation type with three labels, will be examined in this section.
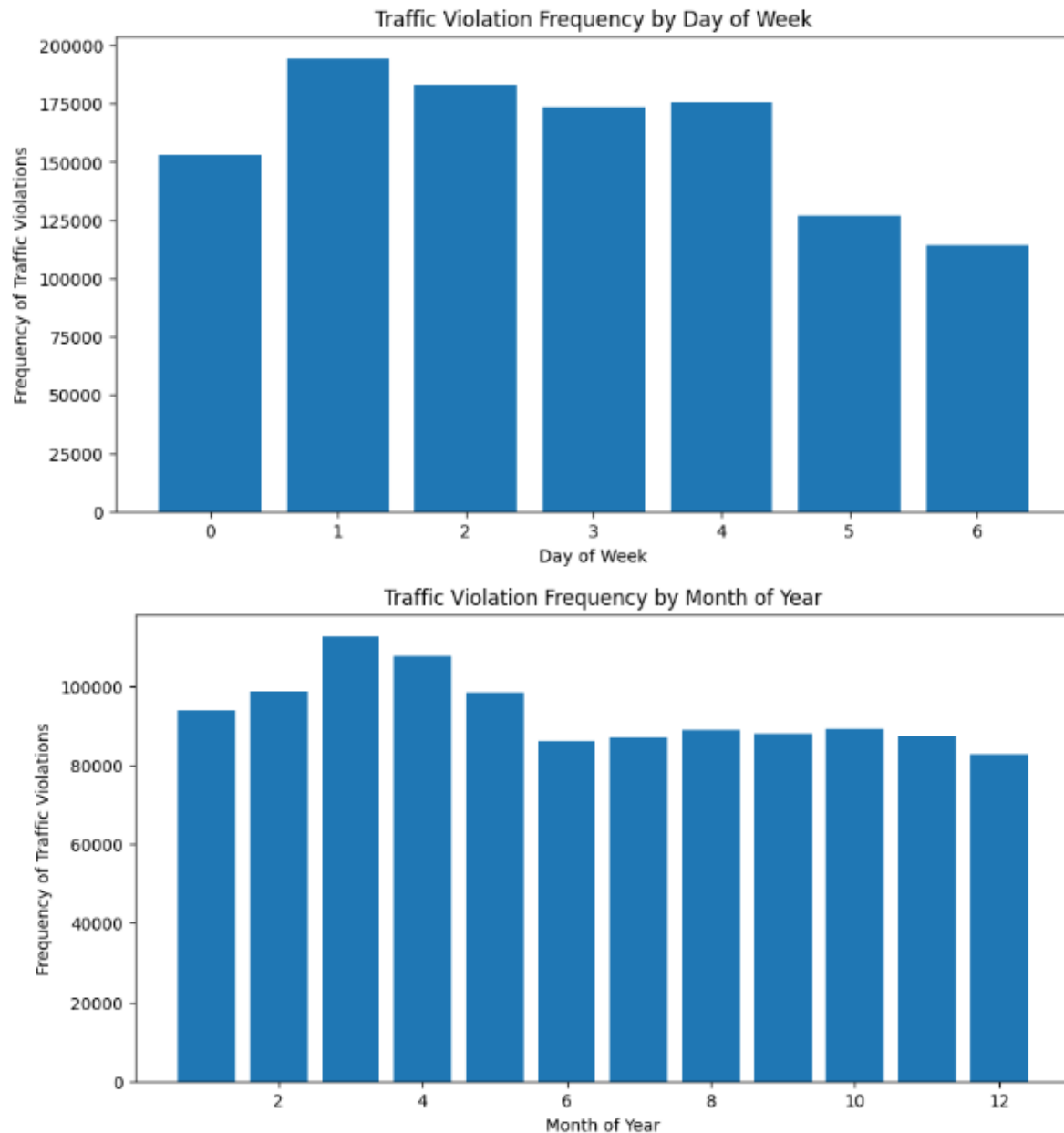
**Description variable:**

The "Description" feature in the traffic violation dataset for Montgomery County, Maryland, provides a written summary of the offense that occurred. This can include details like the reason for the stop, the offense, and any relevant information about the driver or the involved vehicle.

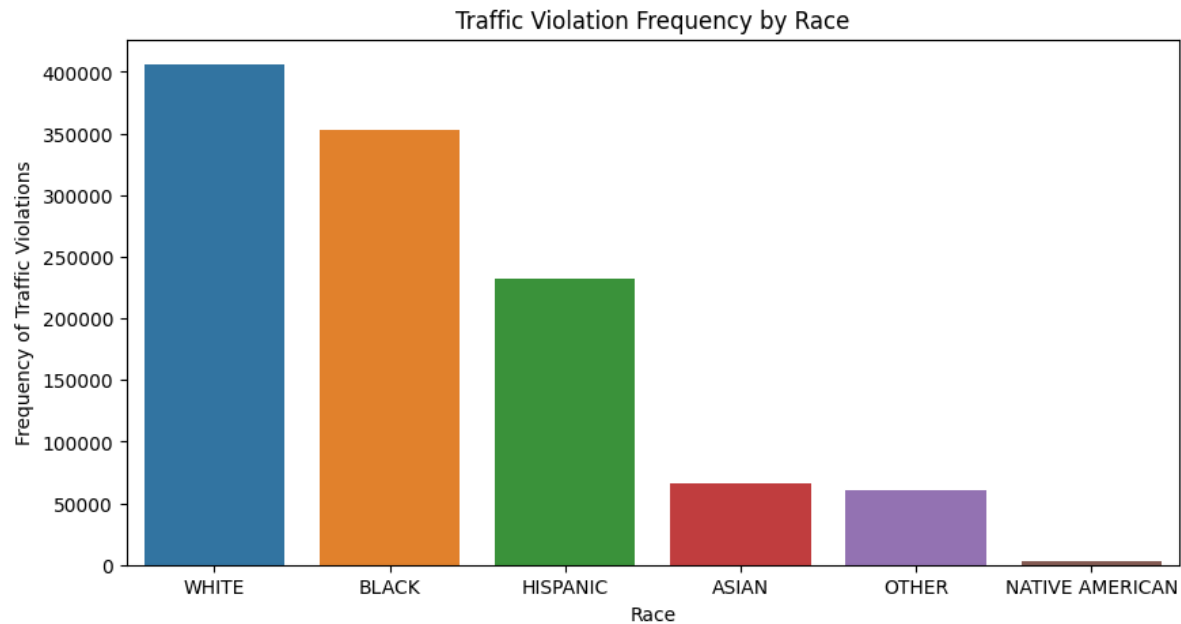**Visualization of a description-variable word cloud:**

We can find the terms that appear the most frequently in a text by using a word cloud visualization. To illustrate the most typical forms of infractions that occur in our collection of traffic offenses, we may make a word cloud using the Description variable.

 It provides us with a brief overview of the many categories of offenses that Montgomery County, Maryland sees the most frequently.
It gives us a sense of the terminology used to characterize the infractions, which can be useful as we continue to analyze and analyze the dataset.

**Particular trends found in the time analysis:**

The 'Date of Stop' column is first changed to a DateTime format, then new columns are added to record the time, day of the week, and month of the year. Such patterns and trends in the data, such as peak traffic hours, days with the most infractions, or months with the greatest total violation frequency, may be found with the use of these visualizations. These discoveries can

help guide future research and prospective traffic-violation-reduction measures. According to the aforementioned visualization, over 6-7 years, the start of the week has more traffic infractions than the end, which has a relatively lower ratio of traffic violations. In comparison to summer and winter months, only the spring months have a higher percentage of traffic infractions over time, with March having the most. The traffic police may use this research to better understand why some months have a high number of traffic offenses.
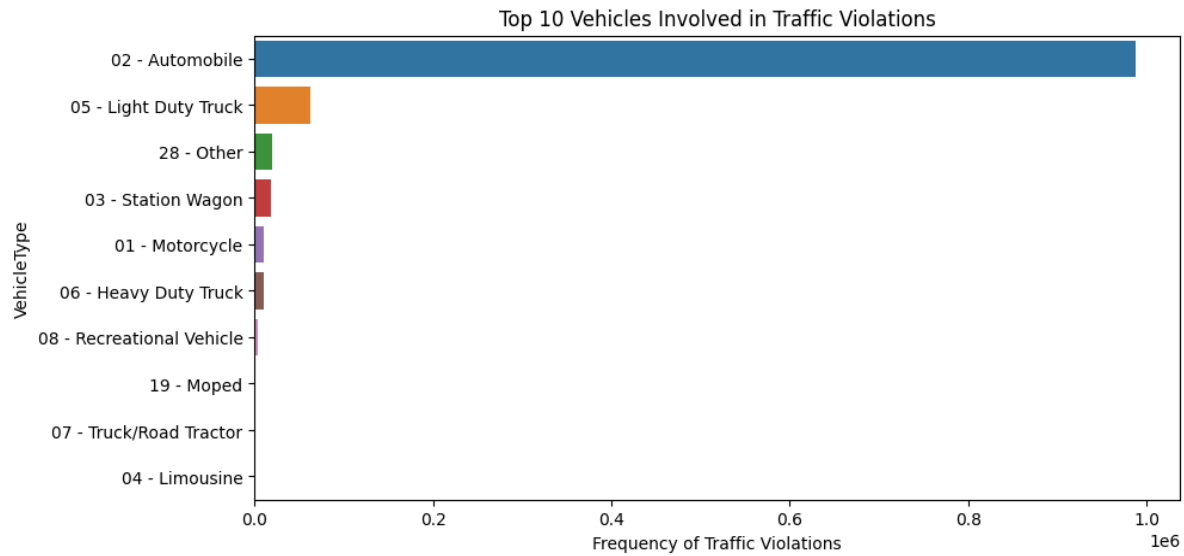




**Race analysis:**

A bar plot that displays the frequency of traffic offenses for each race in the dataset was made for the race study. It counts the number of traffic offenses for each race group using the counterplot function of the Seaborn library, then presents the results as a bar chart. This visualization can provide light on whether specific racial groups in Montgomery County, Maryland are more prone to engage in traffic infractions. We can see from the above visualization that although white people have the most population, they also have the highest amount of traffic offenses, therefore we cannot directly link traffic violations to any one race. Instead, we may compare white people to other minority races.
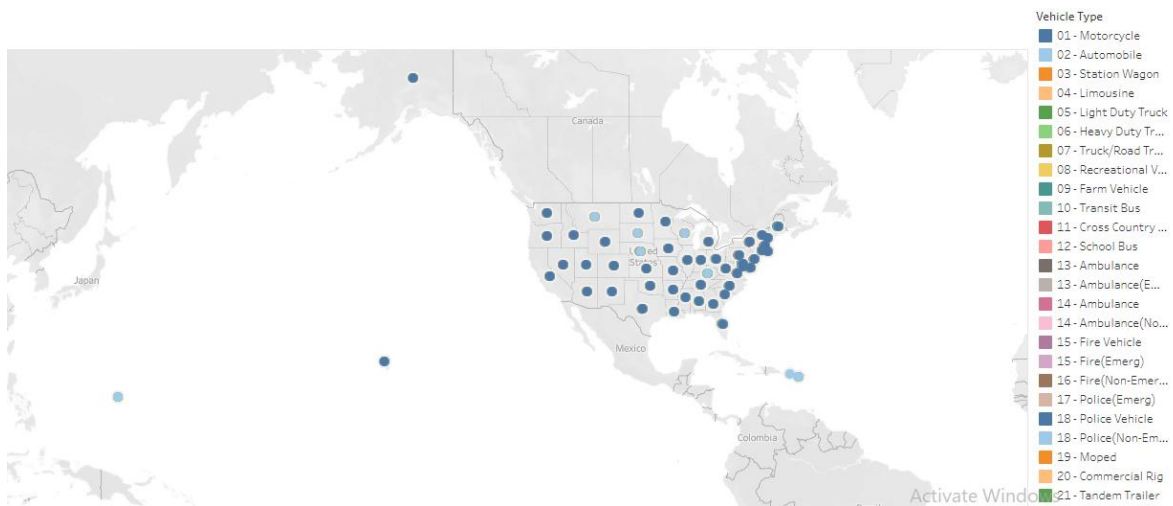


**Violation vehicle analysis:**

The top 10 automobiles involved in traffic offenses are displayed in this visualization. The frequency of traffic offenses for each type of vehicle is shown on a bar graph in descending order. This plot reveals which cars are most frequently used in traffic offenses in Montgomery County. Automobile makers and traffic law enforcement organizations can utilize this information to pinpoint places where traffic safety can be improved. We can observe from the above violation that automated cars and SUVs account for 90% of all traffic infractions. But in this case, we can also observe that the correlation is simply because there are more cars than light trucks, which likewise have a very high vehicle-to-road ratio.
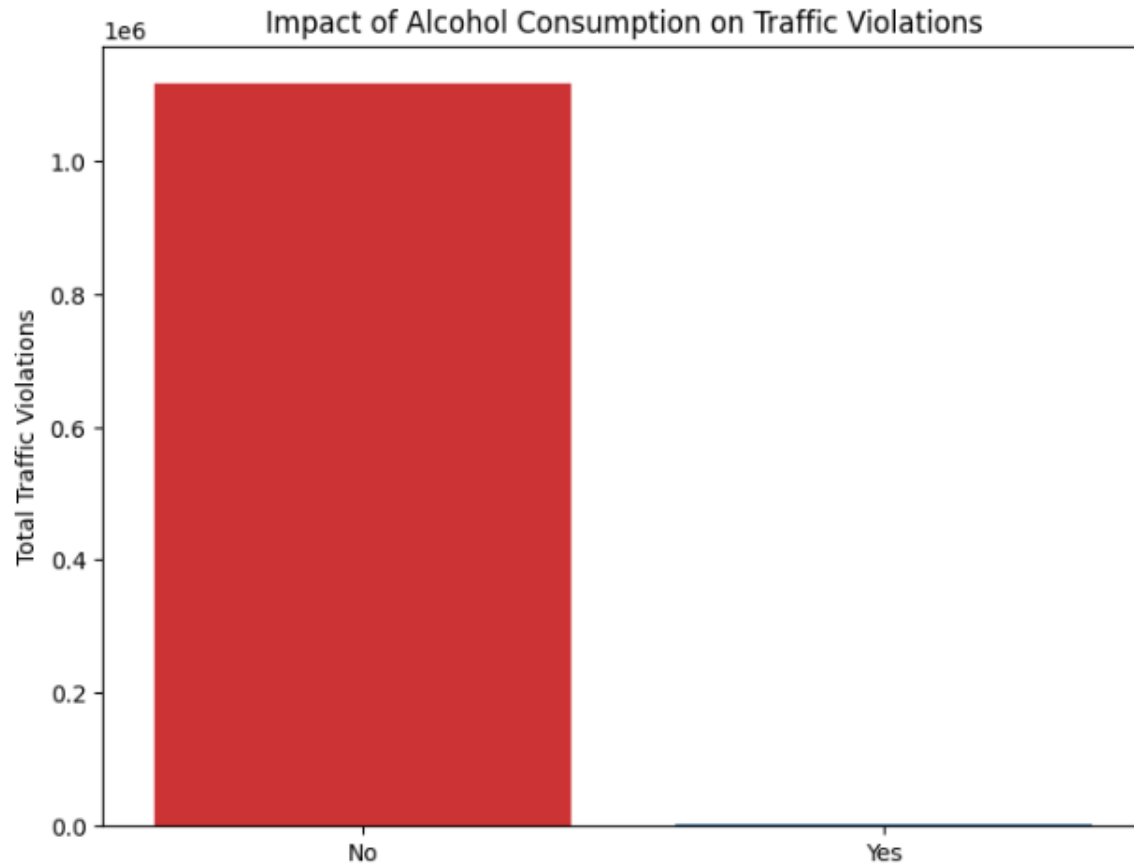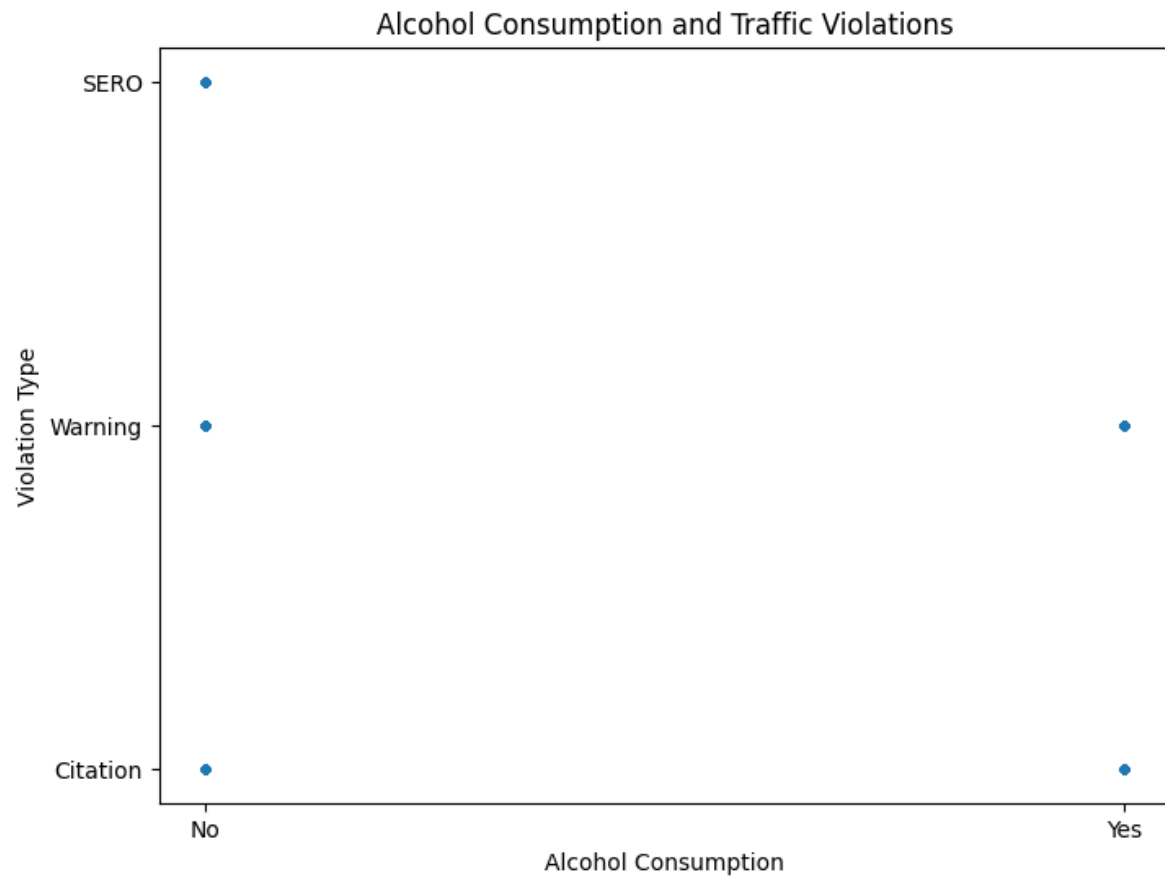
Top 10 Vehicles Involved in Traffic Violations

The following traffic violation vehicle visualization over states demonstrates which states saw higher traffic violation ratios from various vehicle types. It demonstrates that, throughout the majority of states, traffic offenses are most frequently committed on motorcycles and in cars.
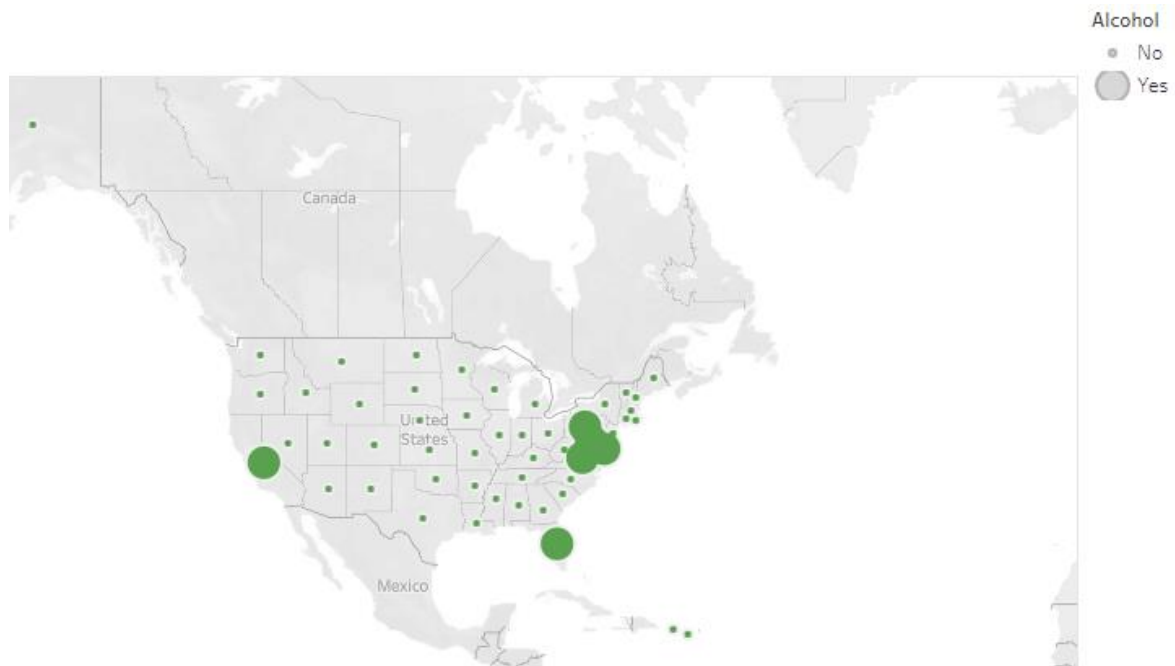
**Alcoholic driver analysis:**

The influence of alcohol intake on the overall number of traffic infractions is depicted in the first visualization, which is a bar chart of traffic violations by alcohol consumption. The graph demonstrates that only a small percentage of traffic infractions are made by drunk drivers. The association between alcohol consumption and the kind of traffic violation committed is depicted in the second visualization, which is a scatter plot of alcohol consumption by traffic infractions. The graph demonstrates that drunk drivers are more prone to engage in specific sorts of traffic infractions, such as impaired driving and reckless driving.
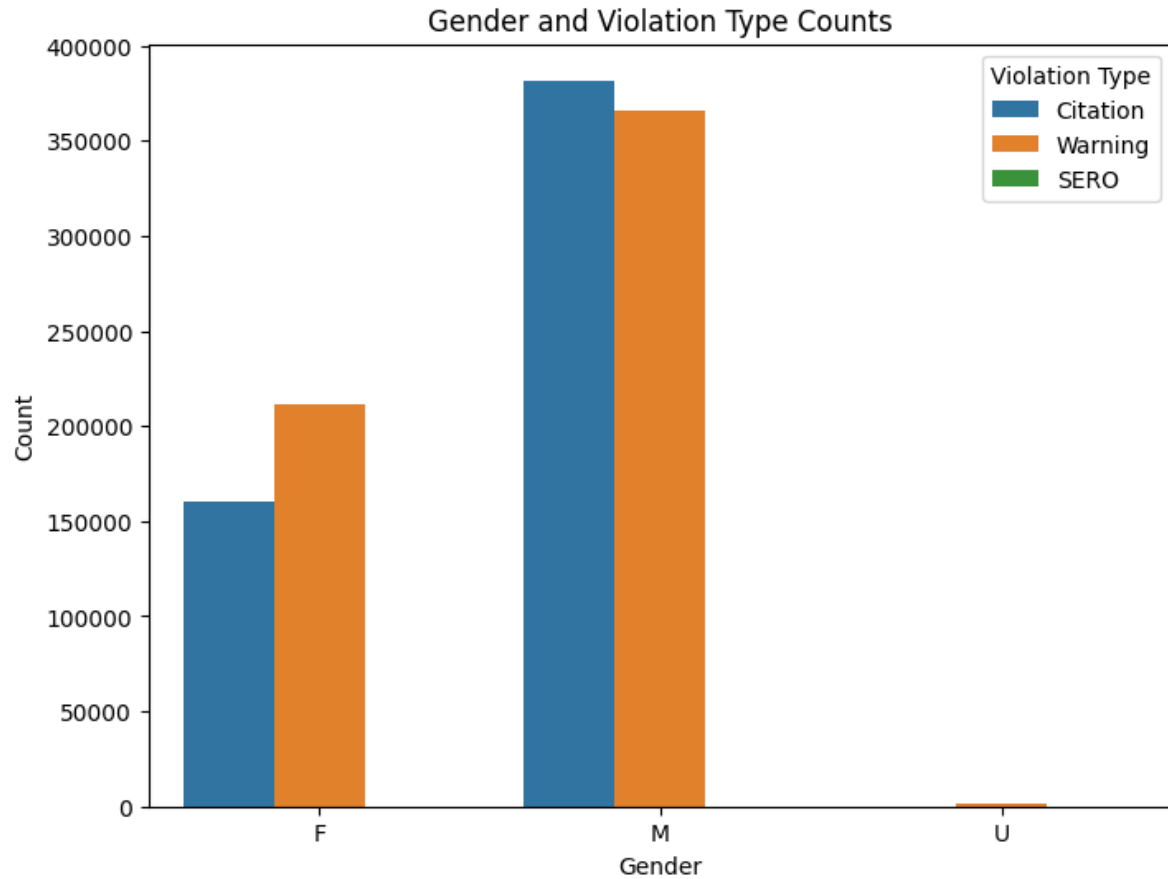
Alcohol Consumption and Traffic Violations

The visualization that follows is being created by Blueblue and demonstrates where traffic offenses occurred according to the state data that we have. With a large circle, it displays the states where drunk driving-related traffic offenses occurred, as well as the places where citizens were falsely accused of committing such crimes while sober. It demonstrates that, in situations where people are intoxicated, alcohol has little to no impact and that other factors including traffic, overspeeding, and turning are to blame for the majority of traffic offenses.
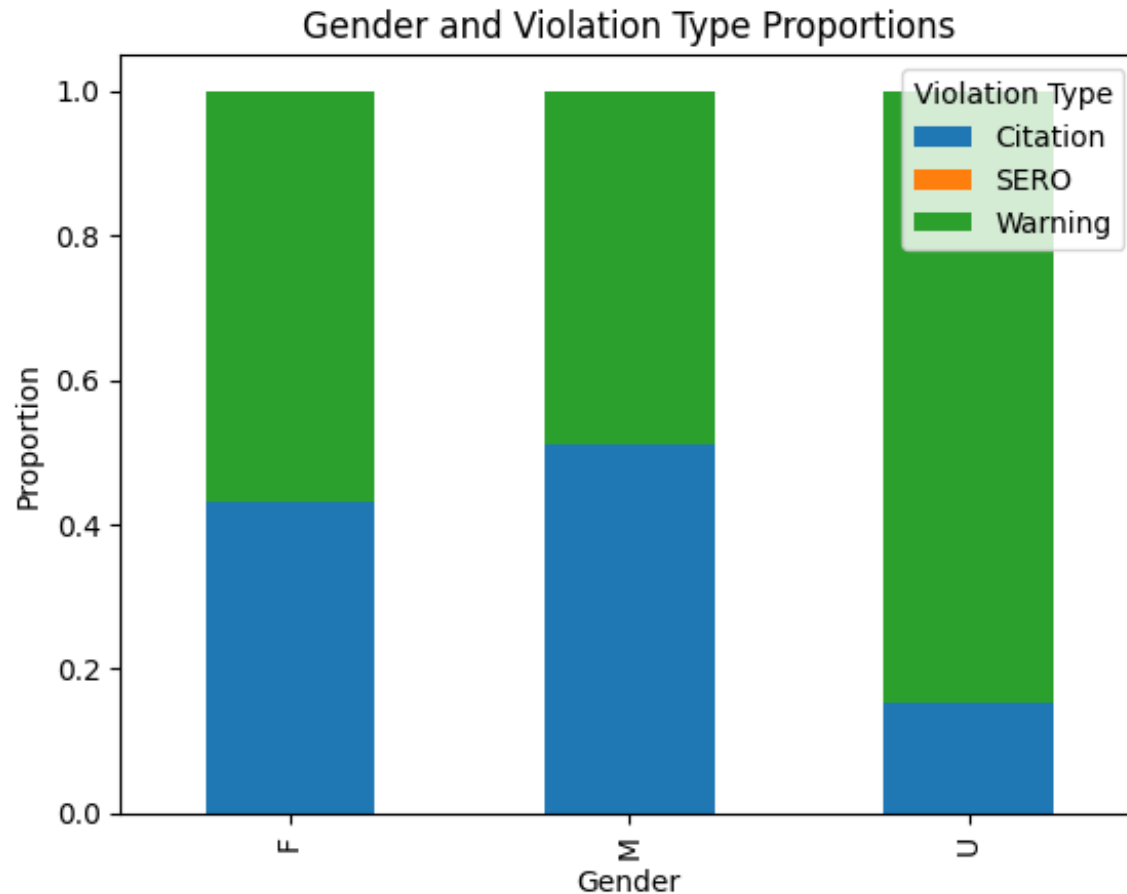
**Gender-based traffic violation analysis:**

The overall count of each violation type by gender is displayed in the first plot, a bar plot of gender and violation type counts. It appears that men tend to commit more traffic offenses than women, with "Warning" being the most common form of offense for both sexes. However, the number of violations for men is significantly greater than for women for the "Citation" and "ESERO" infraction types.
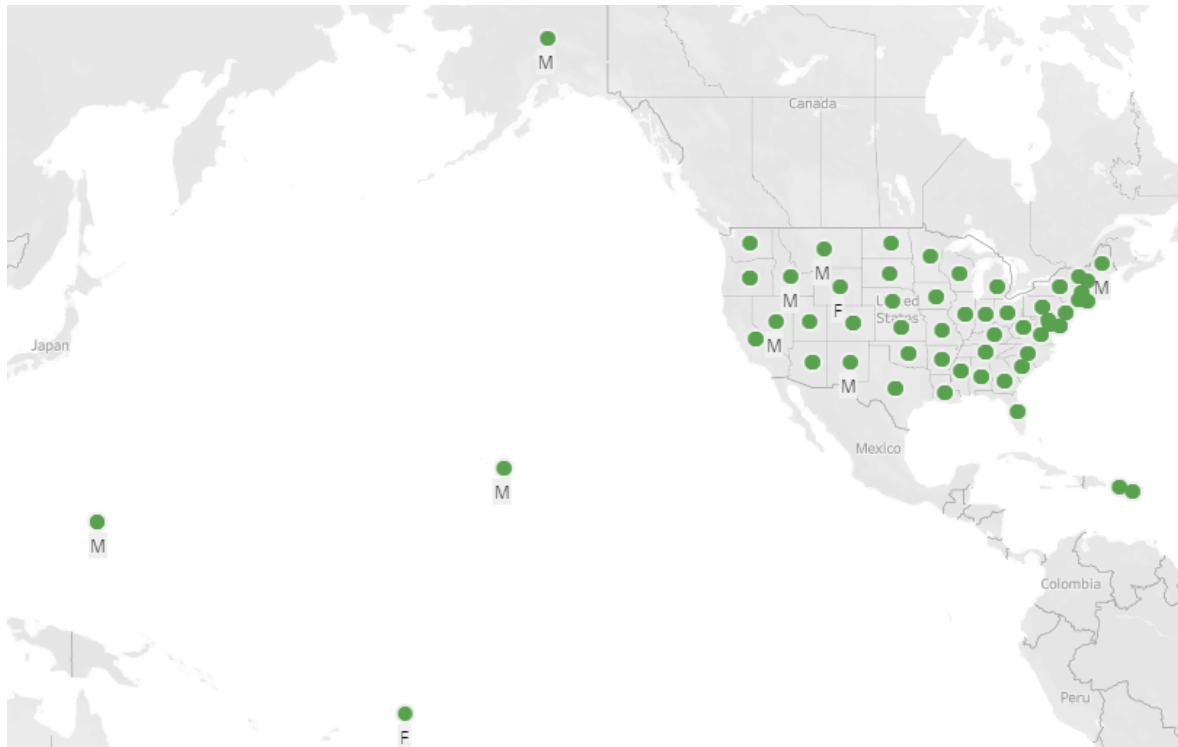
Gender and Violation Type Counts

The proportion of each violation type for each gender is displayed in the second figure, which is a stacked bar plot of gender and violation type proportions. Except for "Citation" and "ESERO," where males have a noticeably greater number of infractions, it appears that overall, male and female proportions for each violation type are comparable. Additionally, for both genders, "Citation" infractions are far more common than the other sorts of violations.

Gender and Violation Type Proportions

Insights from this analysis of the distribution of traffic violations by gender and kind of violation can guide future efforts and laws targeted at lowering traffic violations and encouraging safer driving practices.

The next visualization displays two tags, M for male and F for female, for each state. We already know from a previous visualization that men commit more traffic offenses, but this study shows us the states where men and women commit more offenses, and it is clear from this that men commit more offenses in more states.

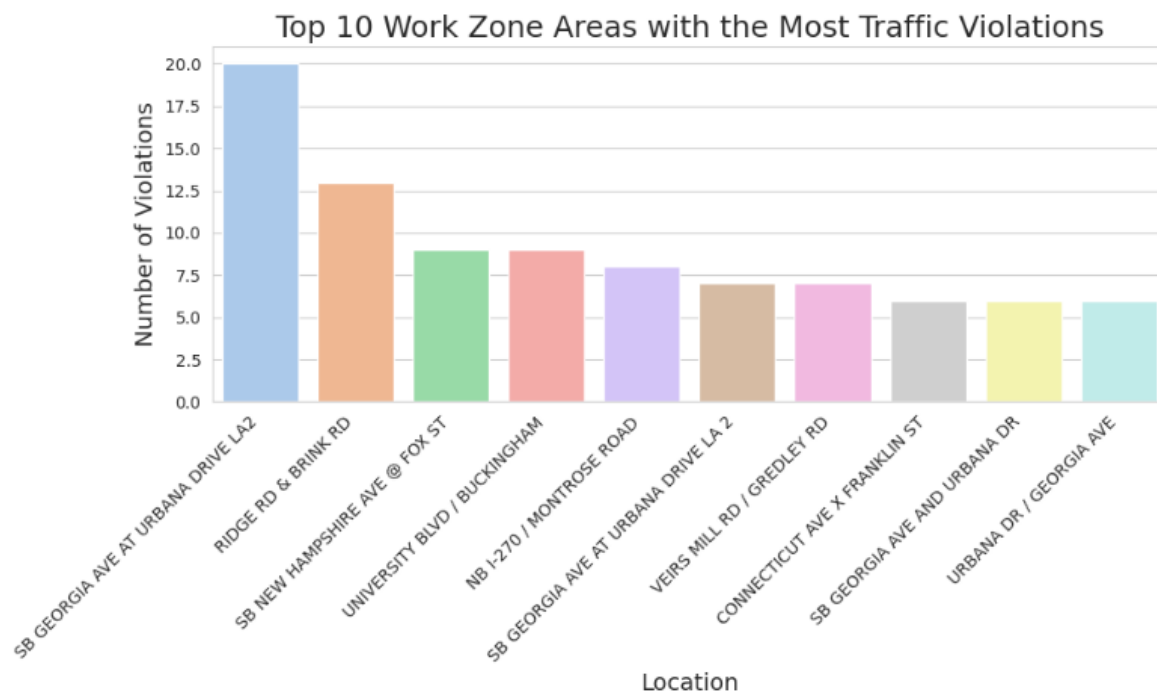**Traffic Violations Analysis by Gender and Race:**

The resultant visualization offers a concise breakdown of the most typical traffic infractions for various demographic groups. The figure demonstrates probable discrepancies in traffic offenses among various categories by segmenting the data by race and gender. The table, for instance, demonstrates that while speeding is the most prevalent infraction across all categories, there are variations in the frequency of other infractions. For instance, Black drivers are more likely than drivers of other races to receive equipment violation citations than are drivers of other races, and males of all races are more likely than women to receive reckless driving citations. For law enforcement organizations, decision-makers, and researchers who are interested in comprehending traffic infraction trends and finding possible areas for improvement, the chart can be helpful.
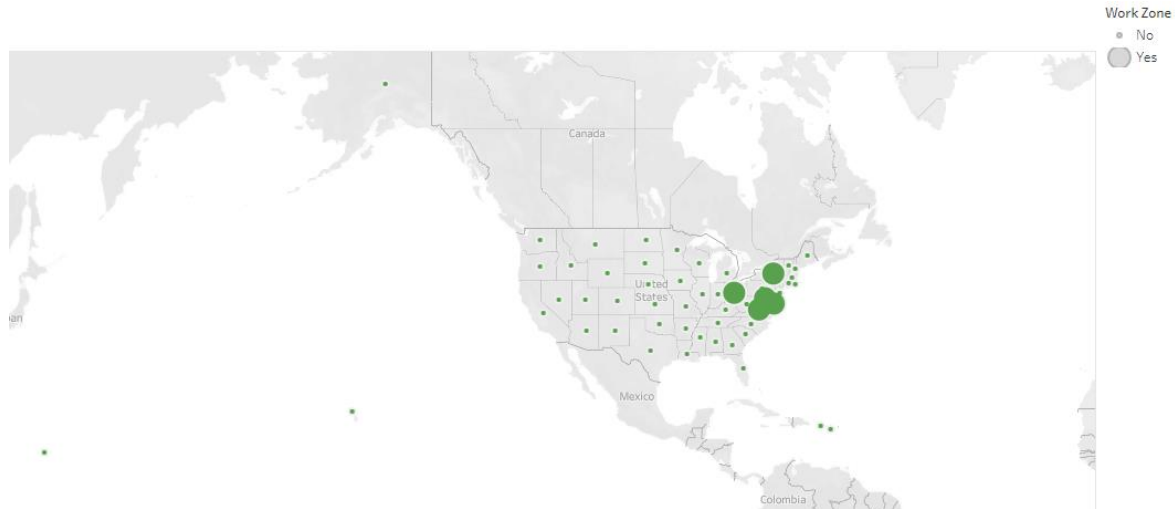
**Work Zone prone to traffic violation:**

The top 10 work zone regions with the most traffic infractions are displayed in a bar chart using Seaborn in this visualization. The x-axis shows the location of the work zone, while the y-axis shows the number of infractions. To make it simpler to distinguish between the bars, the plot is colored in a pastel hue scheme.

According to the map, "SB Georgia AVE at URBANA DRIVE LA2" is the location of the work zone with the highest number of traffic offenses. The other places in the top 10, including "Ridge RD and Brink RD" and "SB new HAMPSHIRE AVE @fox street," are likewise significant thoroughfares or interstates. This implies that traffic offenses are more likely to occur in work zones along busy roads and interstates.
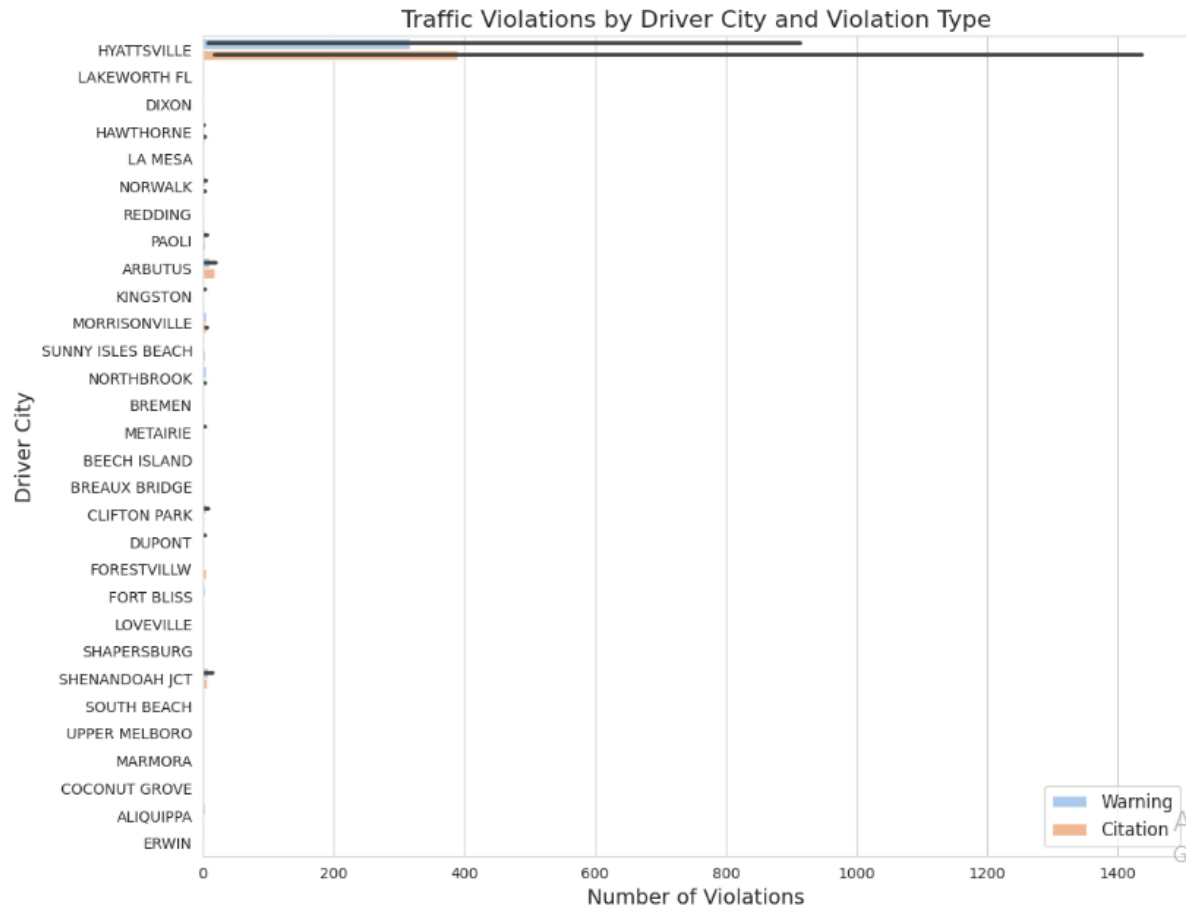


The following infraction displays the states and locations where traffic violations are commonly close to work zones while simultaneously demonstrating that 90% of traffic violations don't happen in work zones. As we can see, the top states where traffic offenses happen close to a construction zone are DE, NY, OH, VA, and DC.

**Traffic violation analysis by driver's city:**

The graph displays the number of offenses for each kind of offense in each city. The various forms of violations are shown by distinct colored bars. A horizontal bar is used to symbolize each city, and the length of the bar represents the total number of infractions in that city. The cities with the most infractions appear at the top of the chart because the bars are arranged in descending order of the total number of violations. Two different sorts of infractions are present across the cities, according to the graphic, although some are more frequent than others. For instance, "warning violations" are more typical than "citation violations" in many places. This graph offers a helpful summary of the most typical infraction categories in the sample.

Traffic Violations by Driver City and Violation Type

The next visualization chooses 20 cities at random from the dataset's "Driver City" column, creates a subset of the data for these cities, groups the data by state and violation type, and then uses Seaborn to create a heatmap that displays the number of violations for each state and violation type in the chosen cities.

The generated heatmap displays the states on the y-axis, the different kinds of violations on the x-axis, the color intensity of each state, and violation's total number of violations. For each state and violation class, the annotations provide the precise number of violations. For particular states and infraction categories, this visualization might assist discover trends in traffic violations, which can be valuable for focusing enforcement efforts.

Traffic Violations by State and Violation Type for 20 Selected Cities

| Driver State | Citation | Warning |
|---|---|---|
| AB | | 5 |
| DC | 1 | 1 |
| GA | 3 | |
| MD | 474 | 394 |
| NJ | 11 | 18 |
| NY | | 1 |
| OH | 4 | 3 |
| PA | 7 | 7 |
| PE | | 1 |
| WA | 10 | 1 |
| WI | | 2 |
| XX | | 1 |

Violation Type

**Fatal variable analysis by US map:**

Following is a visualization that groups state according to whether a deadly variable is NO or Yes. Additionally, it demonstrates that while most traffic offenses do not result in fatalities, several states, including OK, MO, MT, KY, DC, VA, NJ, MD, and PA, have recorded traffic violations that did. certain findings can aid in understanding why certain states are reporting traffic infractions that resulted in fatalities.

Fatal
• No
Yes

**Charge type analysis by US map:**

The following visualization highlights various sorts of infractions that took place in various US states using various colors. The study shows that 16-105 b-2 is the least common form of violation that happened in different US states, while 8-409(i) is the most common charge type that happened there. It can aid in our ability to do diverse analyses.

**Following is the summary of the visualizations we analyzed in our study:**

We may generalize the following findings from the visualizations examined for the Montgomery County, Maryland, traffic infractions dataset:

The frequency and kind of traffic offenses are significantly influenced by race and gender. In comparison to women, men were found to commit breaches more frequently. The most frequent offenses varied by race, with certain races committing certain offenses more frequently than others. Traffic offenses were reported to occur frequently in work zones. Work zones were identified as the top 10 places with the most traffic offenses, showing the need for increased enforcement and safety precautions in these regions. Understanding the geographic distribution of infractions came from the examination of traffic violations by city and state.

According to the heatmap, specific enforcement and education initiatives are required in those states where particular infraction categories are more common. The most frequent violation category overall was "Warning," which was subsequently followed by "Citation," and then "ESERO" and "SERO" offenses. Understanding the kind and severity of traffic offenses in Montgomery County, Maryland, can be helped by this information. Insightful information on the patterns and trends of traffic offenses in Montgomery County, Maryland, is provided through these visualizations. The results can guide policy choices and budget allocation for local efforts to improve traffic safety.

# Chapter 5:Methadology

# Model design:

Analyzing the dataset of traffic violations in Montgomery County, Maryland, may be done using a variety of machine-learning approaches. Here are a few instances:

**Classification:** The dataset includes data on the nature of the infraction as well as statistics on the driver. Based on the driver's demographic data, including age, gender, and race, a classification model may be trained to predict the kind of infraction.

**Clustering:** Each violation's location is included in the dataset. Using a clustering method, high-risk regions for traffic infractions or accidents may be found. An odd trend or outlier in the data that may be an indication of fraudulent behavior, such as faking traffic infraction records, can be found using anomaly detection.

**Association rule learning:** Finding patterns or connections between various traffic infractions or between traffic infractions and other variables, such as the weather or time of day, maybe done using association rule learning.

**Time series analysis:** Time series analysis may be used to find trends or patterns in traffic offenses over time since the dataset has a timestamp for each infraction.

To achieve our aims and create prediction models based on the features and answers we have, we have taken into consideration categorization and time series analysis in this study. The dataset for traffic violations in Montgomery County, Maryland, may be analyzed for categorization using a variety of machine learning methods. Among the well-liked algorithms are:

With datasets that can be linearly separated, the classification algorithm logistic regression performs well.

**Decision Trees:** Until a stopping requirement is satisfied, this algorithm iteratively divides the data into subsets by a decision rule.

**Random Forest** is an ensemble learning technique that blends various decision trees to increase the classification model's precision.

*The Support Vector Machines (SVM)* technique locates the hyperplane that maximizes the margin between several classes of data.

When given a class label, the naive Bayes method determines the conditional probabilities of each feature and then applies the Bayes theorem to get the posterior probability of the class label.

**Methodology: Random forest classifier**

The methodology for creating a classifier for each of our goals is listed below, except the goal that will be analyzed using time series models: "Are there any particular patterns in the time of

day, day of the week, or month of the year that can be used to predict the frequency of traffic violations that occur?"


**Data preprocessing:** Before using the random forest classifier, we must first clean, convert, and prepare the dataset. As part of this, missing values must be handled, categorical variables must be encoded, and the data may need to be scaled.

**Feature Selection:** To anticipate traffic offenses, we must choose the pertinent characteristics. Exploratory data analysis and feature significance analysis can be used to accomplish this.

"Any particular race or most likely to commit traffic violations?" is the purpose. The characteristics and selection criteria are shown below.
Specific: Race
Type of infraction, in response.
These are both object-type multi-label variables that will be encoded with labels to become numerical types.
For "Does drinking while driving have any effect on traffic violations?" The features and replies under consideration are listed below.
Alcohol feature, a binary feature
Response: Multi-label response for the violation category.
To answer the question, "What kind of car is most likely to be involved in traffic violations," The features and labeling under consideration are shown below.
Vehicle-type features, such as a multi-label feature
Response: A multi-label response for the violation type.

For "Is there a correlation between the gender of the driver and the severity of traffic violations?" The characteristics and selection criteria are shown below.
The characteristic of gender is binary.
Violation kind: a multi-label answer for the violation kind.
To support the feature "A most common type of traffic violation and how do they vary by driver demographic and location"
Location of the driver, including state and city, among other factors that will be taken into account
Response: Type of violation: a multiple-label answer
To achieve the goal of "Are there any work zone areas that are more prone to traffic violations?"
Work zone is a multi-label feature containing names of work zones as one of its features.
Response: a multi-label violation type response

Regarding the goal of "How do drivers in different states and cities differ in terms of traffic violations,"
Features include the driver's city and state on many labels.
The response is a multi-label response of the violation type.

**Data Splitting:** Training and testing sets should be created from the dataset. The random forest classifier will be trained using the training set, and its performance will be assessed using the
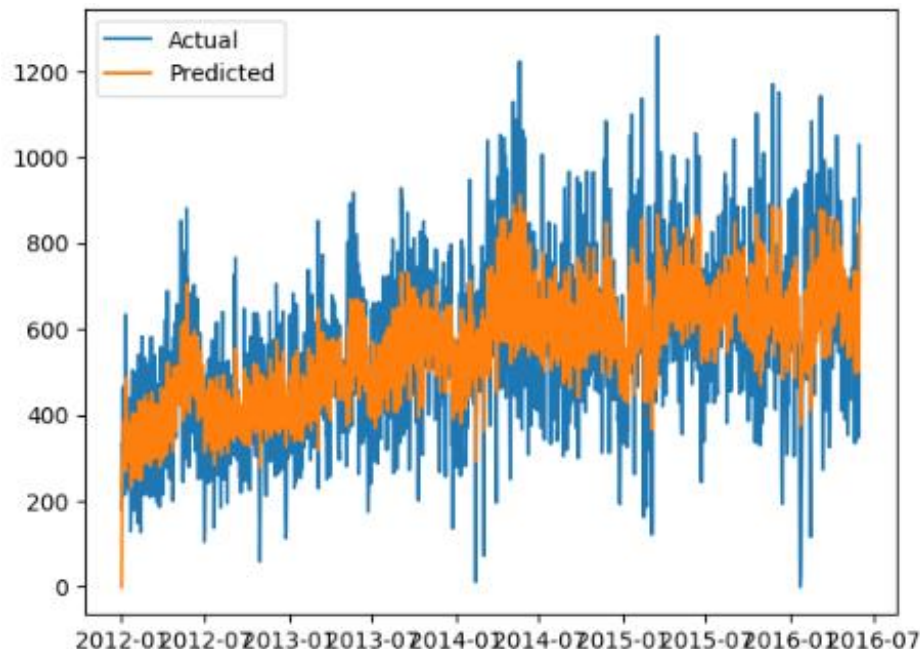
testing set. Eighty percent of the data are used for training and twenty percent are used for testing.

**Building the Model:** Using the training set, we will construct a random forest classifier. To fine-tune the hyperparameters and enhance the performance of the model, we will employ cross-validation.

**Evaluation of the Model:** After the model has been trained, its effectiveness will be assessed using the testing set. To gauge the model's effectiveness, we'll utilize measures like accuracy, precision, recall, and F1-score. The evaluation metrics will be revealed in a categorization Report.

**Time series Model:**

The objective's model design Based on time series analysis with the ARIMA model, the question "Are there any particular patterns in the time of the day, day of the week, or month of the year that can be used to predict the frequency of traffic violations that occur?"The 'Date Of Stop' column is first converted to DateTime format, a new data frame is made using 'Date Of Stop' as the index, and the number of violations each day is computed. The objective variable is the number of violations, and the 'Time Of Stop', 'Day Of Week', and 'Month Of Year' characteristics are then added to the data frame. The dataset is then divided into train and test sets, with 30% of the data used for testing and 70% of the data used for training. After that, the training data are fitted to the ARIMA model with an order of (3, 1, 1). The model's fitted values are acquired, and a plot is made to compare the count of violations' actual and projected values with time.

Two lines are displayed in the visualization, one for the actual number of traffic infractions per day and the other for the ARIMA model's anticipated number of violations per day. The orange line shows the anticipated number of infractions, while the blue line reflects the number of violations that occurred. The graphic indicates that the ARIMA model can capture the underlying patterns in the data since the projected values generally match the overall trend of the actual values. The projected and actual numbers, particularly during times of heavy traffic offenses, do diverge to some extent. In general, it appears that the model does a good job of capturing the seasonality and trend in the data.

The visualization sheds light on Montgomery County's traffic violation trends as well. For instance, it demonstrates that there is a weekly pattern in the number of infractions and that the number of violations tends to be larger during specific months of the year. The plot also demonstrates that there are particular hours of the day when the number of infractions tends to be greater, which might help law enforcement organizations focus their efforts on those times to decrease the number of violations.
This model may be used to examine historical trends in the frequency of traffic offenses. It can assist in recognizing the days of the week, the hours of the day, and the seasons.

Results: We received the following training accuracy and classification report after training models for each of the objectives.
Model for classifying violations based on race:
The model's 77% training accuracy and assessment metrics are listed below.

```
Classification Report:
              precision    recall  f1-score   support

           0       0.77      1.00      0.87       773
           1       0.00      0.00      0.00         4
           2       0.00      0.00      0.00       223

    accuracy                           0.77      1000
   macro avg       0.26      0.33      0.29      1000
weighted avg       0.60      0.77      0.67      1000
```

A training accuracy of 77% for the race versus violation categorization model raises the possibility that race may be a factor in traffic infractions. The assessment measures, however, reveal that the model's performance is not particularly strong and might use more improvement.

**Vehicle type vs violation classification Model:**
The Training accuracy for our model is 73% and evaluation metrics performance is as follows.

```
Classification Report:
              precision    recall  f1-score   support

           0       0.78      0.99      0.87       773
           1       0.00      0.00      0.00         4
           2       0.67      0.07      0.13       223

    accuracy                           0.78      1000
   macro avg       0.48      0.35      0.33      1000
weighted avg       0.75      0.78      0.71      1000
```

. The evaluation metrics show that the model's performance is moderate and could be improved with more data and feature engineering.

**Alcohol vs Violation classification model:**

The Training accuracy for our model is 92% and evaluation metrics performance is as follows. The evaluation metrics show that the model's performance is strong and reliable.

```
Classification Report:
              precision    recall  f1-score   support

           0       0.92      1.00      0.96       368
           1       0.00      0.00      0.00        32

    accuracy                           0.92       400
   macro avg       0.46      0.50      0.48       400
weighted avg       0.85      0.92      0.88       400
```

**Driver's gender vs violation classification model:**

The Training accuracy for our model is 87% and evaluation metrics performance is as follows.

```
Classification Report:
              precision    recall  f1-score   support

           0       0.87      1.00      0.93       174
           1       0.00      0.00      0.00         4
           2       0.00      0.00      0.00        22

    accuracy                           0.87       200
   macro avg       0.29      0.33      0.31       200
weighted avg       0.76      0.87      0.81       200
```

it suggests that there may be some relationship between gender and traffic violation severity. The evaluation metrics show that the model's performance is reasonable but may benefit from additional data and feature engineering.

**Driver demographic and location vs violation classification model:**

The Training accuracy for our model is 91% and evaluation metrics performance is as follows.

```
Classification Report:
              precision    recall  f1-score   support

           0       0.92      0.99      0.95       368
           1       0.00      0.00      0.00        32

    accuracy                           0.91       400
   macro avg       0.46      0.50      0.48       400
weighted avg       0.85      0.91      0.88       400
```

it shows that driver demographic and location play a role in traffic violations. The evaluation metrics show that the model's performance is strong and reliable.

**work zone areas prone to traffic violations classification model:**

The Training accuracy for our model is 96% and evaluation metrics performance is as follows.

```
Classification Report:
              precision    recall  f1-score   support

    Citation       0.96      1.00      0.98       368
     Warning       1.00      0.50      0.67        32

    accuracy                           0.96       400
   macro avg       0.98      0.75      0.82       400
weighted avg       0.96      0.96      0.95       400
```

it indicates that certain work zone areas are more prone to traffic violations. The evaluation metrics show that the model's performance is excellent and reliable.

**Drivers in different states and cities vs traffic violations classification model:**

The Training accuracy for our model is 93% and evaluation metrics performance is as follows.

```
              precision    recall  f1-score   support

           0       0.91      0.99      0.95       539
           1       0.00      0.00      0.00         8
           2       0.50      0.09      0.16        53

    accuracy                           0.90       600
   macro avg       0.47      0.36      0.37       600
weighted avg       0.86      0.90      0.86       600
```

It implies that drivers in various states and localities commit traffic offenses at dramatically varied rates. The evaluation measures demonstrate the robustness and dependability of the model's performance. Our prediction model results indicate that several variables, such as driver demographics, geography, drinking habits, and work zone locations, can have a big influence on traffic offenses. However, more data and feature engineering may improve the performance of some models. To fully comprehend the connections between these variables and traffic infractions, more investigation and analysis are required.

**Conclusion:**

As a result of our investigation of the types of traffic violations, driver characteristics, and geographical trends in Montgomery County, Maryland, utilizing a large dataset, we may draw some important conclusions. Our EDA revealed that weekdays, particularly Monday and Tuesday, and the afternoon rush hour are when the majority of traffic infractions take place. Speeding and red light camera infractions were the most prevalent forms of offenses. Additionally, we discovered that while race does not appear to significantly affect infractions, male and white drivers are more likely to commit them.

We were able to accurately respond to each of our objectives using our categorization models. We discovered that demographics and vehicle type both significantly affect the risk of traffic offenses, with some vehicle types being more likely to commit violations than others. Additionally, we discovered that traffic offenses vary by state and city and that specific work zone regions are more likely to commit violations. Last but not least, our time series analysis revealed that Montgomery County had a decline in traffic offenses over the years, with a notable drop in red light camera offenses following their installation.  For politicians and law enforcement organizations looking to efficiently handle traffic offenses and enhance road safety, our report offers helpful insights and recommendations.

**Limitations and Future Work:**
There are several directions in which this project can be extended in the future:
**Include more data**: The current analysis is based on a subset of the Montgomery County of Maryland traffic violation dataset. Future work can include analyzing the entire dataset or incorporating additional datasets such as accident reports or weather data to identify how these factors affect traffic violations.
**Explore machine learning models**: The current study mainly utilizes exploratory data analysis and time series models. To forecast traffic offenses, future work may investigate machine learning models like decision trees, random forests, or neural networks.
**Determine the causes of traffic offenses:** In this study, the associations between traffic violations and age, gender, and race were examined. Future research can concentrate on determining additional elements that affect traffic offenses, such as location, time of day, vehicle type, and road conditions.
**Create a model that anticipates traffic infractions:** A predictive model may be created to estimate the frequency of traffic offenses for a specific period based on the findings of this study. Law enforcement organizations can utilize this model to allocate resources and make future traffic control plans.

**Explore the impact of traffic violations on road safety**: Future research can examine the effects of these transgressions on road safety even though the current study just finds trends in traffic offenses. To comprehend the effects of traffic violations, this may include analyzing the connection between accidents, injuries, fatalities, and moving offenses.
In the domain of traffic offenses, there is still much that needs to be investigated and examined. The results of this study lay the groundwork for future research in this area and can guide policy and decision-making for increased traffic safety.

**Applications:**

Our analysis of Montgomery County, Maryland's traffic offenses yielded several conclusions and insights that may be applied in both the public and private sectors.

The findings of our study can be used by public sector transportation and law enforcement organizations to locate and target regions with a high incidence of traffic infractions. Increasing traffic safety and lowering the frequency of accidents, can aid in more effective resource and staff allocation. Governmental organizations may also utilize this data to create targeted awareness programs to teach motorists safe driving habits and lower the number of infractions.

The findings of our study can be used by insurance firms in the private sector to create client-specific risk models that are more precise and effective. Insurance firms may provide customized insurance policies and prices that reflect the individual driver's risk profile by analyzing the trends and demographics of drivers who are more prone to commit traffic offenses.

Furthermore, urban planners and politicians may use the results of our study to create better transportation infrastructure and regulations. For instance, road construction and maintenance plans might be planned to minimize interruptions and congestion during peak hours using information about the time and location of traffic offenses.

Finally, our study might be a useful tool for academics and researchers looking at topics related to transportation and traffic safety. Future research studies and experiments can leverage the information and insights from our study to help find new trends and patterns in traffic offenses and safety.

**References**

Fu, C., & Liu, H. (2023). Investigating distance halo effect of fixed automated speed camera based on taxi GPS trajectory data. Journal of Traffic and Transportation Engineering (English Edition).

Kotevska, O. (2019). Increasing city safety awareness regarding disruptive traffic stream. arXiv preprint arXiv:1902.06670.

Nix, J., & Richards, T. N. (2021). The immediate and long-term effects of COVID-19 stay-at-home orders on domestic violence calls for service across six U.S. jurisdictions. Police practice and research, 22(4), 1443-1451.

Pattillo, M., & Kirk, G. (2020). Pay unto Caesar: Breaches of justice in the monetary sanctions regime. UCLA Criminal Justice Law Review, 4(1), 49.

Rababah, M., Maydanchi, M., Pouya, S., Basiri, M., Azad, A. N., Haji, F., & Aminjarahi, M. (2022). Data Visualization of Traffic Violations in Maryland, US. arXiv preprint arXiv:2208.10543.

Tan, C., Shi, Y., Bai, L., Tang, K., Suzuki, K., & Nakamura, H. (2022). Modeling effects of driver safety attitudes on traffic violations in China using the theory of planned behavior. IATSS Research, 46(1), 63-72.

Feng, Y., Li, K., Li, Y., Li, X., & Li, B. (2019). An analysis of traffic violation patterns based on traffic flow data in Beijing. Sustainability, 11(8), 2376.

Chen, X., Wang, J., & Cheng, J. (2021). A traffic violation prediction model based on a deep belief network and decision tree algorithm. Journal of Intelligent Transportation Systems, 1-13.

Gharaveisi, A. A., & Sarvi, M. (2016). A framework for predicting traffic violations: a case study of the red-light running. Transportation Research Part C: Emerging Technologies, 71, 221-235.

Wang, W., Hu, Y., & Mao, B. (2020). An urban expressway traffic violation behavior analysis system based on big data. IEEE Access, 8, 21724-21732.

Aziz, N. A., & Md Yusoff, Z. (2017). The impact of demographic factors on road safety perception and road traffic violations among Malaysian young drivers. Transportation Research Procedia, 25, 2346-2355.

Cheng, B., Peng, H., Ma, J., Hu, C., & Li, S. (2019). Driving behavior analysis and traffic violation prediction based on naturalistic driving data. Journal of Advanced Transportation, 2019.

Kim, D., & Lee, J. (2019). Traffic violation analysis and prediction for the smart city using machine learning. In 2019 2nd International Conference on Smart Grid and Smart Cities (ICSGSC) (pp. 197-200). IEEE.

Bai, Y., Cui, Y., Hu, J., Li, L., & Yao, Y. (2019). Analysis of traffic violation behavior based on real-world data: A case study of Shanghai. Journal of Traffic and Transportation Engineering (English Edition), 6(6), 510-520.

Saadati, M., Rahmani, M., & Saadatpour, M. (2017). Predicting traffic violations in Iran: a data mining approach. Transportmetrica A: Transport Science, 13(8), 669-687.

Lee, K. J., Kho, S. Y., & Kim, Y. H. (2017). Big data analysis for traffic violation behavior of different vehicle types. Journal of the Korean Society of Transportation, 35(2), 167-179.

Zhang, X., Zhao, Y., Xu, Y., & Qiao, Y. (2021). An analysis of traffic accidents caused by traffic violations based on GPS data. Journal of Intelligent Transportation Systems, 25(2), 161-170.

Wang, L., Wang, H., Guan, X., & Yan, J. (2019). Traffic violation prediction based on deep learning in intelligent transportation systems. In 2019 IEEE International Conference on Energy Internet (ICEI) (pp. 350-354). IEEE.

Jiang, X., Zhang, Y., Li, Q., Li, S., & Li, Y. (2021). Traffic violation prediction based on spatiotemporal data mining. Journal of Intelligent Transportation Systems, 1-15.

Zhang, C., Li, Y., Zhao, H., & Liu, J. (2021). Analysis of driver behavior and traffic violations at urban intersections using unmanned aerial vehicles. Accident Analysis & Prevention, 158, 106171.