

# Maryland State traffic Collision Dataset

## Introduction:

This dataset provides comprehensive information on motor vehicle operators involved in traffic collisions that occur on county and local roadways within Montgomery County, Maryland. The data is collected through the Automated Crash Reporting System (ACRS) of the Maryland State Police and reported by various law enforcement agencies including the Montgomery County Police, Gaithersburg Police, Rockville Police, and the Maryland-National Capital Park Police. The dataset contains detailed records of each collision, including information about the drivers involved.

Please note that the collision reports included in this dataset are based on preliminary information provided to the Police Department by the reporting parties. Consequently, the data may include:

1. Information that has not yet been verified through further investigation.
2. A mix of verified and unverified collision data.
3. Preliminary collision classifications that may be subject to change based on subsequent investigations.
4. Data that may contain errors due to mechanical or human error.

## Data Source:

<https://catalog.data.gov/dataset/crash-reporting-drivers-data>

## Dataset Structure and Columns:

The dataset comprises numerous columns providing specific details about each collision incident. Key columns include:

- Report Number: Unique identifier for each collision report.
- Crash Date/Time: Date and time when the collision occurred.
- Injury Severity: Severity of injuries sustained by individuals involved in the collision.

- Driver At Fault: Indicates whether the driver was at fault in the collision.
- Weather: Weather conditions at the time of the collision.
- Surface Condition: Condition of the road surface at the collision site.
- Light: Lighting conditions at the time of the collision.
- Traffic Control: Type of traffic control present at the collision site (e.g., traffic signals, stop signs).
- Driver Substance Abuse: Whether substance abuse was a contributing factor for the driver.
- Vehicle Year, Make, Model: Details about the vehicles involved in the collision.
- Location: Latitude and longitude coordinates of the collision location.

### **Potential Use Cases:**

- Analyzing trends and patterns in traffic collisions to identify high-risk areas and improve road safety measures.
- Studying the impact of various factors such as weather, road conditions, and driver behavior on collision rates.
- Developing predictive models to forecast collision probabilities and allocate resources effectively for prevention and response.

### **Algorithms Intended to Apply:**

- 1. Predictive Models:** Develop models to predict the likelihood of crashes based on various factors like time, location, weather, and road conditions.
  - a. **Logistic Regression:** This algorithm can be utilized to predict the likelihood of crashes by modeling the relationship between the dependent variable (crash occurrence) and various independent variables such as time, location, weather, and road conditions. It's particularly effective when the target variable is binary, making it suitable for predicting crash/no-crash scenarios.
  - b. **Random Forest:** Random Forest is an ensemble learning technique that combines multiple decision trees to create a robust predictive model. It can handle both numerical and categorical data effectively and can capture

complex relationships between predictor variables and crash likelihood. Random Forest can handle missing data and outliers, making it suitable for real-world datasets with diverse characteristics.

- c. **Gradient Boosting Machines (GBM):** GBM is another ensemble learning method that builds predictive models by sequentially adding weak learners (decision trees) to minimize the loss function. It's highly effective in capturing interactions between predictors and can handle both regression and classification tasks. GBM tends to perform well with large datasets and can provide accurate predictions for crash likelihood based on various factors.

## 2. **Cluster Analysis:** Perform cluster analysis to identify groups of similar crash characteristics or patterns.

- a. **K-Means Clustering:** K-Means clustering is a widely used unsupervised learning algorithm for grouping similar data points into clusters. In the context of crash analysis, K-Means can identify groups of crashes with similar characteristics or patterns based on features such as time, location, weather, and road conditions. This can help in identifying distinct clusters representing different types of crashes or hotspots with similar contributing factors.
- b. **Hierarchical Clustering:** Hierarchical clustering is another clustering algorithm that organizes data into a hierarchy of clusters. It can be applied to identify nested clusters of crash characteristics, allowing for a more detailed understanding of similarities and differences between different crash groups. Hierarchical clustering is useful for exploring the hierarchical structure of crash data and identifying subgroups with shared characteristics.
- c. **DBSCAN (Density-Based Spatial Clustering of Applications with Noise):** DBSCAN is suitable for identifying clusters of varying shapes and densities in spatial data. It can be utilized to detect clusters of crashes occurring near each other, indicating potential accident-prone zones or areas with high crash densities. DBSCAN can handle noise and outliers effectively, making it robust for real-world crash data analysis.