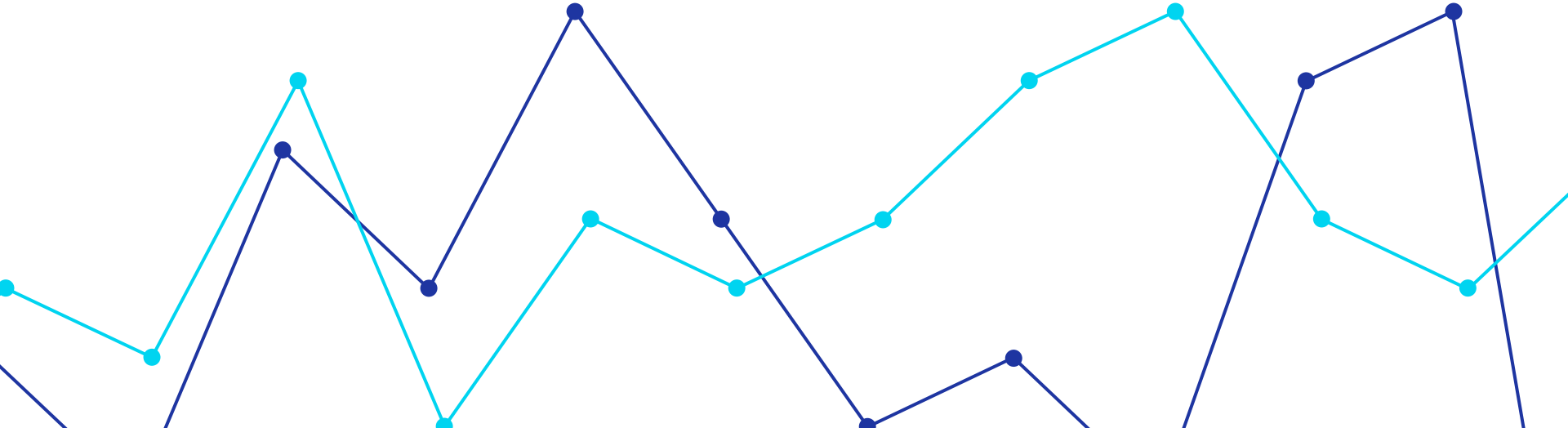


US Accidents Analysis

Ahmed Osama Helmy
Aliaa Gheis

Abdallah Ahmed
Omar Mahmoud



Problem Definition

- Road safety is a critical concern, and understanding accident patterns can help cities improve traffic management and reduce accident rates.
- This project aims to analyze accident data to identify high-risk locations, contributing factors, and potential mitigation strategies.
- By leveraging big data processing, we will extract valuable insights for transportation authorities and urban planners.

From Data to Insights

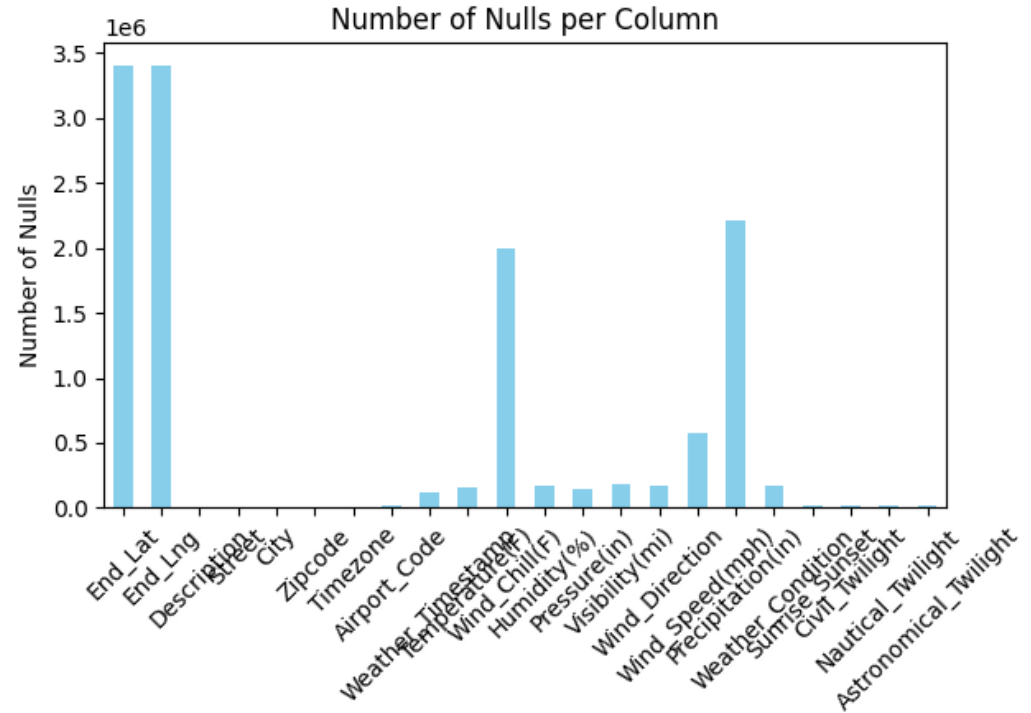
- **Data Exploration**
- **Data Visualization**
- **Data Cleaning**
- **Feature Engineering**
- **Aggregation or using models**
- **Visualizing & Extracting Insights**

Data Exploration/Visualization

Nulls in Data



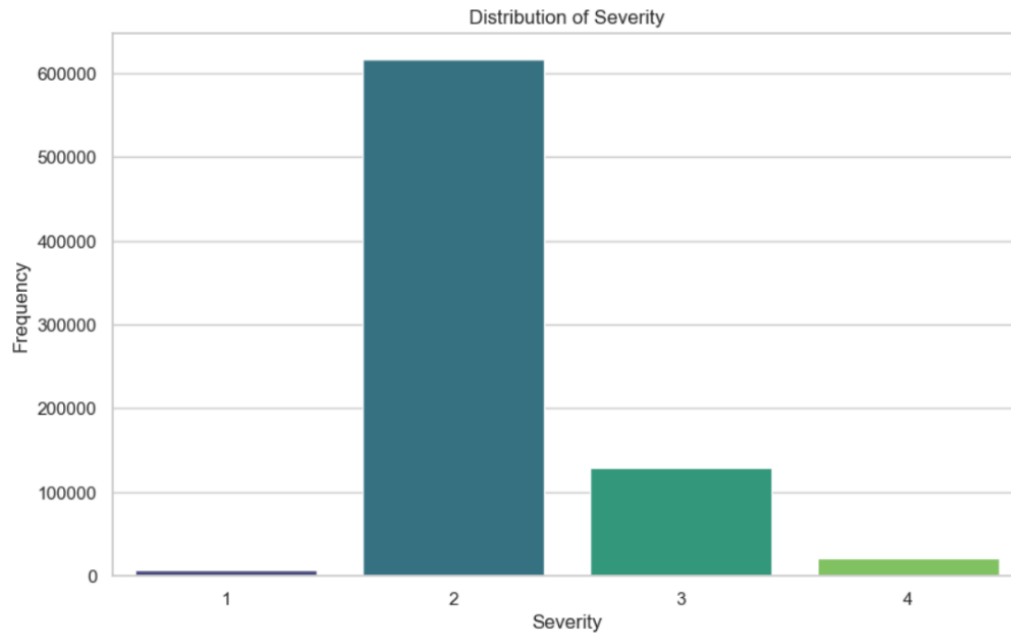
Some Columns have a lot of missing values



Severity Distribution



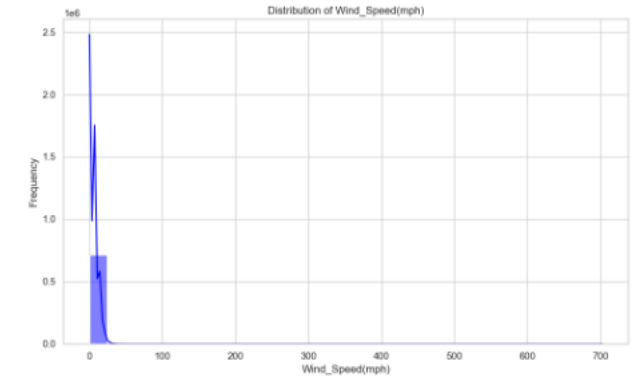
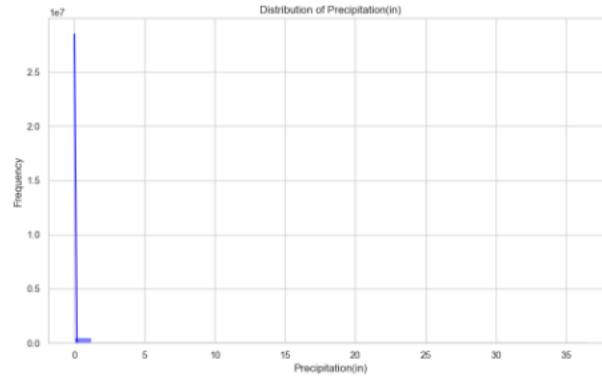
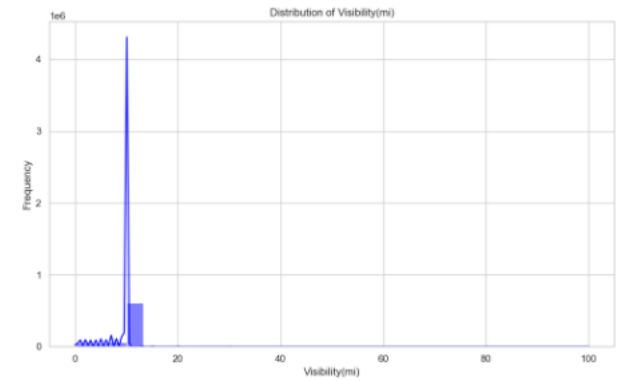
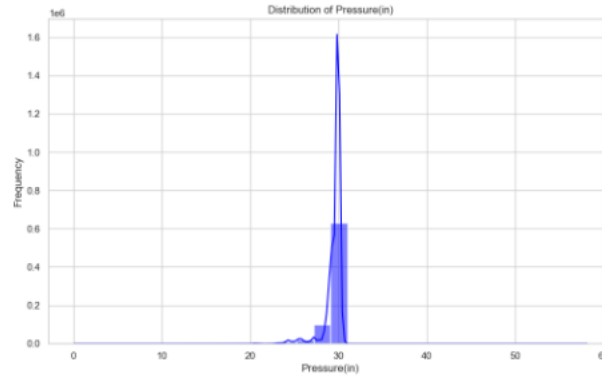
Severity 2 is
dominating



Outliers

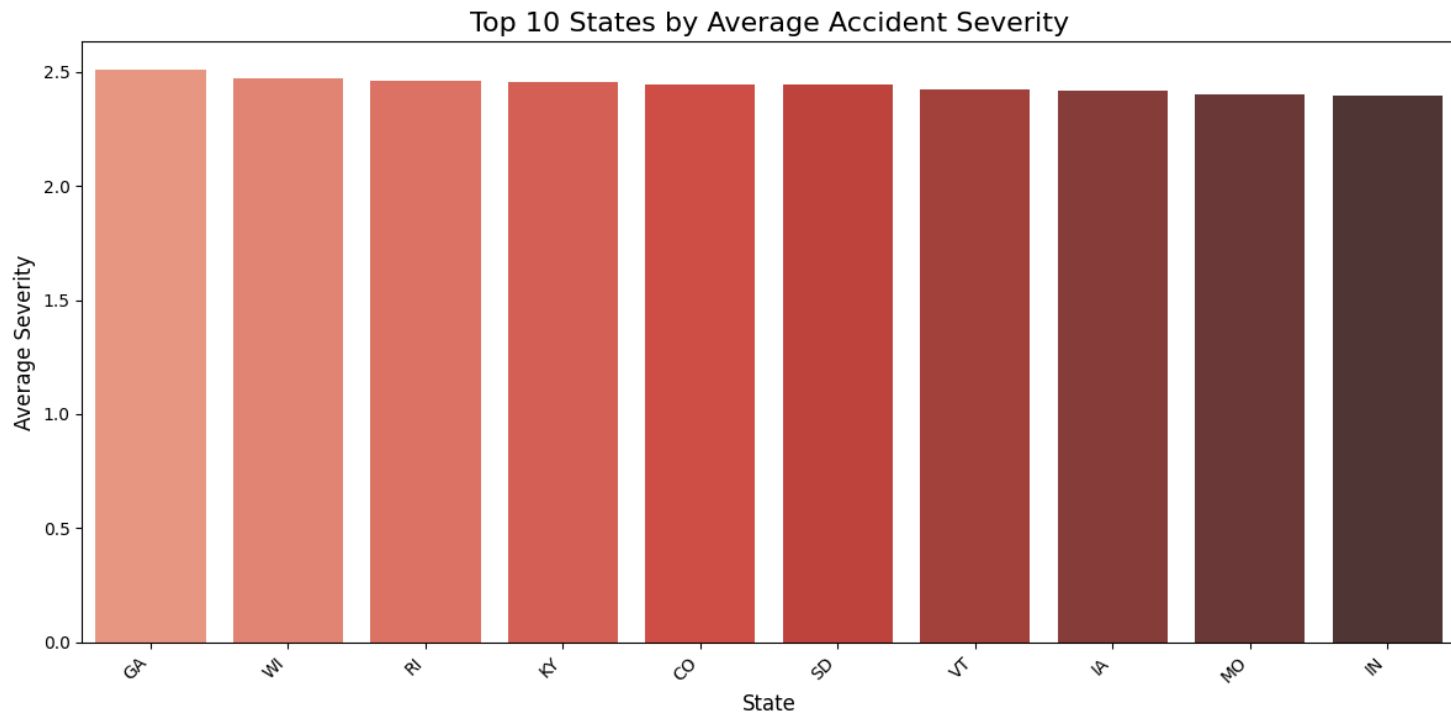


Some columns
have outliers



Avg Severity per state

Uniform
Distribution



Data Cleaning

Irrelevant Columns



Eliminates columns that don't provide meaningful data for modeling or analysis.



Columns such as ID, Source, Description, Street, City, Zipcode, Airport_Code, etc., are dropped as they are not useful for analysis or prediction tasks.

Handling Missing Values



Dropping Columns like End_Lat & End_Lng as the percentage of missing values was greater than 40%



Imputing missing values in numeric columns by inserting the mean value



Imputing missing values in categorical columns by inserting the mode value

Outliers



Eliminating records with temperature higher than 56.7 C as reported in this [article](#) that the maximum US temperature was 134.4°F (56.7°C)



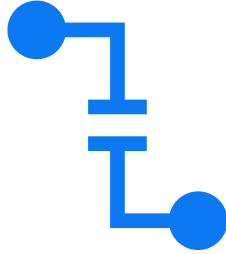
For Wind Speed values, we identified outliers by considering the maximum observed wind speeds. According to the World Meteorological Organization, the highest recorded wind speed was 254 mph (408 km/h). We decided to remove any records with wind speeds exceeding this threshold to eliminate extreme outliers.



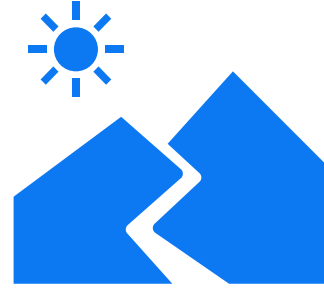
We used IQR in other columns

Feature Engineering

Adding Road-Related Features



A Boolean variable `Is_Complex_Road` was added to interpret whether the road is complex by utilizing the other variables like (Junction, Railway, Crossing)



This will help in giving insights into the effect of complexity of roads.

Risk Score



The dataset was aggregated at the state level to simplify the analysis and provide actionable insights at a regional scale.



Key metrics such as average accident severity, total accident count, and risk score were computed for each state.



The **Risk score** was calculated as the product of average severity and accident count, capturing both the frequency and severity of accidents.



The **Risk score** was then normalized to be from 0 to 1.



Also, A Boolean variable **Is_High_Risk** was added to detect if a state was high risky or not by using the 75th quartile.

Time Related Features

- **Hour of the Day:** Captures the time of day when accidents occur (e.g., morning rush hours, nighttime)
- **Day of the Week:** Identifies whether accidents are more frequent on weekdays or weekends.
- **Month:** Highlights seasonal trends in accident occurrences (e.g., higher rates during winter months due to adverse weather conditions).
- **Year:** Tracks long-term trends in accident frequency over multiple years.
- **Duration:** Tracks the duration of the accident in minutes by subtracting the start time from the end time
- **Season:** Determines the Season when the accident happened (Summer, ...)

Outliers

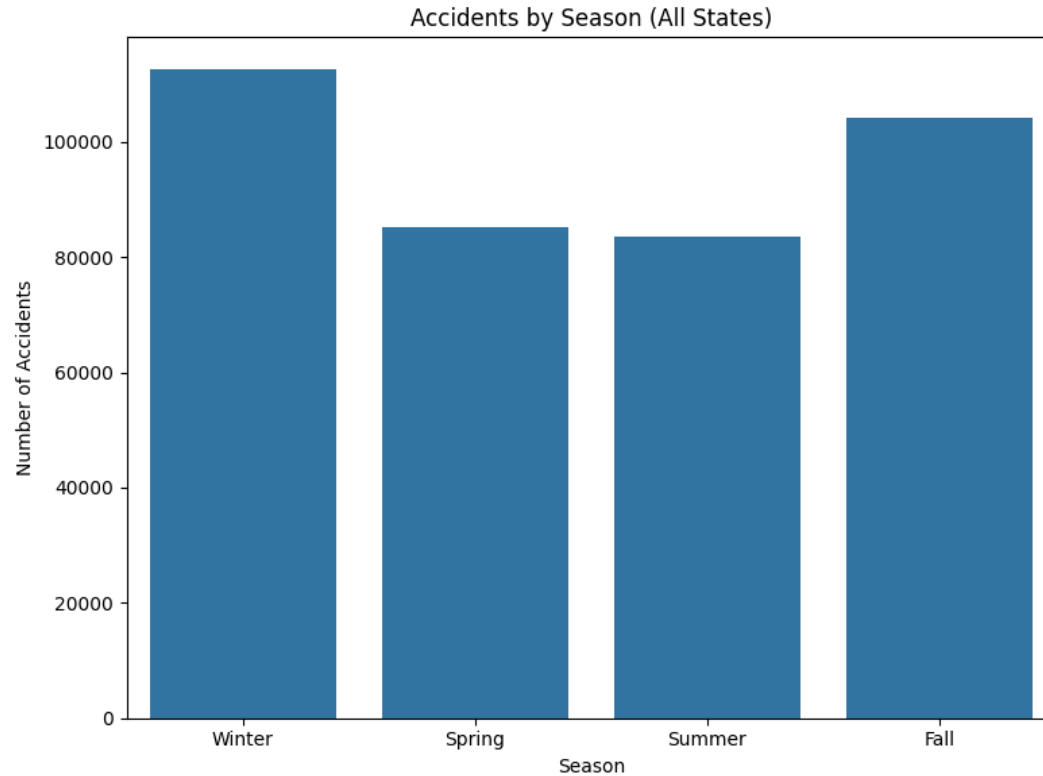
- Eliminating records with temperature higher than 56.7 C as reported in this [article](#) that the maximum US temperature was 134.4°F (56.7°C)
- For Wind Speed values, we identified outliers by considering the maximum observed wind speeds. According to the World Meteorological Organization, the highest recorded wind speed was 254 mph (408 km/h). We decided to remove any records with wind speeds exceeding this threshold to eliminate extreme outliers.
- We used IQR in other columns

Outliers

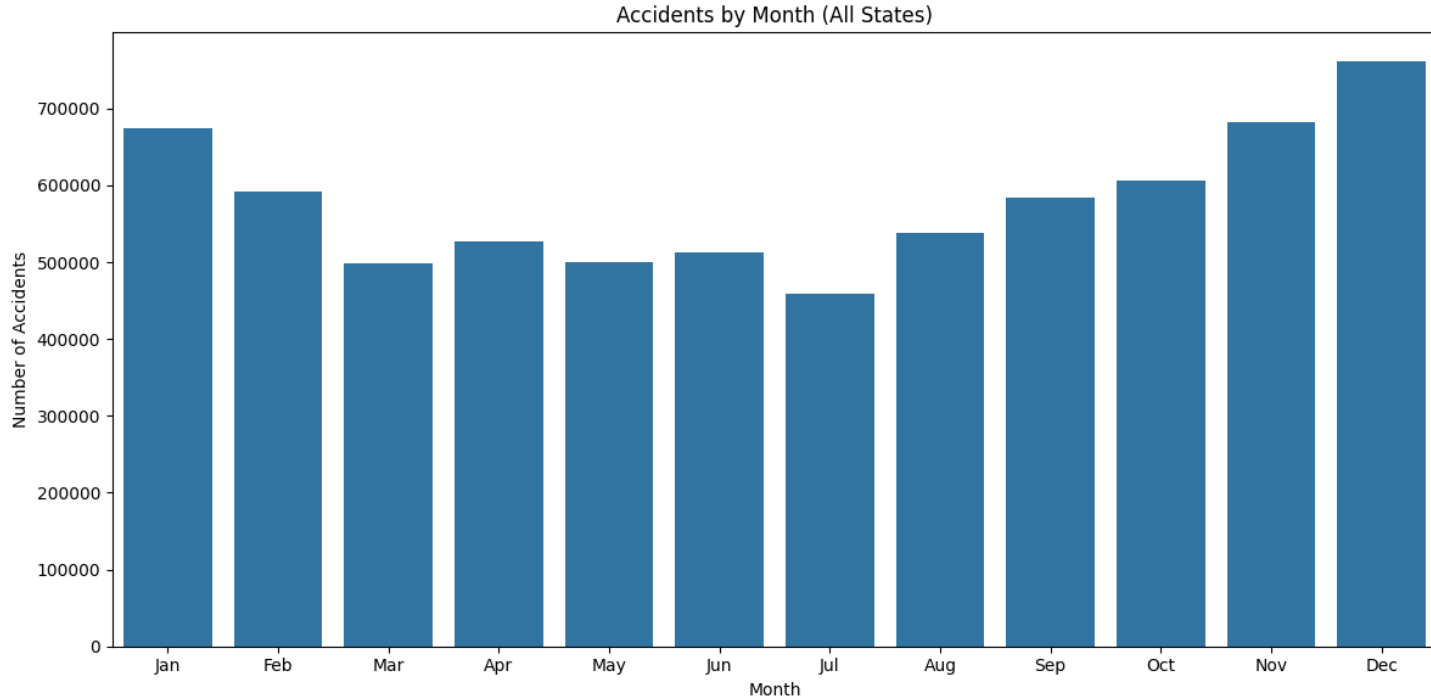
- Eliminating records with temperature higher than 56.7 C as reported in this [article](#) that the maximum US temperature was 134.4°F (56.7°C)
- For Wind Speed values, we identified outliers by considering the maximum observed wind speeds. According to the World Meteorological Organization, the highest recorded wind speed was 254 mph (408 km/h). We decided to remove any records with wind speeds exceeding this threshold to eliminate extreme outliers.
- We used IQR in other columns

Accidents based on Time

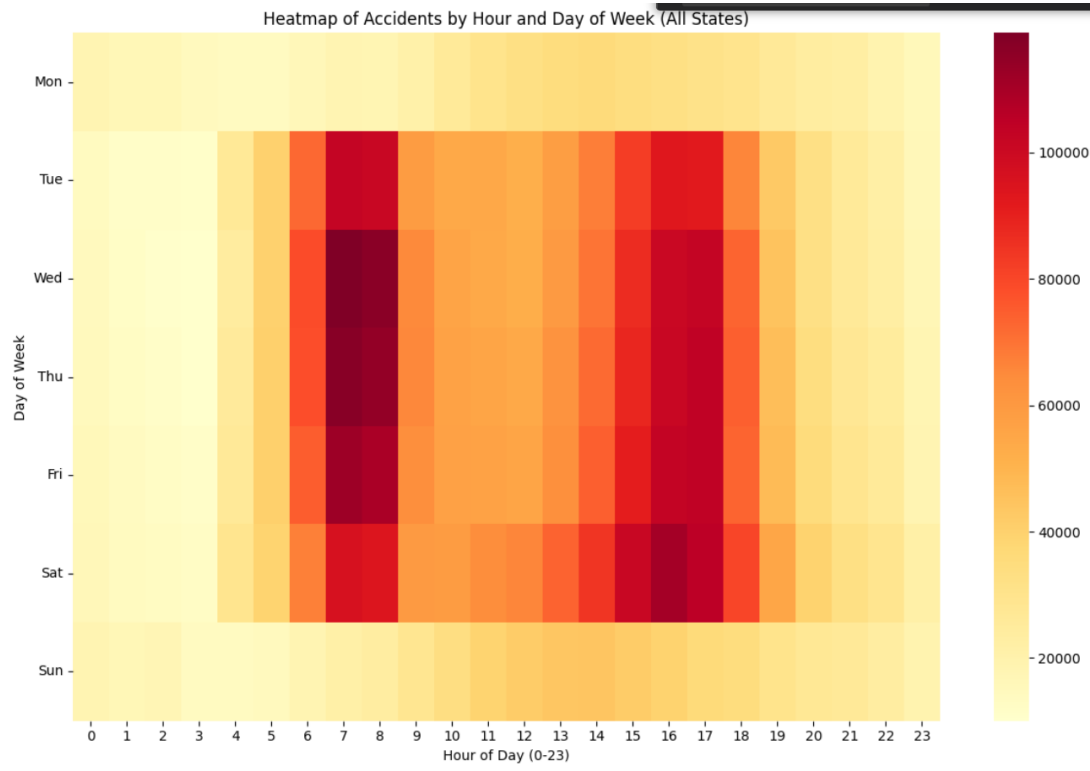
Accident Trends by Season.



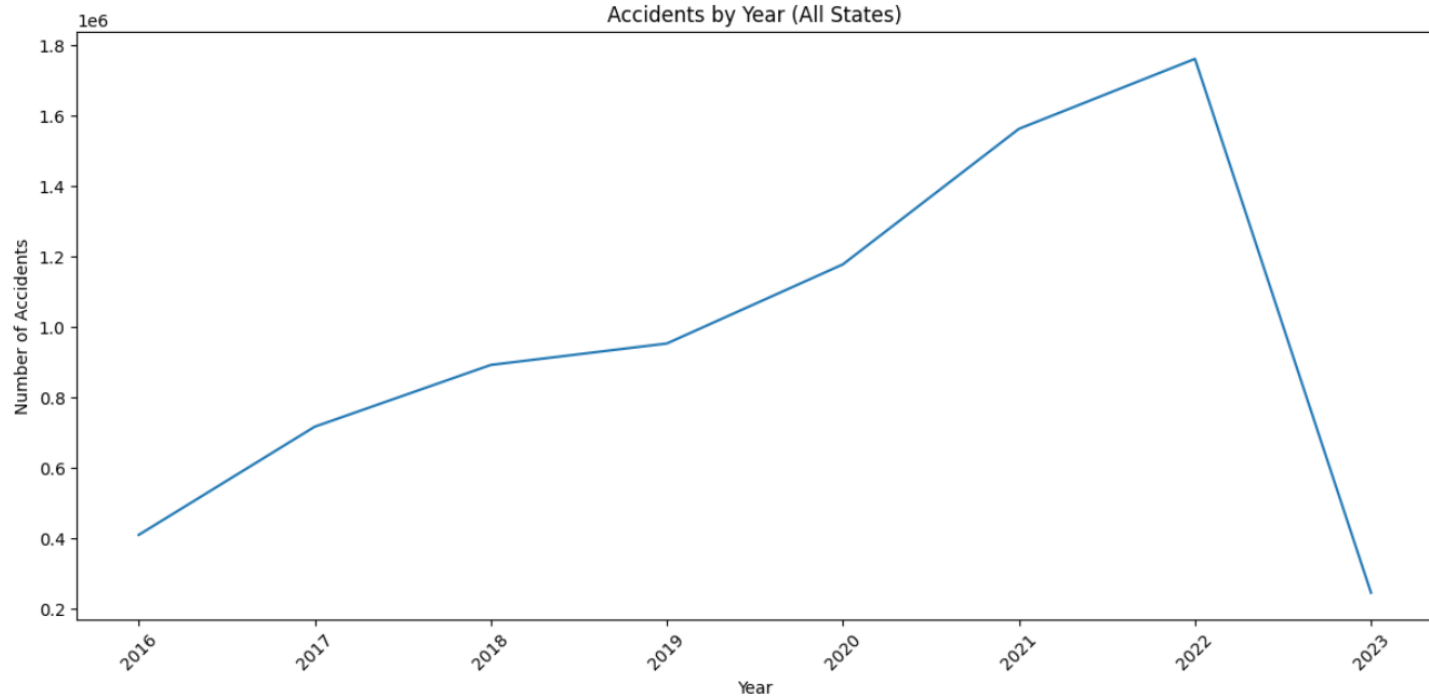
Accident Trends by Month.



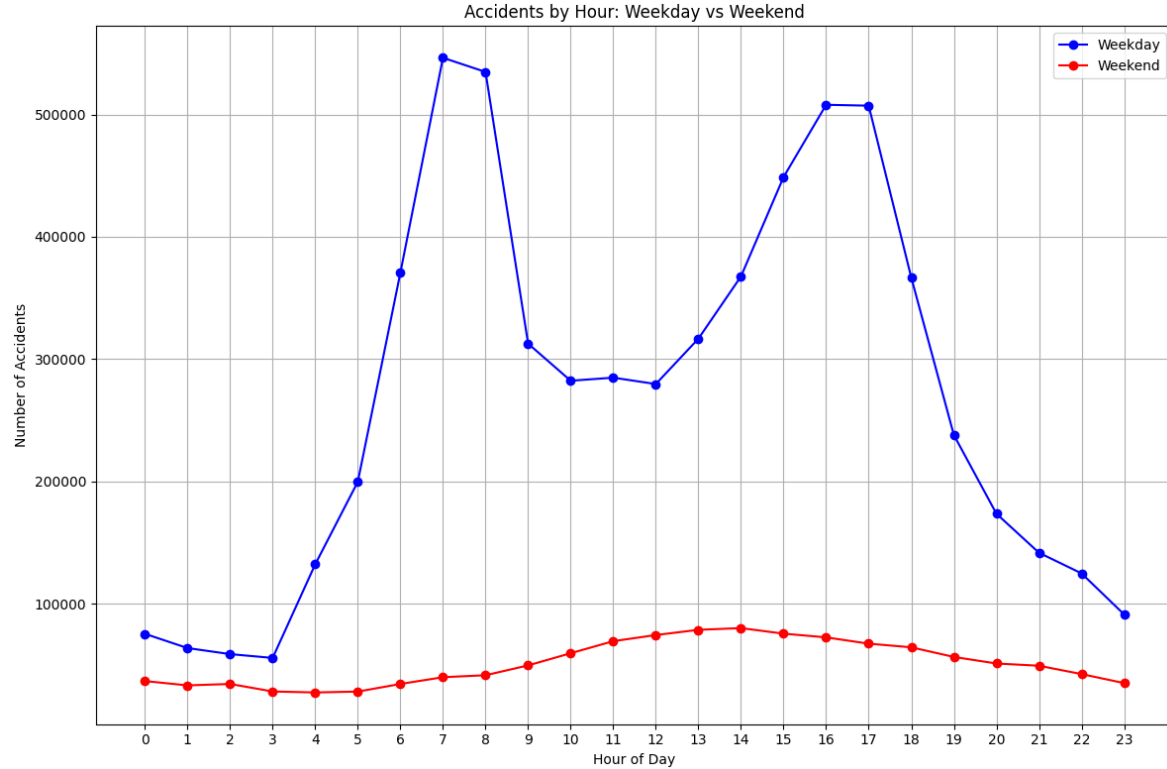
Accident Trends by Hour.



Accident Trends by Year.

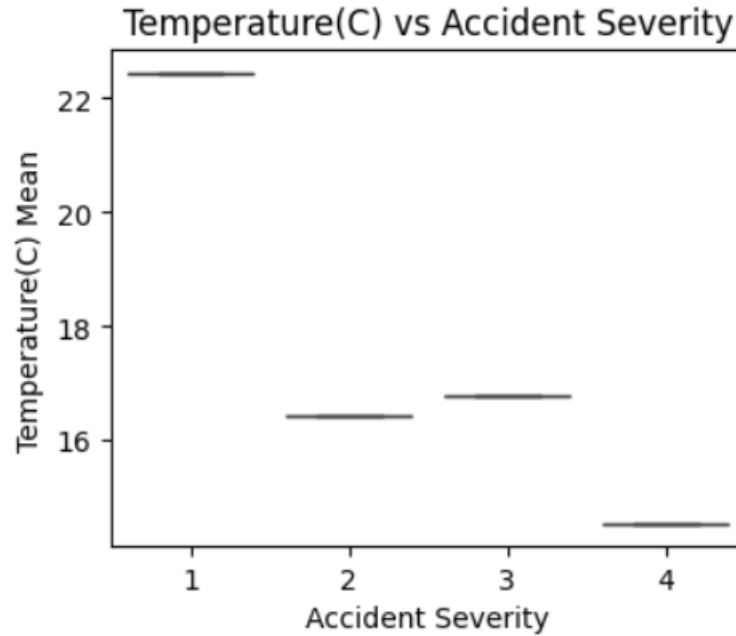


Weekends and Weekdays behavior

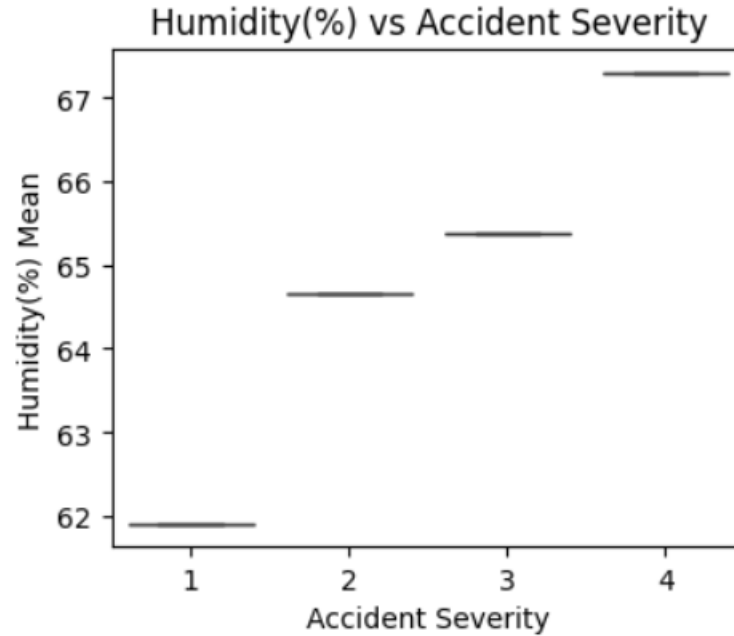


Weather Effect on Accidents

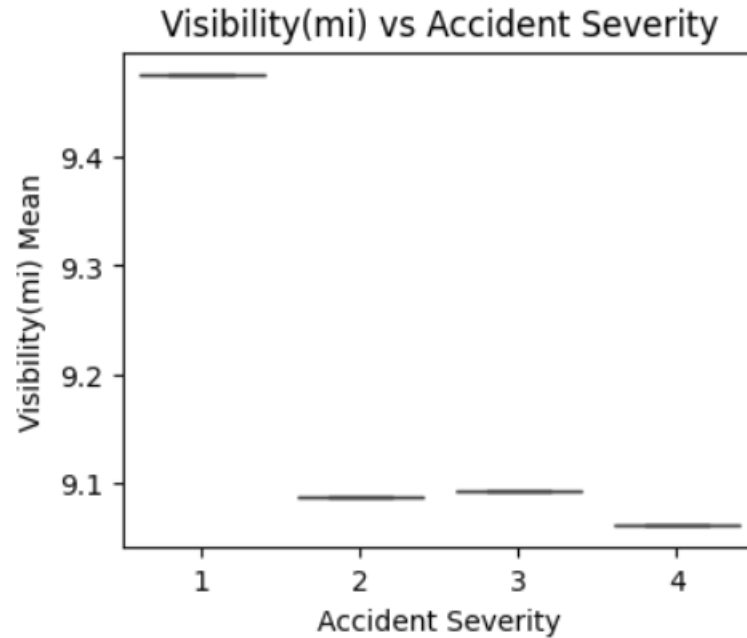
Temperature & Severity



Humidity & Severity

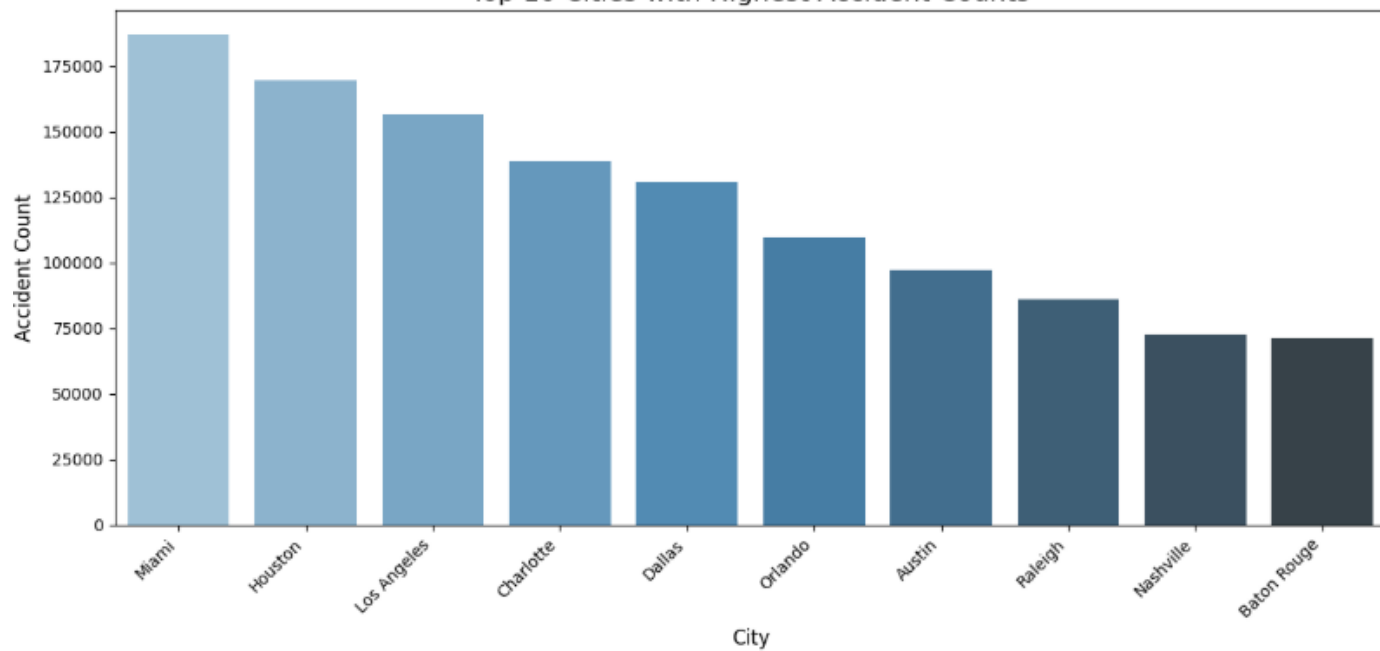


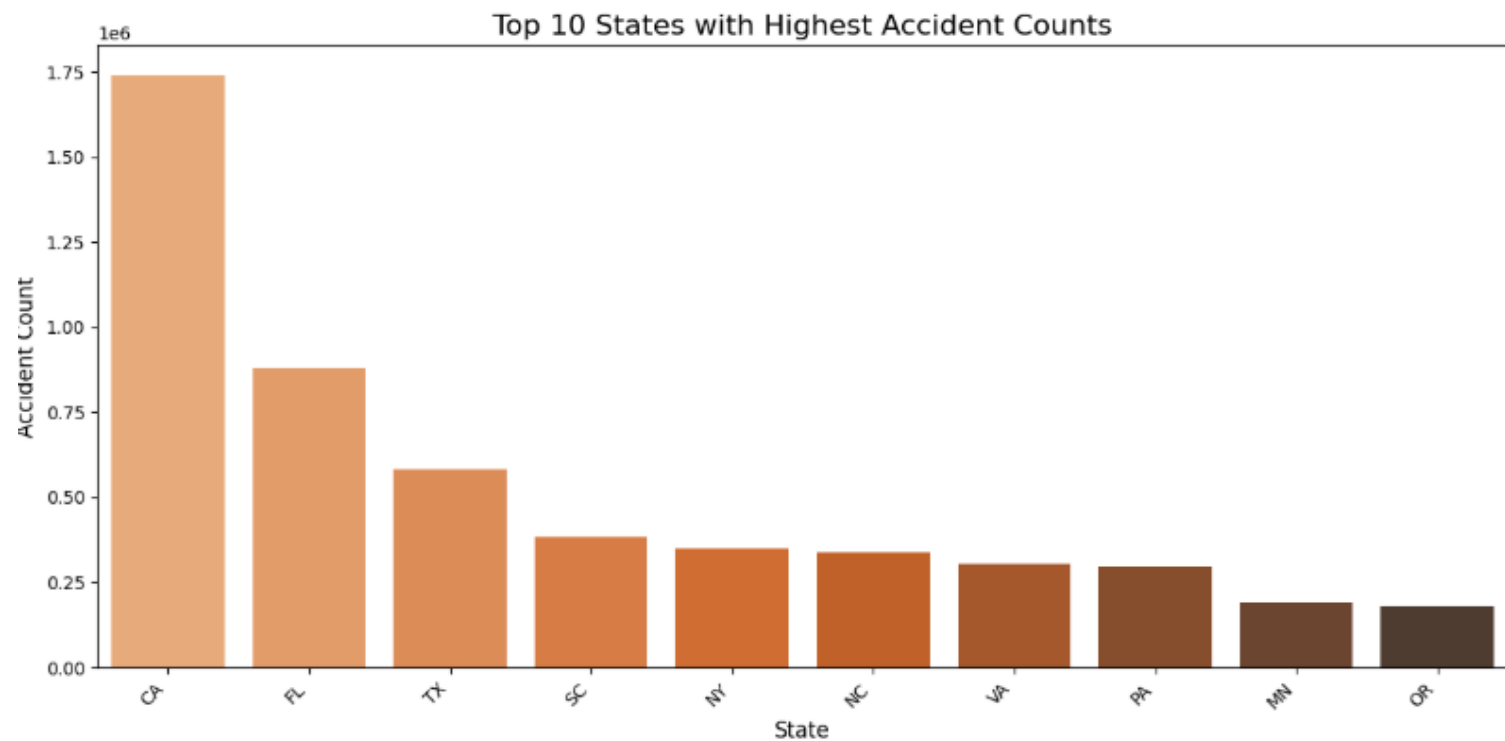
Visibility & Severity



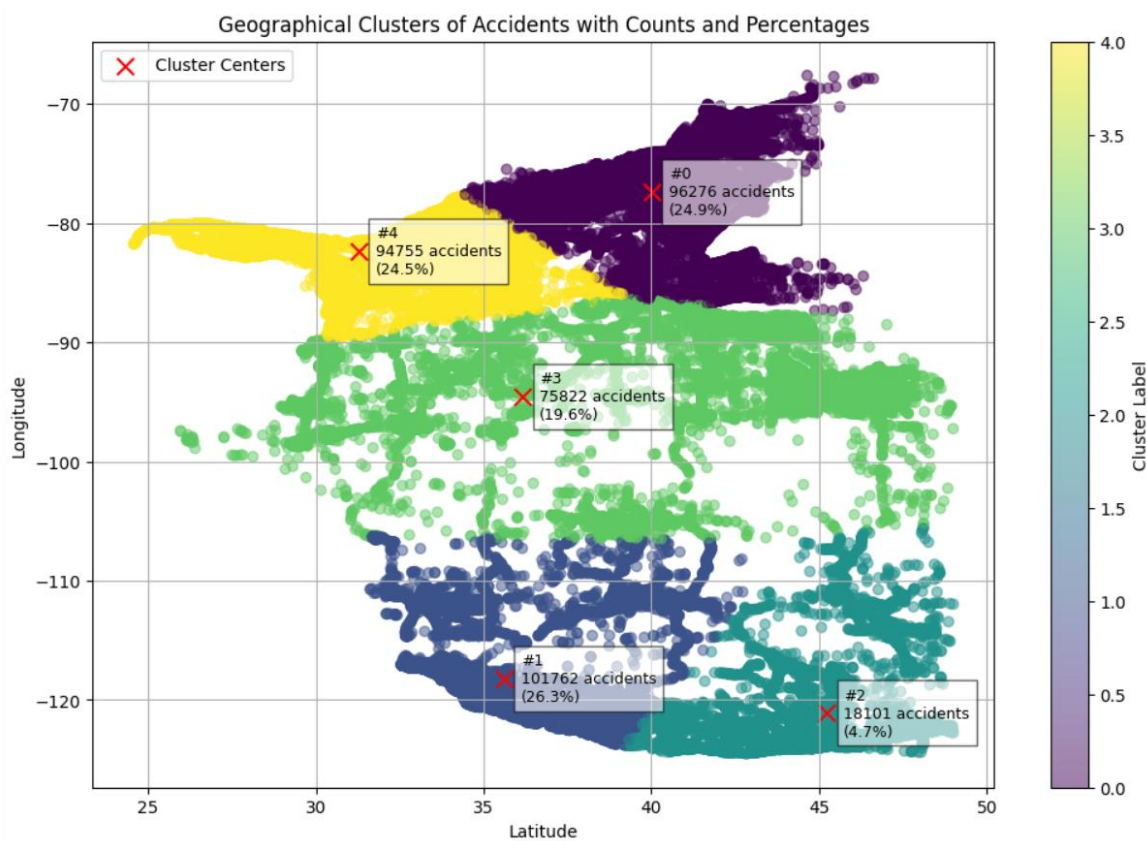
Accidents based on Location

Top 10 Cities with Highest Accident Counts

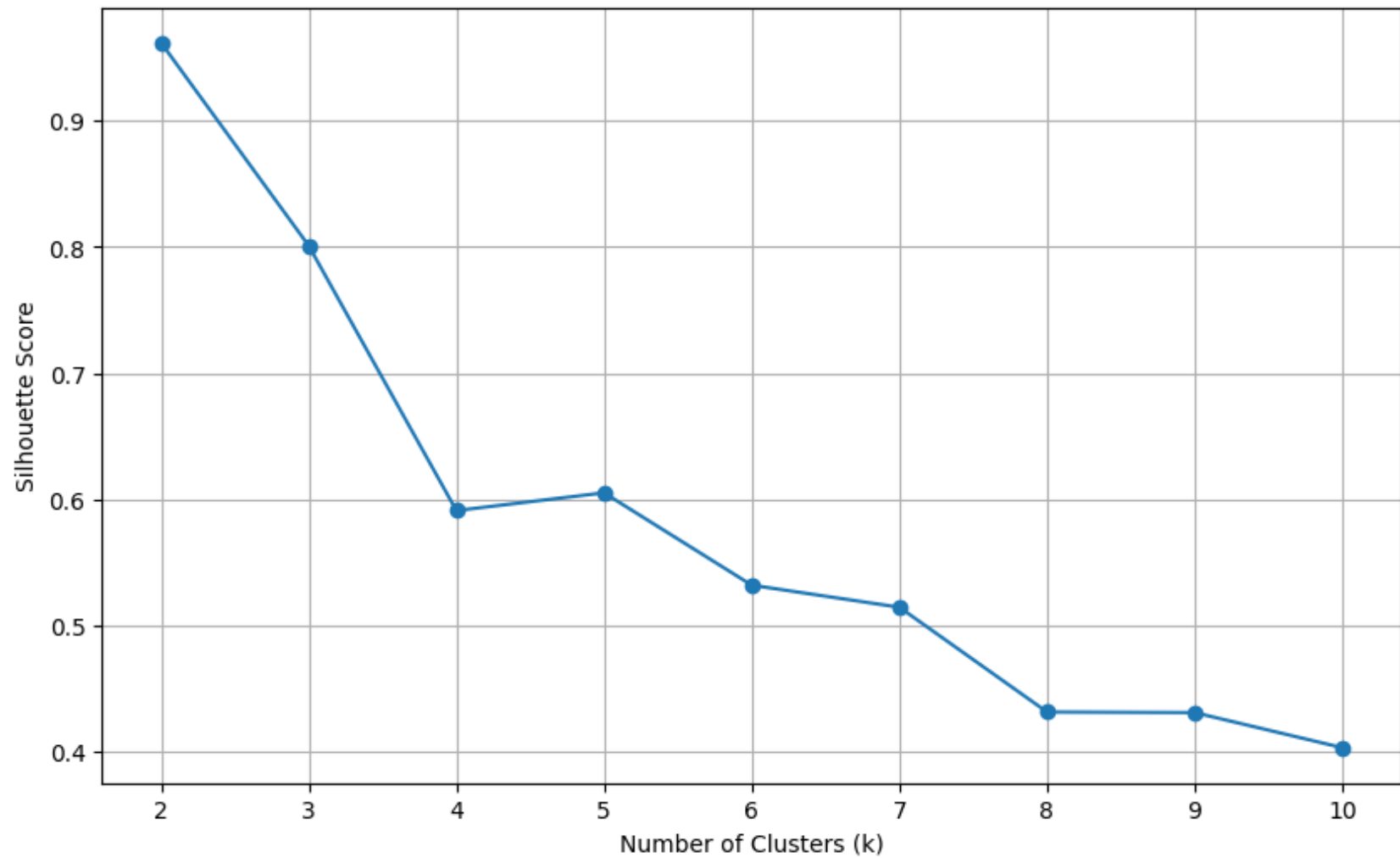




Top 5 Clusters in US with Accidents

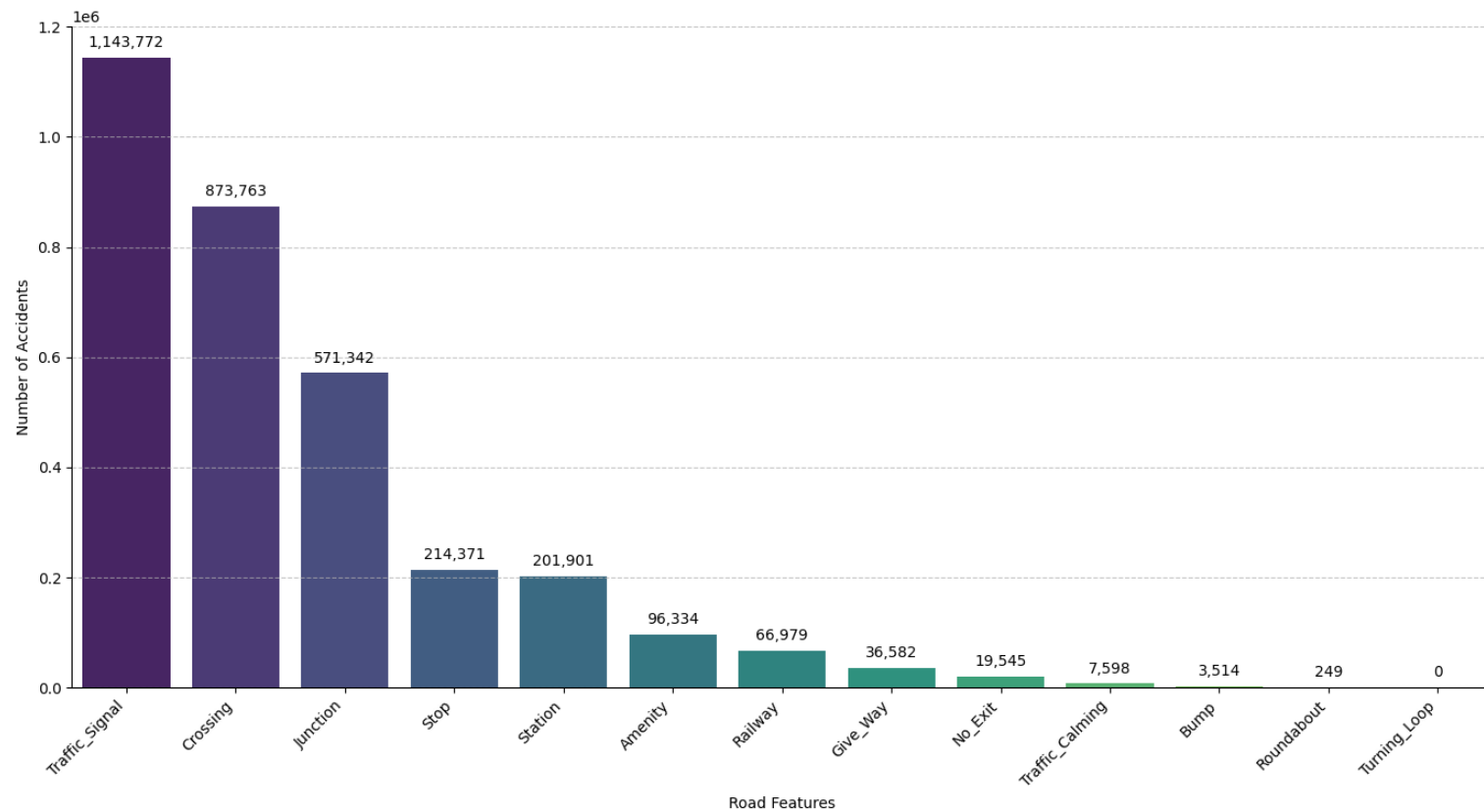


Silhouette Score vs. Number of Clusters

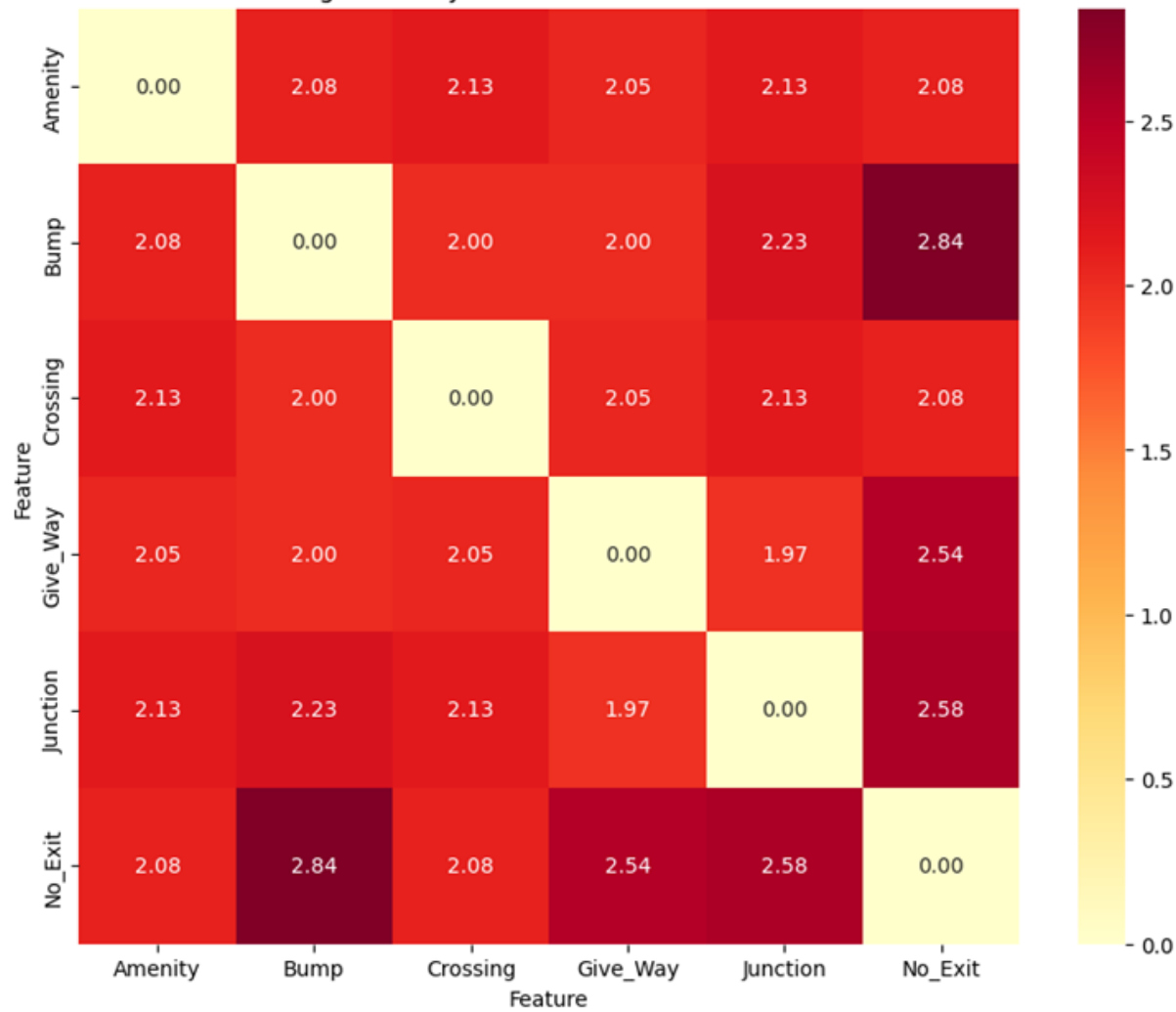


Road Conditions & Accidents

Accidents Associated with Road Features

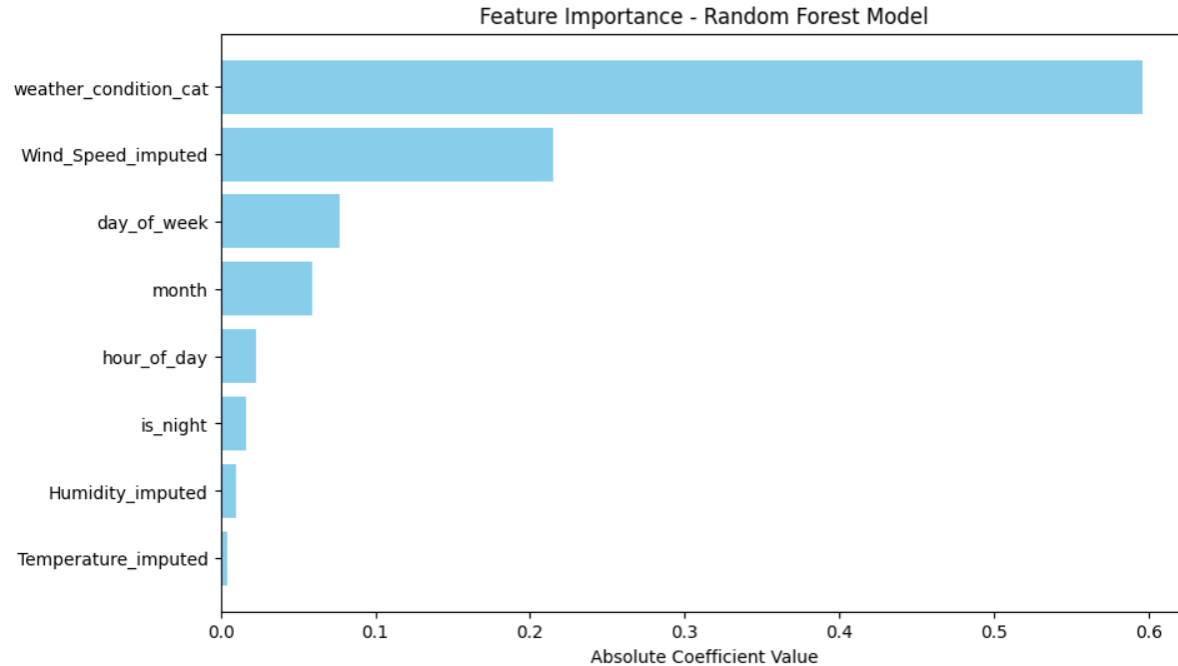


Average Severity When Road Features Co-occur



Predicting Severity & Risk

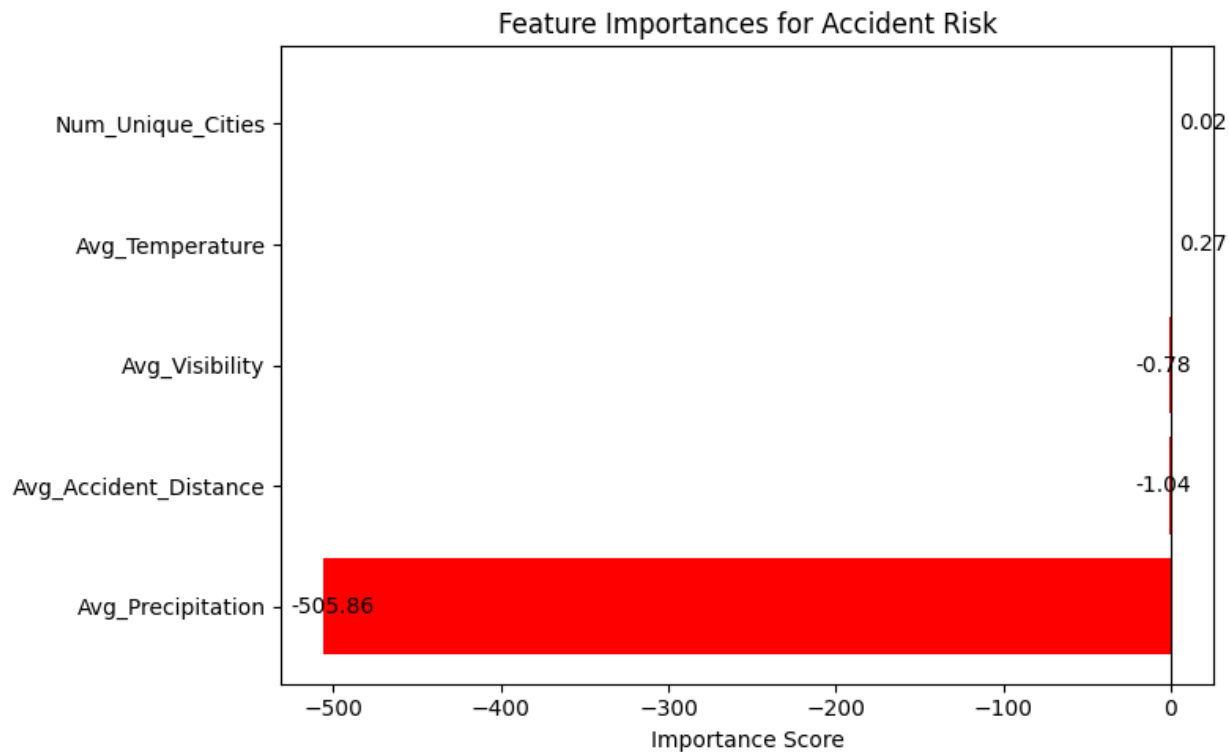
Multiclass Classification



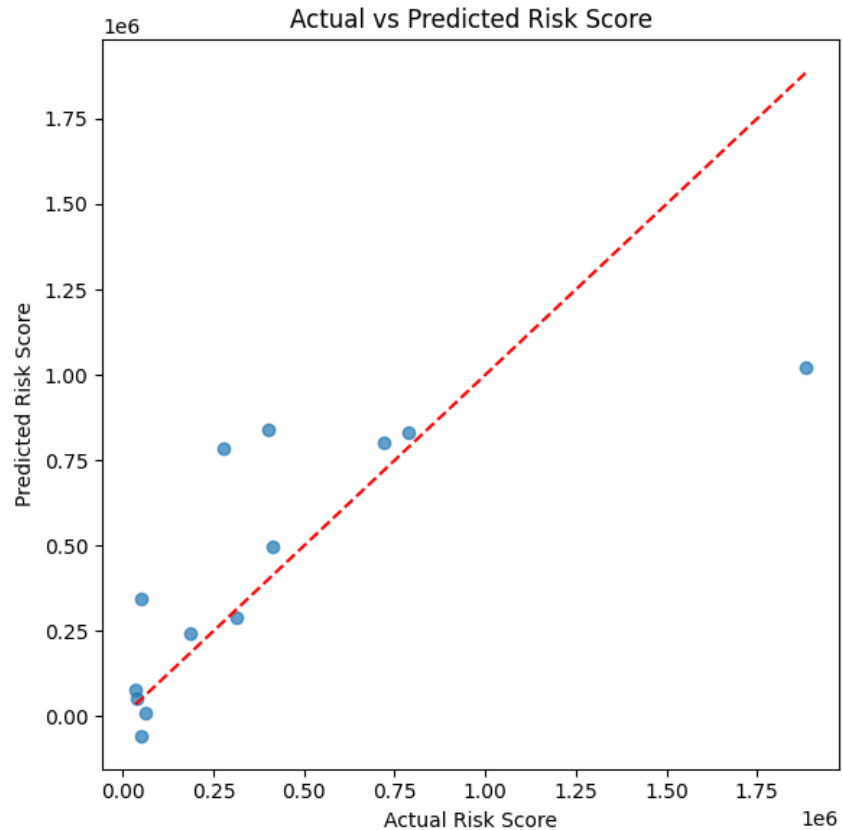
Binary Classification for Severity

Model	AUC	F1	Accuracy
Logistic Regression	0.7297	0.7492	0.8077
Random Forest	0.7509	0.7184	0.8053
Gradient-Boosted Trees	0.8089	0.7660	0.8192
Decision Tree	0.6452	0.7438	0.8112

Risk Classification for States



Risk Prediction using Regression



Azure Databricks

The screenshot shows the Azure Databricks Main Notebook interface. The top bar includes the Databricks logo, a search bar, and the text "bigdata-databricks-workspace". The notebook is titled "Main Notebook" and is in "Python" mode. The left sidebar shows the "Workspace" view with a file tree containing "Main Notebook" and "us_accidents_sample.csv". The main area displays a code editor with the following Python code:

```
08:49 PM (<1s) 1

StructField("Junction", BooleanType(), True),
StructField("No_Exit", BooleanType(), True),
StructField("Railway", BooleanType(), True),
StructField("Roundabout", BooleanType(), True),
StructField("Station", BooleanType(), True),
StructField("Stop", BooleanType(), True),
StructField("Traffic_Calming", BooleanType(), True),
StructField("Traffic_Signal", BooleanType(), True),
StructField("Turning_Loop", BooleanType(), True),

StructField("Sunrise_Sunset", StringType(), True),
StructField("Civil_Twilight", StringType(), True),
StructField("Nautical_Twilight", StringType(), True),
StructField("Astronomical_Twilight", StringType(), True)
})
return spark.read.option("header", "true").schema(schema).csv("file:///"+os.getcwd()+"/us_accidents_sample.csv")
```

▼ (3) Spark Jobs

- ▶ Job 1  [View](#) (1 stage)
- ▶ Job 2  [View](#) (1 stage)
- ▶ Job 3  [View](#) (1 stage)

gettig accuracy

Test RDD has 77081 elements and 4 partitions.

Fully Distributed Mode

Compute > New compute > Simple form: OFF

Fully-Distributed Cluster

Policy: Unrestricted

☒ Multi node ☐ Single node

Access mode: Dedicated (formerly: Single user) Single user or group access: Ahmed Osama

Performance

Databricks runtime version: Runtime: 15.4 LTS (Scala 2.12, Spark 3.5.0)

☒ Use Photon Acceleration

Worker type: Standard_DC4as_v5 16 GB Memory, 4 Cores Workers: 2 ☒ Spot instances

Driver type: Same as worker 16 GB Memory, 4 Cores

☐ Enable autoscaling

☒ Terminate after 120 minutes of inactivity

Tags

Add tags

Key	Value	Add

> Automatically added tags

Advanced options

Summary

2 Workers

32 GB Memory

8 Cores

1 Driver

16 GB Memory, 4 Cores

Runtime

15.4.x-scala2.12

Unity Catalog

Photon

Standard_DC4as_v5

6 DBU/h