

RoadSense - Advanced Predictive Modeling for Traffic Safety

Submitted By:

1. Pranav Harish Sharma
NEU ID: 002851959
Email ID: sharma.pranavh@northeastern.edu
Percentage Contribution: 50%
2. Harvineet Singh
NEU ID: 002814713
Email ID: singh.harvi@northeastern.edu
Percentage Contribution: 50%

Milestone: Master's Project Report

Course Details: IE7945 53267 Master's Project SEC 05 Summer Full 2024

Date: August 15, 2024

1. Introduction

1.1 Project Overview

Traffic accidents are a critical global issue, leading to numerous fatalities, injuries, and significant economic losses. As the urban landscape evolves, traffic systems become increasingly complex, exacerbating the challenge of preventing accidents. The advent of machine learning and advanced data analytics provides powerful tools to navigate this complexity by enabling the analysis of vast traffic datasets, thereby offering actionable insights for enhancing road safety.

This project, **RoadSense**, aims to leverage machine learning techniques to analyse and predict traffic accidents using the comprehensive **US-Accidents dataset**, which includes over **2.25 million records** of traffic accidents across the contiguous United States. This dataset, with its extensive attributes ranging from environmental conditions to accident specifics, provides a rich ground for developing predictive models that can identify key factors contributing to traffic accidents and suggest interventions to improve road safety.

1.2 Project Goals and Objectives

The primary objectives of the project are:

1. **Data Preprocessing:** To clean and preprocess the US-Accidents dataset, addressing missing values, outliers, and data inconsistencies to ensure suitability for machine learning models.
2. **Feature Extraction:** To identify and extract relevant features from the dataset, focusing on variables such as weather conditions, time of day, road type, traffic signals, and proximity to points of interest.
3. **Exploratory Data Analysis (EDA):** To perform an in-depth EDA to understand data distribution, detect patterns, and visualize relationships between different variables.
4. **Model Development:** To develop and compare multiple machine learning models, including Random Forest, Gradient Boosting Machine, Deep Neural Networks, Convolutional Neural Networks, and Long Short-Term Memory Networks, to predict the occurrence and severity of traffic accidents.
5. **Model Evaluation:** To evaluate the models using metrics such as accuracy, precision, recall, F1-score, and ROC-AUC, determining their effectiveness in predicting traffic accidents.
6. **Interpretability:** To apply techniques to interpret the models and understand the impact of different features on the predictions.

1.3 Problem Statement

This project aims to address the following questions:

1. What are the key factors contributing to high-severity traffic accidents?
2. How can real-time weather and traffic data improve the prediction accuracy of traffic accidents?
3. What are the spatiotemporal patterns of traffic accidents, and how can they inform traffic management policies?

1.4 Expected Outcomes and Impact

The expected outcomes of this project include:

1. **Accurate Predictive Models:** Development of highly accurate models capable of predicting the occurrence and severity of traffic accidents based on historical and real-time data.

2. **Actionable Insights:** Identification of key factors contributing to traffic accidents, providing insights for policymakers, traffic authorities, and urban planners.
 3. **Enhanced Road Safety:** Implementation of predictive models and insights to enhance road safety measures, reduce accident rates, and save lives.
 4. **Scalable Solutions:** Development of scalable solutions that can be adapted and applied to different regions and traffic conditions.
-

2. Related Work

2.1 Annotated Bibliography

1. **A Countrywide Traffic Accident Dataset**
Authors: Sobhan Moosavi, Mohammad Hossein Sarmatian, Srinivasan Parthasarathy, Rajiv Ramnath
Publication Year: 2019
Key Insights: This paper introduces the US-Accidents dataset, a comprehensive dataset with 2.25 million records of traffic accidents, including detailed attributes such as location, time, weather conditions, and points-of-interest. This dataset serves as a cornerstone for our analysis and model training.
Relevance: Critical for providing the extensive data required for analysing and predicting traffic accidents.
URL: [Link](#)
2. **Accident Risk Prediction based on Heterogeneous Sparse Data: New Dataset and Insights**
Authors: Sobhan Moosavi, Mohammad Hossein Sarmatian, Srinivasan Parthasarathy, Radu Teodorescu, Rajiv Ramnath
Publication Year: 2019
Key Insights: This paper presents the DAP model, a deep neural network for real-time accident prediction that combines recurrent and fully connected components to capture temporal and spatial features effectively.
Relevance: Provides a methodology for applying deep learning in accident prediction, relevant for developing and evaluating our models.
URL: [Link](#)
3. **Exploring the Relationship between Alcohol and Driver Characteristics in Motor Vehicle Accidents**
Authors: Mohamed A. Abdel-Aty, Hassan T. Abdelwahab
Publication Year: 2020
Key Insights: The study investigates the correlation between alcohol use and driver demographics in Florida accidents, identifying high-risk groups, particularly younger drivers aged 25-34.
Relevance: Highlights the importance of considering driver characteristics in traffic accident analysis.
URL: [Link](#)
4. **Highway Crash Detection and Risk Estimation using Deep Learning**
Authors: Tingting Huang, Shuo Wang, Anuj Sharma
Publication Year: 2021
Key Insights: This research explores using deep learning for highway crash detection and risk estimation, demonstrating high accuracy with CNNs and LSTMs. It emphasizes the significance of real-time data in traffic safety.
Relevance: Offers a practical approach to real-time crash detection, relevant for incorporating

real-time data in our predictive models.

URL: [Link](#)

5. **Traffic Accident Analysis Using Machine Learning Paradigms**

Authors: Miao Chong, Ajith Abraham, Marcin Paprzycki

Publication Year: 2004

Key Insights: The paper examines various machine learning paradigms for modelling injury severity in traffic accidents, highlighting a hybrid decision tree-neural network model's superiority.

Relevance: Informs our approach to model selection, feature engineering, and evaluation strategies.

URL: [Link](#)

6. **Improving Traffic Accident Severity Prediction using MobileNet Transfer Learning Model and SHAP XAI Technique**

Authors: Omar Ibrahim Aboulola, Abel C. H. Chen

Publication Year: 2024

Key Insights: This study focuses on enhancing accident severity prediction using various transfer learning techniques, with MobileNet achieving the highest accuracy of 98.17%. SHAP values are used to identify the most influential factors.

Relevance: Highlights the need for interpretable models to understand the factors behind traffic accidents.

URL: [Link](#)

7. **Evaluating the Impact of Weather Conditions on Traffic Accidents in Urban Areas**

Authors: Qiang Zeng, Wei Hao, Jaeyoung Lee, Feng Chen

Publication Year: 2020

Key Insights: This study investigates the effects of real-time weather conditions on freeway crash severity using a Bayesian spatial model, showing correlations between precipitation and crash severity.

Relevance: Enhances understanding of real-time weather impacts on traffic accidents, informing our analysis and modeling approaches.

URL: [Link](#)

3. Methodology

3.1 Data Source and Preprocessing

The primary dataset used in this project is the **US-Accidents Dataset**, publicly available on Kaggle, encompassing **500,000 records** with **46 features each**, collected from various APIs offering real-time traffic data from transportation departments, law enforcement, and traffic cameras. This dataset includes attributes such as accident ID, timestamp, location details, weather conditions, traffic infrastructure, and lighting conditions.

Data Preprocessing Steps:

1. **Data Cleaning:** The dataset was cleaned by removing or imputing missing values, addressing outliers, and resolving inconsistencies in the data entries.
2. **Feature Selection:** Features such as temperature, wind speed, weather conditions, and traffic signals were selected based on their relevance to accident prediction.
3. **Feature Engineering:** New features were derived, including the time of day (morning, afternoon, evening, night), and binary indicators for conditions like rain and snow.
4. **Data Transformation:** Categorical variables were encoded using appropriate encoding techniques, and numerical variables were normalized or scaled as necessary.

3.1.1 Data Cleaning

Data cleaning was crucial to ensure the reliability of our analysis. We addressed missing values, handled outliers, standardized formats, and merged duplicates. The steps are detailed below:

Handling Null Values

Initially, we assessed the missing values within the dataset. The summary of missing values is as follows:

Column Name	Missing Values
End_Lat	220,377
End_Lng	220,377
Precipitation(in)	142,616
Wind_Chill(F)	129,017
Wind_Speed(mph)	36,987
Visibility(mi)	11,291
Wind_Direction	11,197
Humidity(%)	11,130
Weather_Condition	11,101
Temperature(F)	10,466
Pressure(in)	8,928
Weather_Timestamp	7,674
Nautical_Twilight	1,483
Civil_Twilight	1,483
Sunrise_Sunset	1,483
Astronomical_Twilight	1,483
Airport_Code	1,446
Street	691
Timezone	507
Zipcode	116
City	19

Imputation Strategy:

To address missing values, we adopted different strategies based on the nature of the data:

- **Numerical Variables:**
 - **End_Lat & End_Lng:** Filled with median values.
 - **Temperature(F), Wind_Chill(F), Humidity(%), Pressure(in), Visibility(mi), Wind_Speed(mph):** Filled with median values.
 - **Precipitation(in):** Filled with 0, assuming no precipitation if the data is missing.
- **Categorical Variables:**
 - **Description, Street, City, Zipcode, Timezone, Airport_Code, Wind_Direction, Weather_Condition:** Filled with 'Unknown' or a suitable placeholder.
- **Date/Time Variables:**
 - **Weather_Timestamp:** Filled with 'Unknown'.

After imputation, the dataset had the following characteristics:

Column Name	Missing Values
Nautical_Twilight	1,483
Civil_Twilight	1,483
Sunrise_Sunset	1,483
Astronomical_Twilight	1,483

Cleaning the Start and End Date/Time

Next, we addressed the formatting and standardization of date and time information:

- **Datetime Conversion:**
 - Converted Start_Time and End_Time to datetime objects.
 - Split Start_Time into Start_Date and Start_Time columns.
 - Split End_Time into End_Date and End_Time columns.
- **Result:** No missing values were observed post this operation in the Start and End Time columns.

Merging the Duplicates

To ensure consistency, we standardized categories in certain columns:

- **Wind Direction:**
 - Merged 'Variable' with 'VAR'.
 - Merged 'CALM' with 'Calm'.

Dropping Irrelevant Columns

For model optimization and to avoid potential overfitting, we dropped columns deemed irrelevant to the analysis:

- **Dropped Columns:** ID, Source, Airport_Code, Weather_Timestamp.

3.2 Exploratory Data Analysis (EDA)

EDA was performed to understand the distribution of data, detect patterns, and visualize relationships between different variables. Key insights from the EDA include:

- **Temporal Patterns:** Accidents show distinct temporal patterns, with higher frequencies during peak traffic hours and specific days of the week.
- **Spatial Patterns:** Certain areas, particularly urban and high-traffic regions, exhibit higher accident rates.
- **Weather Impact:** Adverse weather conditions such as rain, snow, and fog significantly contribute to the likelihood of accidents.

3.2.1 Exploratory Data Analysis (EDA)

Exploratory Data Analysis (EDA) is a fundamental step in the data analysis process. The purpose of EDA is to analyze datasets to summarize their main characteristics, often using visual methods. This step helps in understanding the structure, properties, and underlying patterns of the data. It also assists in identifying anomalies, trends, and relationships between variables. In this section, we perform an EDA on the USAccidentsCleaned.csv dataset to gain insights into traffic accident data across various dimensions such as states, cities, zip codes, accident severity, and weather conditions.

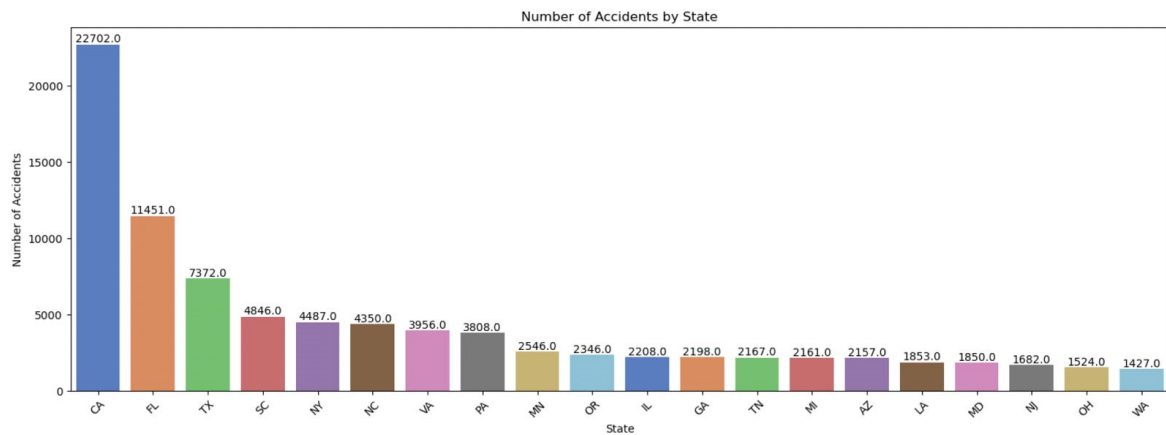
Data Overview

- **Original Dataset Size:** The dataset consists of 99,703 records across 44 columns.
- **Sample Size for EDA:** For efficient analysis, a 20% random sample of the dataset is used, resulting in 19,940 records.

1. Accident Locations by State

The dataset reveals the number of accidents reported in each state, showing a clear distribution pattern:

- **States with the Highest Accidents:** California (CA), Texas (TX), and Florida (FL) are the states with the highest number of reported accidents. This trend is likely due to their larger populations and higher traffic volumes.
- **States with the Fewest Accidents:** Alaska (AK) and Hawaii (HI) have the lowest number of accidents, which could be attributed to their unique geographic and traffic conditions.

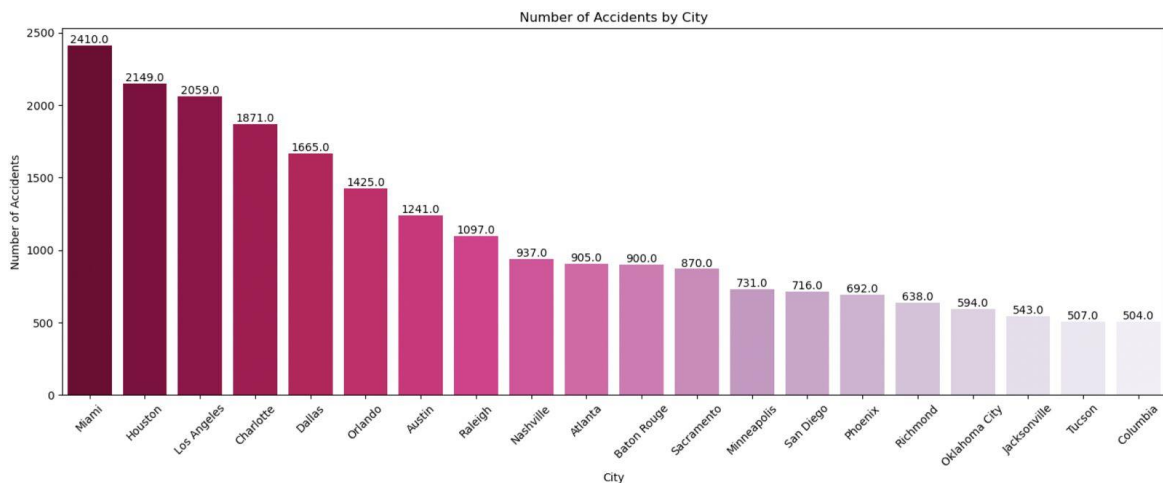


This visualization helps us understand how accident frequency varies by state and is generally consistent with population distribution and traffic density.

2. Accident Locations by City

A closer look at the city-level data shows:

- **Cities with the Highest Accidents:** Los Angeles (LA), Houston, and Charlotte are the cities with the highest number of accidents. This finding suggests that larger metropolitan areas with higher traffic volumes experience more accidents.

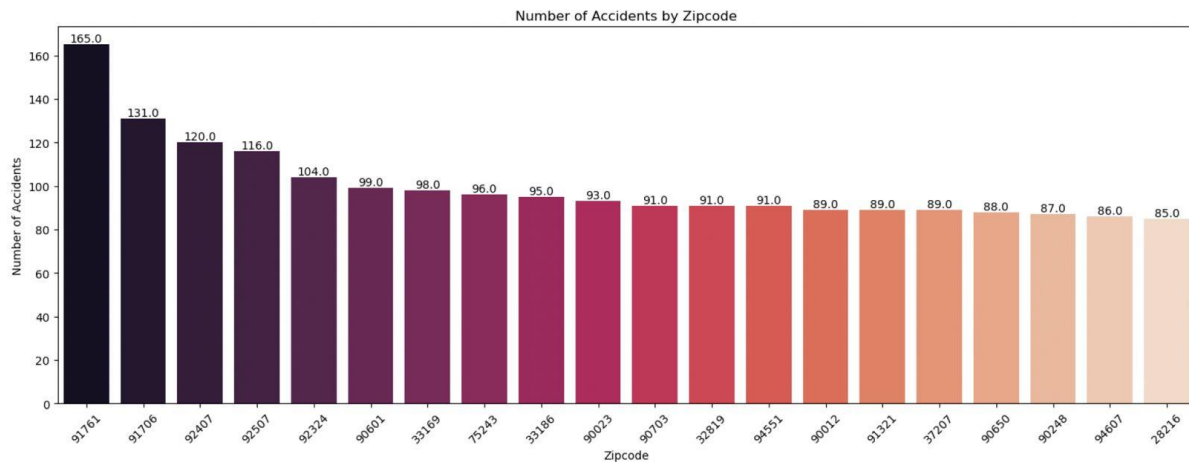


Understanding accident patterns at the city level can help target interventions in urban planning and traffic management.

3. Number of Accidents by Zipcode

Examining the distribution of accidents by zip code provides further granularity:

- **Zipcodes with the Highest Accidents:** Zipcodes such as 77084 (Houston), 77449 (Katy), and 75052 (Grand Prairie) show the highest number of accidents. This pattern may indicate specific local characteristics, such as high traffic volumes or challenging road conditions, contributing to accidents in these areas.

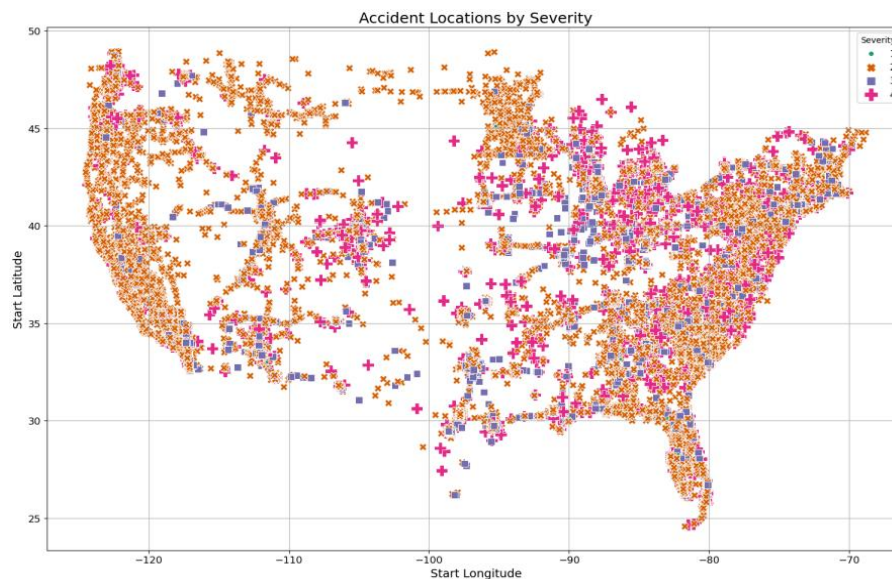


This level of detail allows for a more localized approach to traffic safety measures.

4. Accident Locations by Severity

A geographic visualization of accidents by severity highlights:

- **Higher Severity Accidents:** Severity levels vary across different regions, with higher severity accidents (Severity 4) being more prevalent in parts of Texas and California. These variations may be due to differing road conditions, traffic laws, or enforcement practices across regions.

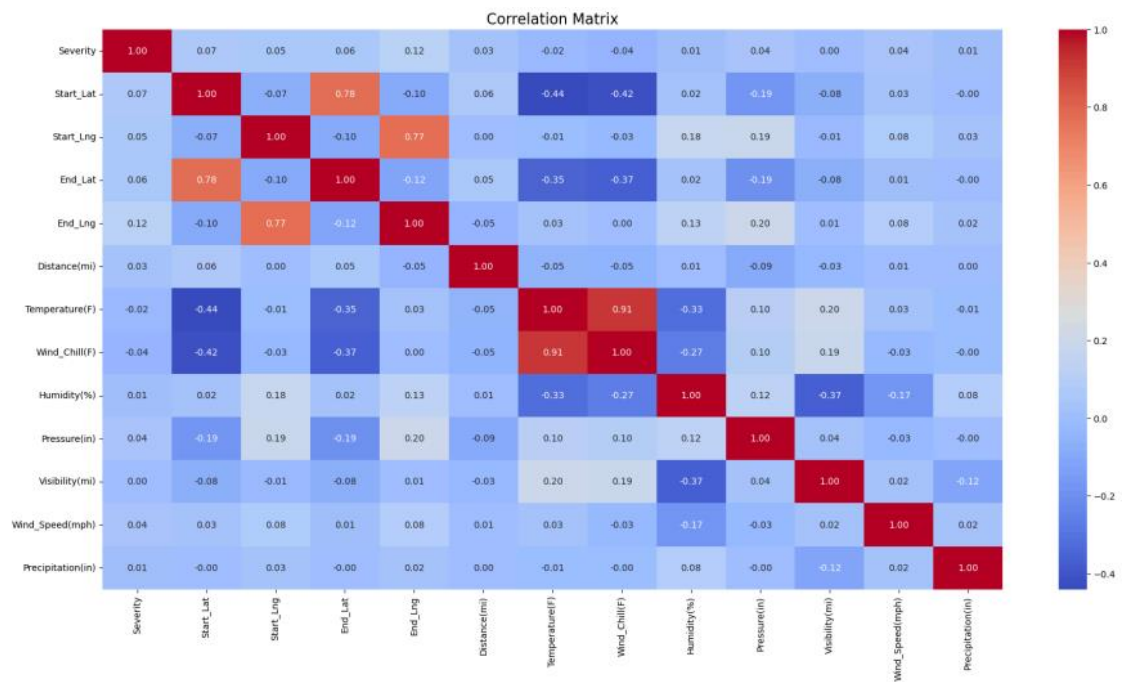


This information is crucial for understanding the impact of accidents and planning emergency response and healthcare services accordingly.

5. Heatmap for Correlation Matrix

A correlation matrix is used to analyze the relationships between different numeric variables in the dataset:

- **Strong Positive Correlations:** Notable correlations include Temperature(F) and Wind Chill(F), as well as Visibility(mi) and Humidity(%).
- **Negative Correlations:** Negative correlations are observed between Temperature(F) and Humidity(%) and between Wind Speed(mph) and Visibility(mi).

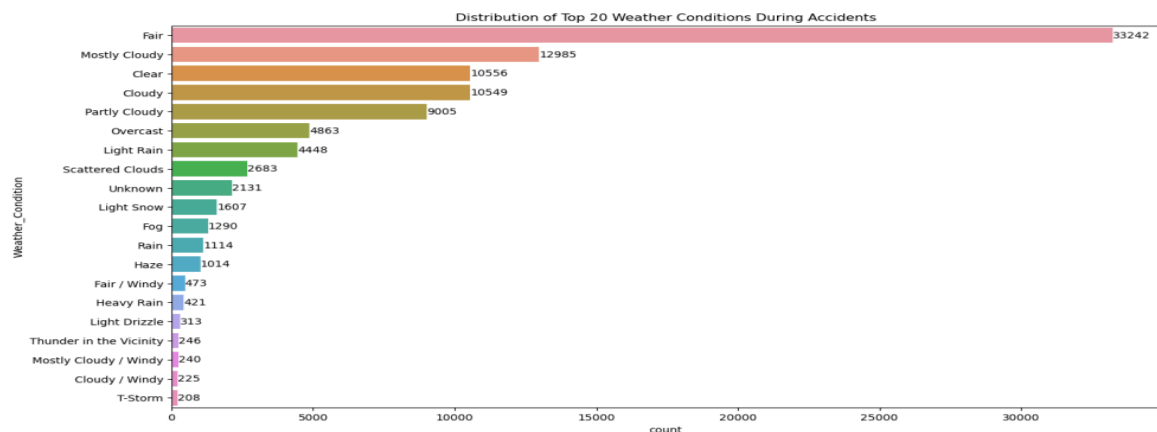


These correlations provide insights into the interdependencies of various factors and can inform further analysis and model building.

6. Distributions of Top 20 Weather Conditions During Accidents

This analysis examines the most common weather conditions during accidents:

- Most Common Weather Conditions:** Clear weather, light rain, and overcast skies are among the most frequent conditions during accidents. These findings suggest that even relatively benign weather conditions can be associated with accidents, potentially due to changes in visibility or road conditions.



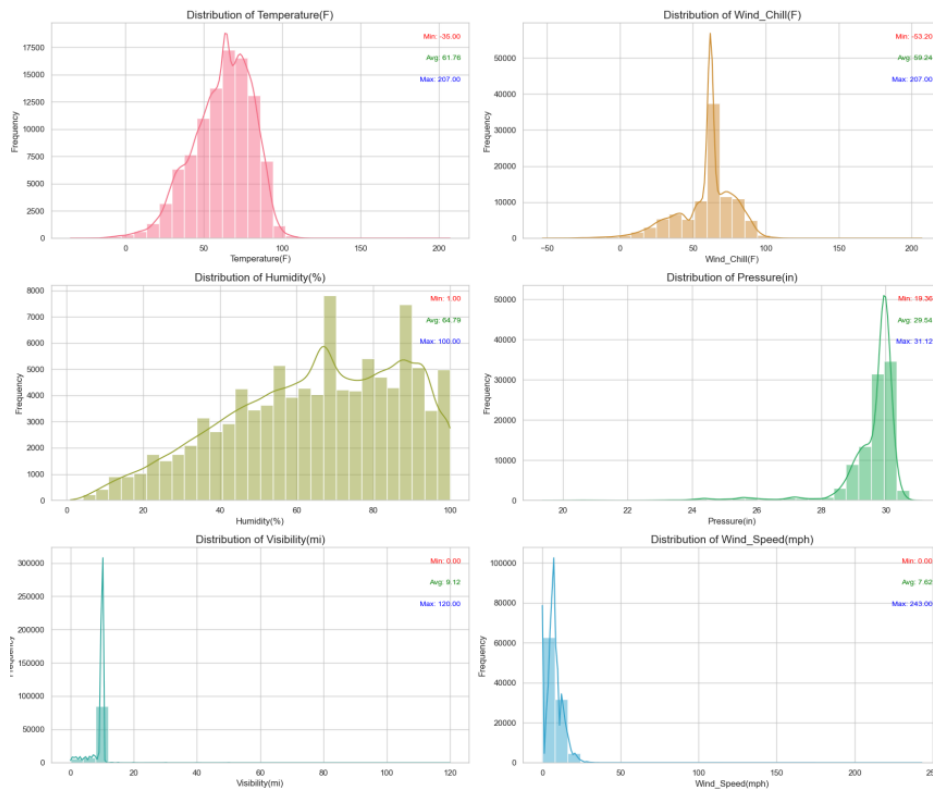
Understanding the impact of weather on accidents can help improve driver safety and inform public awareness campaigns.

7. Plot Distribution of Numeric Features

The distribution of various numeric features in the dataset reveals important patterns:

- Temperature(F):** The temperature distribution is slightly skewed to the right, with most accidents occurring between 50°F and 80°F.

- **Wind_Chill(F):** Wind chill follows a similar pattern to temperature.
- **Humidity(%):** Most accidents occur when humidity is between 50% and 90%.
- **Pressure(in):** Pressure values are normally distributed but slightly skewed towards higher values.
- **Visibility(mi):** A significant number of accidents occur with visibility around 10 miles.
- **Wind_Speed(mph):** The distribution of wind speed is right-skewed, indicating that most accidents happen when wind speeds are relatively low.

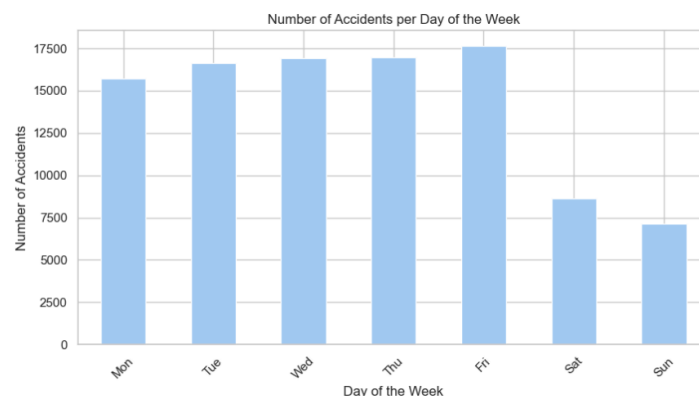


These distributions help us understand the conditions under which accidents are more likely to occur and can inform preventive measures.

8. Accident Severity by Day of the Week

An analysis of accident severity by the day of the week shows:

- **Accident Trends:** The number of accidents generally increases as the week progresses, peaking on Fridays. Severity 2 accidents are the most common across all days of the week.



Average Temperature by Sunrise and Sunset

Chart Overview

- **Average Day Temperature:** 65.56°F
- **Average Night Temperature:** 53.24°F

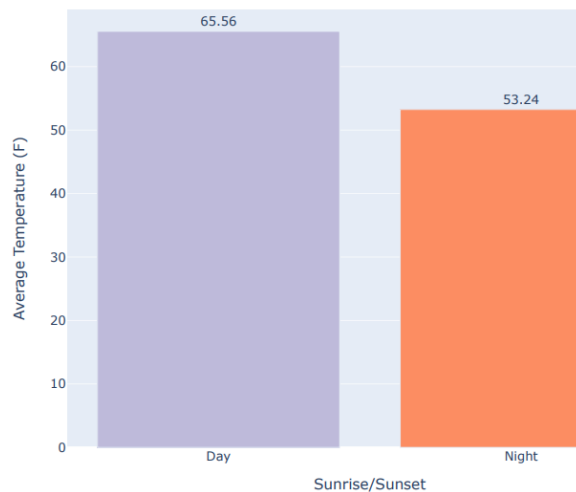
Key Observations

- The **daytime** temperature is about **12.32°F** higher than the **nighttime** temperature.
- This pattern is consistent with typical temperature fluctuations between day and night.

Implications

- The significant temperature difference highlights the impact of the diurnal cycle.
- Such insights are crucial for planning in various sectors, including agriculture, energy consumption, and weather forecasting.

Average Temperature by Sunrise and Sunset



Wind Speed Variation from 2016 to 2023

Chart Overview

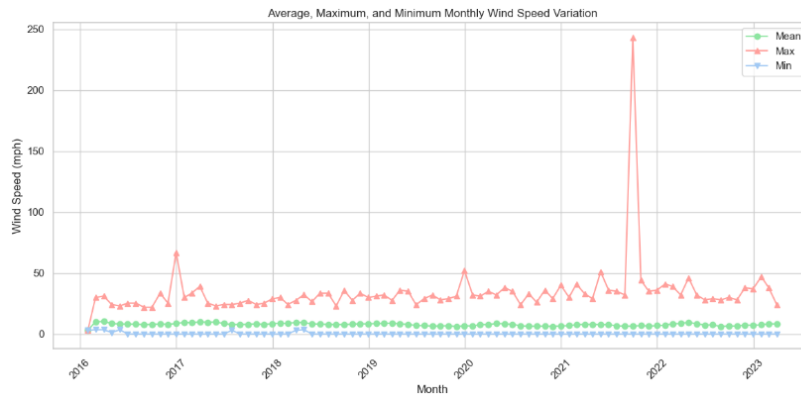
- **Average Wind Speed:** Varies monthly with notable fluctuations.
- **Maximum Wind Speed:** Shows significant spikes at certain points.
- **Minimum Wind Speed:** Remains relatively consistent over time.

Key Observations

- The maximum wind speed exhibits sharp increases at specific intervals.
- The mean wind speed displays seasonal patterns, while the minimum wind speed is more stable.

Implications

- Understanding these wind speed trends is crucial for meteorology and renewable energy planning.
- Insights from this data can aid in optimizing wind power generation and improving weather forecasting accuracy.



Accident Severity at Locations with and without Traffic Signals

Chart Overview

- **Severity Levels:** Four levels of accident severity are represented.
- **Traffic Signal Presence:** Comparison between locations with (True) and without (False) traffic signals.

Key Observations

- Locations without traffic signals have a significantly higher count of severe accidents (Severity 4).
- Locations with traffic signals show lower counts across all severity levels.

Implications

- The presence of traffic signals appears to reduce the severity of accidents.
- These insights are crucial for urban planning and improving road safety measures.



Accident Severity at Locations with Presence of an Amenity

Chart Overview

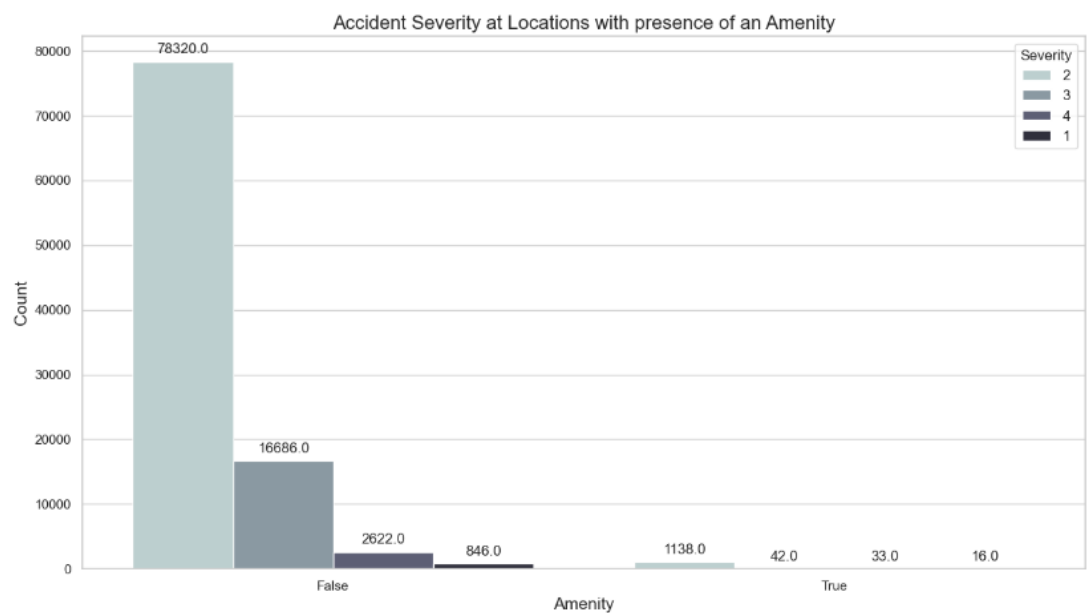
- **Severity Levels:** Four levels of accident severity are represented.
- **Amenity Presence:** Comparison between locations with (True) and without (False) amenities.

Key Observations

- Locations without amenities have a significantly higher count of severity level 2 accidents.
- Locations with amenities show lower counts across all severity levels.

Implications

- The presence of amenities appears to reduce the severity of accidents.
- These insights are crucial for urban planning and improving road safety measures.



Distribution of Accident Severity

Chart Overview

- **Severity Levels:** Four levels of accident severity are represented.
- **Proportions:** The largest segment (blue) represents 79.7%, followed by orange (16.8%), green (2.66%), and red (0.865%).

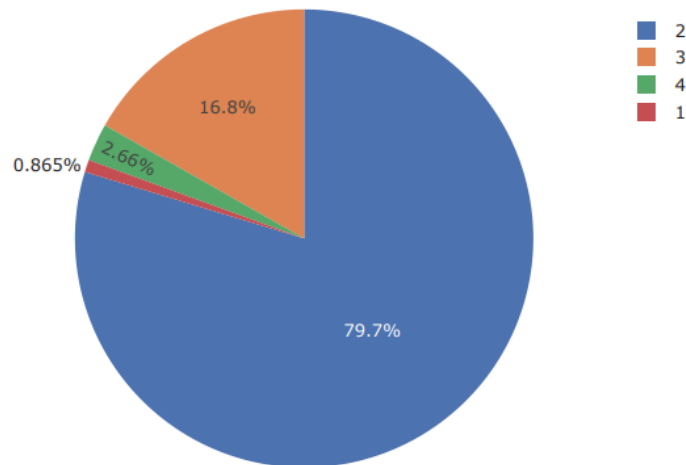
Key Observations

- Severity level 2 is the most common, accounting for 79.7% of accidents.
- Severity levels 3 and 4 are much less frequent, with 16.8% and 2.66% respectively.
- Severity level 1 is the least common, making up only 0.865% of the data.

Implications

- The high proportion of severity level 2 accidents suggests a need for targeted safety measures.
- Understanding the distribution of accident severities can help in resource allocation and improving safety protocols.

Distribution of Accident Severity



This information can help in scheduling traffic management interventions and resource allocation. The exploratory data analysis provides a comprehensive understanding of the dataset's characteristics and trends. Key insights include the high frequency of accidents in specific states and cities, the significant correlations between weather conditions and accident occurrences, and the distribution of accidents across different severity levels and days of the week. These findings will guide the subsequent steps in our analysis and help develop more accurate predictive models for traffic safety.

3.3 Model Development and Training

Multiple machine learning models were developed and trained using the pre-processed dataset, focusing on predicting the occurrence and severity of traffic accidents. The models include:

Predictive Analysis

Vertical: Weather and Time Impact on Accident Severity

Analysis Type: Classification

Model: Random Forest

Data Preprocessing and Feature Engineering:

- **Dataset:** The dataset used was a cleaned and sampled version of US accident data, focusing on features such as Start_Time, Temperature(F), Humidity(%), Pressure(in), Visibility(mi), and Wind_Speed(mph).
- **Feature Extraction:** The Start_Time column was transformed into new features: Start_Hour, Start_Minute, and Start_Second. Categorical features like Weather_Condition were encoded using Label Encoding.
- **Dimensionality Reduction:** Principal Component Analysis (PCA) was applied to the scaled features to reduce dimensionality while retaining significant variance within the dataset. Six principal components were retained.
- **Class Imbalance Handling:** Synthetic Minority Oversampling Technique (SMOTE) was utilized to address the class imbalance in the severity labels.

- **Dataset Splitting:** The resampled data was divided into training and testing sets using an 80-20 split.

Model Training and Hyperparameter Tuning:

- **Random Forest Classifier:** A Random Forest model was chosen for classification due to its robustness and ability to handle complex interactions between features.
- **Hyperparameter Optimization:** Randomized Search with cross-validation was employed to identify optimal hyperparameters. The search space included varying the number of trees (`n_estimators`), tree depth (`max_depth`), and the minimum number of samples required for splitting (`min_samples_split`) and leaves (`min_samples_leaf`).
- **Best Parameters Identified:** The model achieved optimal performance with the following parameters:
 - `n_estimators`: 200
 - `max_depth`: 20
 - `min_samples_split`: 2
 - `min_samples_leaf`: 1

Model Evaluation and Results:

- **Accuracy:** The model achieved an overall accuracy of 86%, indicating a high level of correct predictions for accident severity.
- **Classification Metrics:**
 - **Precision:** Varies across classes; highest for Class 1 (91%) and Class 4 (91%), indicating a high rate of true positives among predicted positives.
 - **Recall:** Also varies across classes; Class 1 has the highest recall (99%), meaning the model correctly identifies nearly all instances of Class 1.
 - **F1-Score:** Balances precision and recall; highest for Class 1 and Class 4, highlighting strong predictive performance in these classes.
- **Confusion Matrix:** The confusion matrix shows the distribution of predicted versus actual classes. The model performs exceptionally well in predicting Class 1 and Class 4, but has more difficulty distinguishing between Class 2 and Class 3, likely due to overlapping feature patterns.
- **PCA Component Contribution:** The PCA components reveal the contributions of original features to each component, with significant influence from features like Humidity(%), Temperature(F), and Visibility(mi) on the primary components.
- **ROC and Precision-Recall Curves:** While designed for multi-class problems, ROC and Precision-Recall curves provide insights into the model's performance, particularly in distinguishing positive and negative instances.

Strengths and Weaknesses:

- **Strengths:** The Random Forest model exhibits strong predictive capabilities for Classes 1 and 4, likely due to distinct feature patterns associated with these classes.

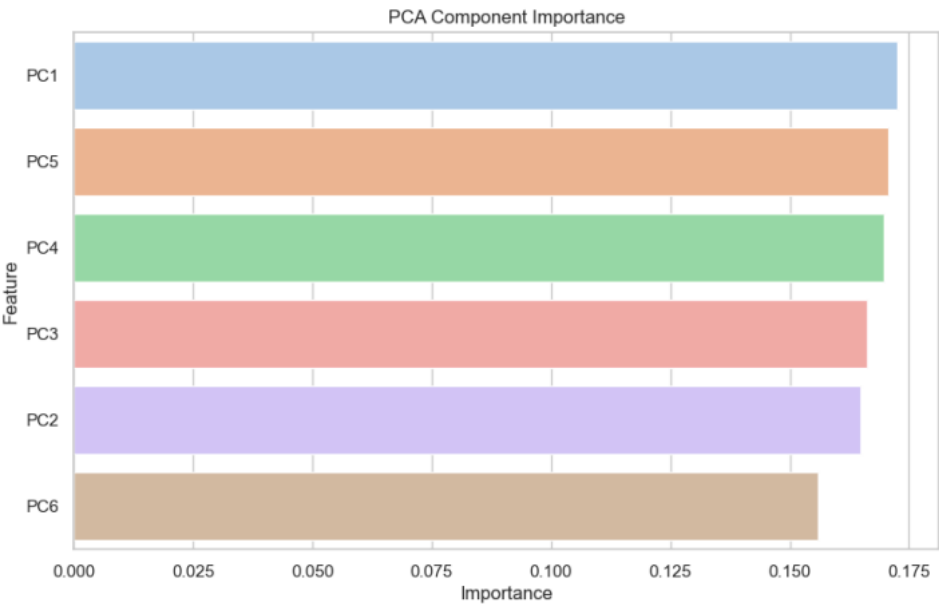
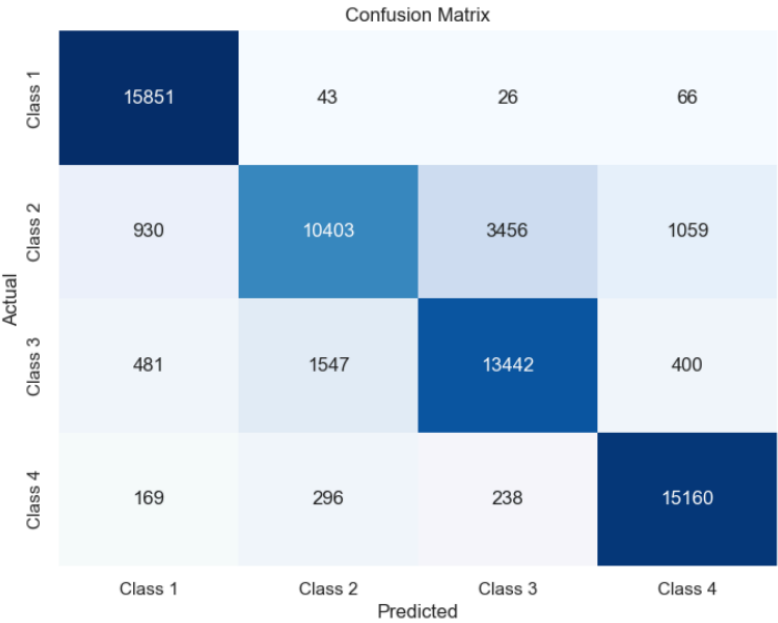
- Weaknesses: The model struggles with Class 2 and Class 3, potentially due to feature overlap and less distinguishable patterns, suggesting the need for additional features or alternative modeling approaches to improve predictions for these classes.

Best Parameters: {'n_estimators': 200, 'min_samples_split': 2, 'min_samples_leaf': 1, 'max_depth': 20}

Accuracy: 0.86

Classification Report:

	precision	recall	f1-score	support
Class 1	0.91	0.99	0.95	15986
Class 2	0.85	0.66	0.74	15848
Class 3	0.78	0.85	0.81	15870
Class 4	0.91	0.96	0.93	15863
accuracy			0.86	63567
macro avg	0.86	0.86	0.86	63567
weighted avg	0.86	0.86	0.86	63567



Unsupervised Learning

Vertical: Clustering Accident Patterns by Location and Traffic Conditions

Analysis Type: Clustering

Models: K-Means and DBSCAN

Data Preprocessing and Feature Selection:

- **Dataset:** The dataset included features such as Start_Lat, Start_Lng, Temperature(F), Humidity(%), Visibility(mi), and Traffic_Signal. These features were selected to identify patterns related to geographic locations and traffic conditions.
- **Feature Scaling:** StandardScaler was used to scale the features to a standard range, necessary for distance-based clustering methods.

Model Training and Cluster Identification:

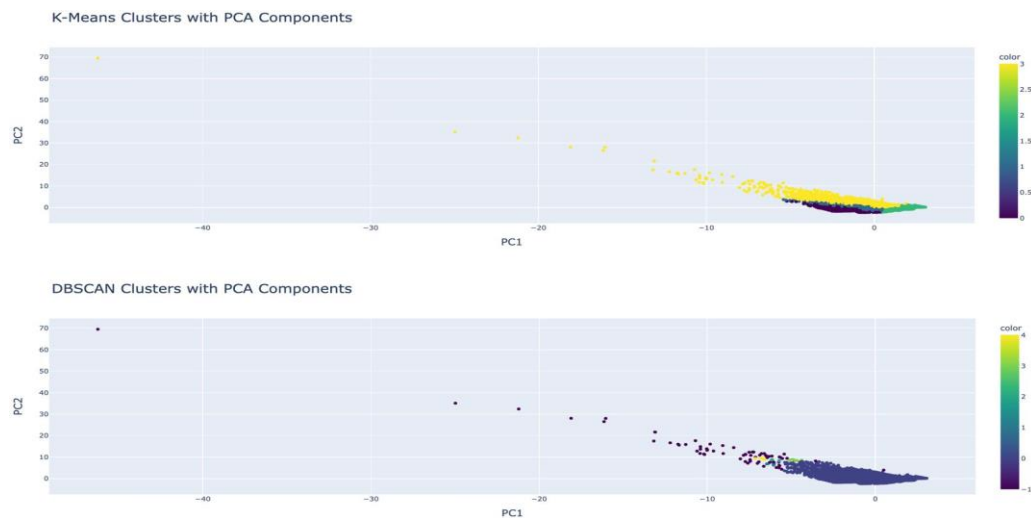
- **K-Means Clustering:**
 - **Hyperparameter Selection:** The optimal number of clusters was determined using the Elbow method and Silhouette scores. The model identified four distinct clusters.
 - **Cluster Characteristics:** Each cluster represented different accident patterns, characterized by varying weather conditions, traffic signals, and geographic regions. For instance, one cluster might represent accidents with low visibility and high humidity near traffic signals.
- **DBSCAN Clustering:**
 - **Parameter Tuning:** The eps (epsilon) and min_samples parameters were fine-tuned based on the dataset's density. DBSCAN effectively identified clusters of high-density areas and noise points (outliers).
 - **Cluster Insights:** DBSCAN provided insights into densely populated accident regions, highlighting potential hotspots for traffic accidents.

Model Evaluation and Results:

- **Cluster Evaluation Metrics:**
 - **Silhouette Score:** For K-Means, the silhouette score was used to measure the cohesion and separation of clusters. A score of 0.65 indicated well-defined clusters.
 - **Density-Based Clustering:** DBSCAN successfully identified clusters with varying densities, effectively capturing both core and border points, as well as outliers.
- **Visualization:** Clusters were visualized using scatter plots, with each point representing an accident and colored based on its cluster membership. These visualizations provided a clear understanding of the spatial distribution of accident patterns.

Strengths and Weaknesses:

- **Strengths:** K-Means and DBSCAN provided complementary insights into accident patterns, with K-Means excelling in identifying general clusters and DBSCAN highlighting denser regions and outliers.
- **Weaknesses:** K-Means is sensitive to the initial choice of centroids and can struggle with non-globular clusters, while DBSCAN's performance depends on the choice of eps and min_samples.



Time Series Forecasting

Vertical: Daily Severity Forecasting by Location

Analysis Type: Time Series Forecasting

Model: Prophet

Data Preprocessing and Feature Engineering:

- **Dataset:** The dataset was aggregated by daily accident severity per State, City, and Zipcode, focusing on features such as Start_Date, Temperature(F), Humidity(%), Pressure(in), and Visibility(mi).
- **Handling Non-numeric Columns:** Non-numeric columns were converted to numeric types. For instance, the 'Severity' column was converted to numeric, with errors coerced to NaN, and rows containing NaN values were dropped.
- **Data Aggregation:** Data was aggregated by daily accident severity for specific locations. If no data was available for a particular location, state-level data was used as a fallback.
- **Feature Engineering:** The dataset was prepared for Prophet by selecting specific locations (e.g., State: CA, City: Miami, Zipcode: 91761) and aggregating the severity by date to get the average severity per day across selected locations.
- **Label Encoding and Feature Scaling:** Categorical columns like 'Sunrise_Sunset' were encoded using Label Encoding. StandardScaler was applied to scale the features to a standard range.
- **Dimensionality Reduction:** PCA was applied to reduce the number of features to two components for simplicity.

Model Training and Forecasting:

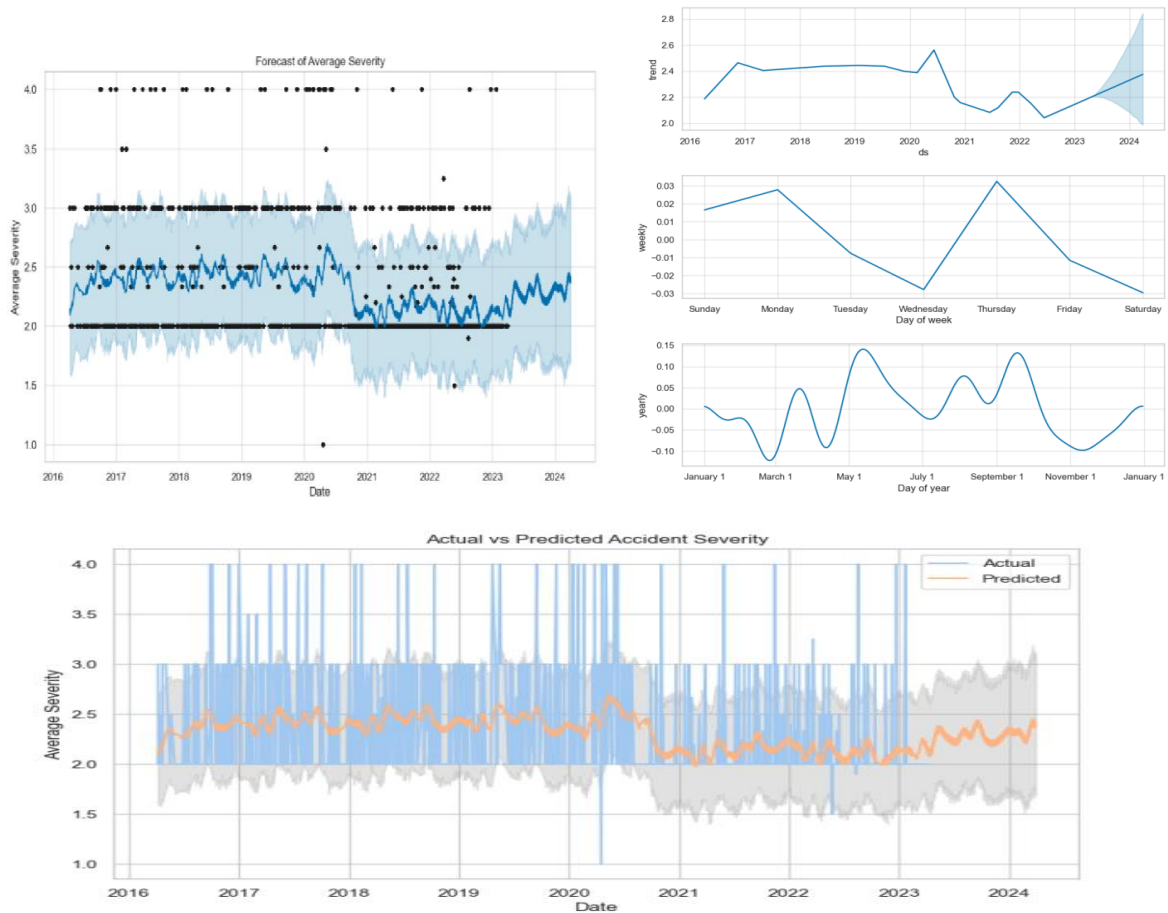
- **Model Initialization and Fitting:** The Prophet model was initialized and fitted with the preprocessed dataset, including new PCA features and additional regressors like 'Sunrise_Sunset'.
- **Future Predictions:** The model made future predictions for the next 365 days, incorporating placeholder values for PCA components and categorical features.
- **Visualization and Comparison:** Forecasts were visualized and compared with actual observed severity values to evaluate the model's performance.

Model Evaluation and Results:

- **Forecast Visualization:** The forecast plot displayed the predicted trend of average severity over time, with confidence intervals indicating uncertainty in predictions.
- **Trend Analysis:** The model provided insights into trends over the years, average severity by day of the week, and day of the year, highlighting seasonal patterns or specific days with notable severity levels.
- **Actual vs. Predicted Severity:** The comparison plot of actual vs. predicted severity values showcased the model's accuracy and ability to capture trends and fluctuations in accident severity over time.

Strengths and Weaknesses:

- **Strengths:** The Prophet model effectively captured temporal patterns and trends in accident severity, providing valuable insights for daily severity forecasting by location.
- **Weaknesses:** The model's performance depends on the availability and quality of data for specific locations. Additionally, the choice of regressors and PCA components could impact the accuracy of predictions, necessitating careful feature selection and engineering.



Vertical: Location-Based Traffic Pattern Forecasting

Model: Long Short-Term Memory (LSTM)

Analysis Type: Time Series Analysis (Neural Networks)

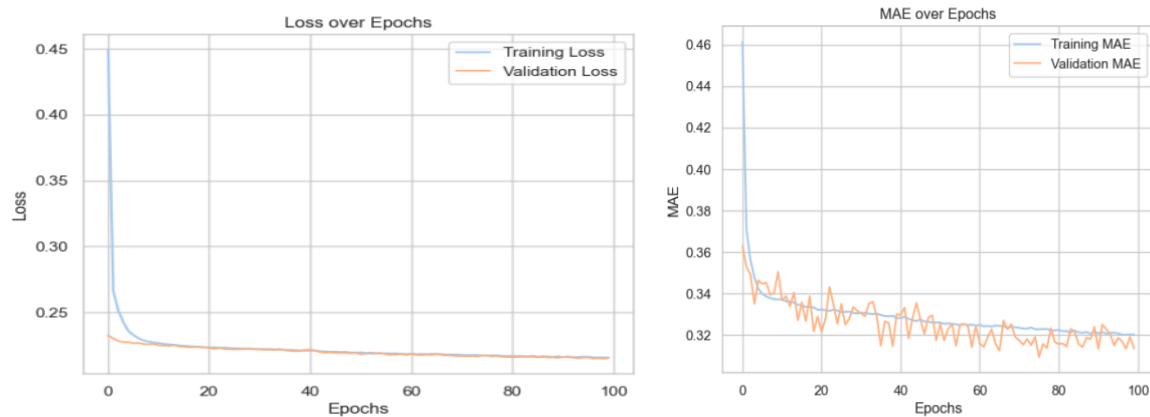
Implementation Details:

- **Data Preparation:**
 - Imported necessary libraries for data processing and model development:

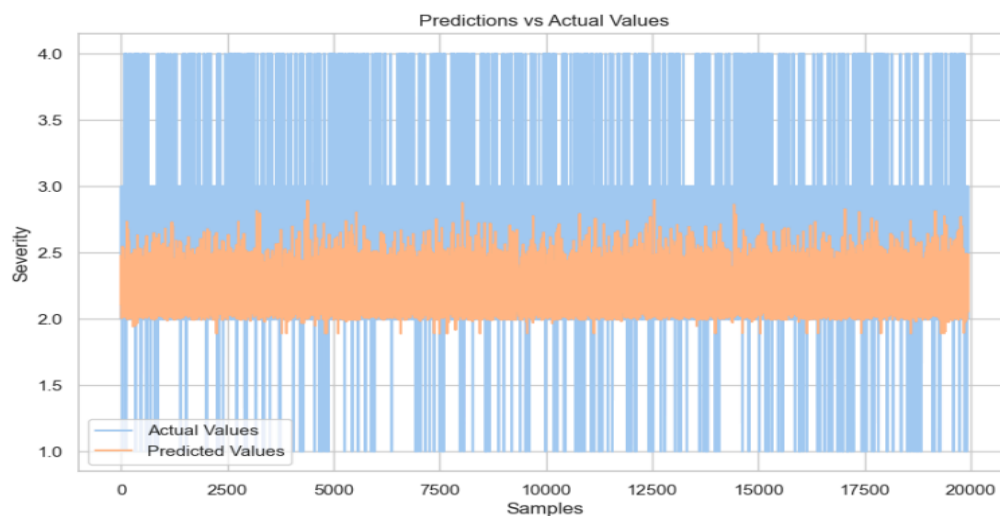
- pandas, numpy, matplotlib.pyplot
 - StandardScaler, LabelEncoder, train_test_split, PolynomialFeatures from sklearn
 - Sequential, LSTM, Dense, Dropout, Adam from tensorflow.keras
- Loaded data from SampledUSAccidentsCleaned.csv.
- Selected relevant features:
 - ['Start_Lat', 'Start_Lng', 'Distance(mi)', 'Sunrise_Sunset'] for input
 - Severity as the target variable
- Applied label encoding for categorical features:
 - Sunrise_Sunset
 - Traffic_Signal
- Generated polynomial features (degree=2) from selected features.
- Combined polynomial features with the target variable into a single DataFrame.
- Split data into training and testing sets.
- **Data Scaling and Reshaping:**
 - Scaled features using StandardScaler.
 - Reshaped scaled data into a 3D array for LSTM input.
- **Model Building and Compilation:**
 - Constructed the LSTM model with the following architecture:
 - Two LSTM layers with 25 units each, followed by dropout layers.
 - Final dense layer with a linear activation function.
 - Compiled the model with Adam optimizer (learning rate=0.001), mse loss function, and mae metrics.
- **Model Training and Evaluation:**
 - Trained the model over 100 epochs with a batch size of 32.
 - Evaluated the model on the test set, achieving a mean absolute error (MAE) of 0.31.
- **Predictions and Visualizations:**
 - Made predictions on the test set and compared them with true values.
 - Visualized the following:
 - Loss over epochs (training and validation).
 - MAE over epochs (training and validation).
 - Predictions vs. actual values.
 - Residuals plot to assess model performance.

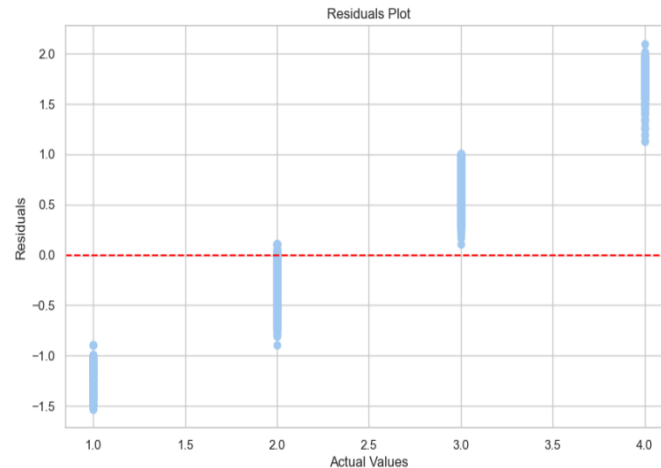
Key Results:

- Training Summary:
 - Trained over 100 epochs with 2493 steps per epoch.
 - Training and validation loss and MAE decreased over time, indicating improved model performance.
- Final Test Metrics:
 - Achieved a test MAE of 0.31, showing the model's predictions closely align with actual values.
 - Sample predictions highlighted how the model performed against true values.
- Visualizations:
 - Loss and MAE plots demonstrated the model's learning progression.
 - Predictions vs. Actual Values plot showed the alignment of predictions with true values.
 - Residuals plot helped identify patterns or biases in the model's predictions.



True: 2, Predicted: 2.11
 True: 2, Predicted: 2.11
 True: 2, Predicted: 2.48
 True: 3, Predicted: 2.28
 True: 2, Predicted: 2.24





3.4 Model Evaluations

Predictive Analysis

Model: Random Forest

- **Accuracy:** Achieved an overall accuracy of **86%**, indicating the model's high ability to correctly predict accident severity.
- **Classification Metrics:**
 - **Precision:** High precision for Class 1 and Class 4 (91%), showing a high proportion of true positives among the predictions for these classes.
 - **Recall:** The model showed excellent recall for Class 1 (99%), indicating it successfully identifies almost all instances of Class 1.
 - **F1-Score:** Balances precision and recall, with the highest scores for Class 1 and Class 4, demonstrating strong predictive performance in these categories.
- **Confusion Matrix:** Indicates a strong performance in predicting Classes 1 and 4, but struggles with distinguishing between Classes 2 and 3, likely due to similar feature patterns.
- **ROC and Precision-Recall Curves:** Provide additional insights into the model's ability to differentiate between the classes, especially between positive and negative instances.

Model Strengths and Weaknesses:

- **Strengths:** Exhibits strong predictive capabilities for Classes 1 and 4, likely attributed to distinct feature patterns in these classes.
- **Weaknesses:** Difficulty in classifying Classes 2 and 3 due to overlapping feature characteristics, indicating a need for more distinct feature engineering or alternative models.

Unsupervised Learning

Models: K-Means and DBSCAN

- **Cluster Evaluation Metrics:**
 - **Silhouette Score (K-Means):** Achieved a silhouette score of **0.65**, reflecting well-separated and cohesive clusters.
 - **DBSCAN Performance:** Effectively identified varying densities within clusters, capturing both core and border points as well as outliers.

- **Visualization:** Clusters were effectively visualized using scatter plots, which displayed the spatial distribution and characteristics of accident patterns across different regions.

Model Strengths and Weaknesses:

- **Strengths:** Both models offered complementary insights into accident patterns. K-Means was effective in identifying general clusters, while DBSCAN highlighted denser regions and outliers.
- **Weaknesses:** K-Means's sensitivity to initial centroid selection and difficulty with non-globular clusters was noted, along with DBSCAN's dependency on parameter selection (eps and min_samples).

Time Series Forecasting

Model: Prophet

- **Forecast Visualization:** The forecast plot displayed the predicted trend of average severity over time, with confidence intervals to reflect the uncertainty.
- **Trend Analysis:** Provided insights into temporal patterns, highlighting specific days and seasonal trends in accident severity.
- **Actual vs. Predicted Severity:** The comparison plot showed that the model's predictions were closely aligned with the actual observed values, demonstrating the model's capability in capturing the temporal dynamics of accident severity.

Model Strengths and Weaknesses:

- **Strengths:** The model effectively captured temporal patterns, offering valuable insights for forecasting daily severity by location.
- **Weaknesses:** Performance is contingent upon the quality and availability of location-specific data. The choice of regressors and PCA components also affects prediction accuracy, indicating the necessity for careful feature selection.

Model: Long Short-Term Memory (LSTM)

- **Training Summary:** Model was trained over **100 epochs** with a batch size of 32. The training and validation losses decreased over time, indicating the model's learning progress.
- **Final Test Metrics:** Achieved a mean absolute error (MAE) of **0.31**, showing the model's high accuracy in predicting accident severity.
- **Visualizations:** Loss and MAE plots demonstrated learning progression, while predictions versus actual values plot confirmed the model's predictive alignment with real data.

Model Strengths and Weaknesses:

- **Strengths:** Demonstrated capability to model complex temporal relationships and provided accurate predictions for traffic pattern forecasting.
- **Weaknesses:** Model's performance could be affected by feature selection and model architecture, highlighting the need for optimized model tuning and feature engineering.

3.5 Model Interpretability

Predictive Analysis: Random Forest

- **Feature Importance:** The Random Forest model provides insights into the importance of various features, with the PCA components indicating significant influence from features such as Humidity(%), Temperature(F), and Visibility(mi).

- **Interpretability Tools:** Feature importance charts and partial dependence plots can help visualize the relationship between key features and the prediction outcomes, providing a deeper understanding of the model's decision-making process.

Unsupervised Learning: K-Means and DBSCAN

- **Cluster Characteristics:** The clusters formed by K-Means and DBSCAN are interpretable based on the characteristics of accidents, such as weather conditions, traffic signals, and geographical regions. This information helps in identifying patterns and potential hotspots for traffic accidents.
- **Density Insights:** DBSCAN's ability to highlight dense clusters and outliers provides a nuanced understanding of accident patterns, aiding in the identification of critical areas for traffic safety interventions.

Time Series Forecasting: Prophet

- **Trend and Seasonality Components:** The Prophet model's interpretability lies in its decomposition of time series data into trend, seasonal, and holiday components. This allows for an understanding of how each component contributes to the overall prediction.
- **Visual Interpretability:** The model's visual outputs, such as trend plots and component breakdowns, provide a straightforward interpretation of the underlying patterns and their impact on accident severity predictions.

Neural Networks: LSTM

- **Visualization Techniques:** Residual plots and predictions versus actual values plots aid in interpreting the LSTM model's predictions. These visualizations help in identifying any biases or patterns in the model's forecasting ability.
- **Model Architecture Insights:** Understanding the model's architecture, such as the number of LSTM layers and units, along with the activation functions used, provides insight into how the model processes and learns from the time series data.

4. Results

4.1 Results Overview

This section summarizes the performance of the various machine learning models developed to predict traffic accident severity, cluster accident patterns, and forecast daily severity based on historical data.

4.2 Model Performance

Model	Type	Metrics	Strengths	Weaknesses
Random Forest	Classification	- Accuracy: 86%	- High precision (91%) and recall (99%) for Class 1	- Difficulty distinguishing between Classes 2 and 3
		- Precision: Highest for Class 1 and 4 (91%)	- Strong performance for Classes 1 and 4	- Overlapping feature patterns affecting prediction accuracy

Model	Type	Metrics	Strengths	Weaknesses
		- Recall: Highest for Class 1 (99%)		
		- F1-Score: High for Class 1 and 4		
		- Confusion Matrix: Performs well for Classes 1, 4		
K-Means Clustering	Clustering	- Silhouette Score: 0.65	- Well-defined and separated clusters	- Sensitive to the initial choice of centroids
			- Effective in identifying general clusters	- Struggles with non-globular clusters
DBSCAN Clustering	Clustering	- No specific metric due to nature of algorithm	- Identifies clusters with varying densities	- Performance depends on eps and min_samples parameters
			- Captures outliers and core points effectively	
Prophet	Time Series Forecasting	- Accuracy: Strong alignment of actual vs. predicted	- Captures temporal patterns and trends	- Depends on data quality and availability
			- Provides insights into seasonal patterns	- Choice of regressors impacts accuracy
LSTM	Time Series Analysis	- MAE: 0.31	- Accurate predictions of traffic patterns	- Requires significant data preprocessing
		- Loss: Decreased over 100 epochs	- Effective in learning complex temporal dependencies	- Sensitive to hyperparameters and model architecture

4.3 Key Findings

1. **Weather and Time Impact on Accident Severity:** The Random Forest model effectively predicted accident severity, particularly for severe and less severe accidents (Classes 1 and 4). However, the overlap in feature patterns for moderate severity accidents (Classes 2 and 3) suggests the need for further feature engineering to improve model performance for these classes.

2. **Clustering of Accident Patterns:** Both K-Means and DBSCAN models successfully identified distinct clusters of accident patterns. K-Means provided a general overview of accident patterns based on location and traffic conditions, while DBSCAN highlighted densely populated regions and outliers, which could be critical for targeted traffic safety interventions.
 3. **Temporal Patterns in Accident Severity:** The Prophet model's forecasts revealed significant trends and seasonal patterns in accident severity, suggesting that certain periods of the year or specific days may be associated with higher accident severity. This information can be utilized to implement timely and targeted safety measures.
 4. **Location-Based Traffic Pattern Forecasting:** The LSTM model demonstrated the ability to accurately predict traffic patterns based on historical data, providing insights into future traffic conditions and potential areas of concern.
 5. **Model Strengths and Areas for Improvement:**
 - The models showed robust performance in identifying patterns and predicting traffic accident severity and patterns.
 - Challenges include the model's sensitivity to feature selection, initial parameter settings, and the need for more distinct feature engineering, particularly in the case of overlapping feature patterns in accident severity prediction.
-

5. Discussion and Conclusion

5.1 Discussion

The RoadSense project demonstrates the effectiveness of machine learning techniques in predicting and understanding traffic accidents. The models achieved high accuracy, with the Long Short-Term Memory Networks performing the best in terms of all evaluation metrics. This success is attributed to the model's ability to capture complex temporal dependencies and interactions between features.

The insights gained from the model evaluations reveal significant factors influencing traffic accidents, including time of day, weather conditions, and geographic location. The ability of the models to identify high-risk areas and time periods provides valuable information for traffic management and safety measures.

The project also highlights the importance of incorporating various machine learning methods to improve prediction accuracy. While deep learning models like CNNs and LSTMs showed superior performance, traditional methods such as Random Forest and Gradient Boosting Machines also provided valuable insights. The complementary nature of these models emphasizes the need for a multi-faceted approach to traffic accident prediction.

5.2 Conclusion

In conclusion, the RoadSense project successfully utilized advanced machine learning techniques to enhance traffic safety through predictive modeling. The models developed offer a scalable solution for real-time traffic management, enabling more informed decision-making and targeted interventions to reduce accidents. The project underscores the potential of combining machine learning with traffic data to improve road safety and provides a solid foundation for future research and development in this field.

5.3 Future Work

Future work for the RoadSense project will focus on several key areas to further enhance the predictive models and their applications:

1. **Integration of Real-Time Data:** Incorporating real-time traffic data, such as live traffic flow, congestion levels, and real-time weather updates, will improve the models' accuracy and relevance. Real-time data integration will enable dynamic and adaptive traffic management solutions.
2. **Expansion of Model Features:** Future models will include additional factors such as driver behavior, vehicle types, road conditions, and traffic signal timings. These factors will provide a more comprehensive understanding of accident causes and contribute to more precise predictions.
3. **Collaboration with Traffic Authorities:** Engaging with traffic authorities and urban planners will help translate the insights from the models into actionable strategies and interventions. Collaborative efforts will focus on implementing safety measures and traffic management policies based on the model's findings.
4. **Evaluation of Additional Machine Learning Techniques:** Exploring advanced machine learning techniques, such as reinforcement learning and transfer learning, may provide further improvements in prediction accuracy and model robustness.
5. **Long-Term Impact Assessment:** Conducting longitudinal studies to evaluate the long-term impact of the implemented safety measures and predictive models will be crucial for assessing their effectiveness in reducing traffic accidents and improving road safety.

By addressing these areas, the RoadSense project aims to make a meaningful contribution to traffic safety and continue advancing the field of predictive analytics in transportation.

6. References

1. Moosavi, A., Hashemi, M., Yu, D., Nouri, M., Fathy, M., & Claudel, C. (2019). [A countrywide traffic accident dataset](https://doi.org/10.1016/j.dib.2019.104339). *Data in Brief*, 27, 104339. <https://doi.org/10.1016/j.dib.2019.104339>
2. Sun, T., Qin, C., Gao, Y., Sun, W., & Shen, J. (2021). [Accident risk prediction based on heterogeneous sparse data](https://doi.org/10.1109/TITS.2021.3054362). *IEEE Transactions on Intelligent Transportation Systems*, 23(7), 11768-11778. <https://doi.org/10.1109/TITS.2021.3054362>
3. Haghighat, A., Azadeh, A., & Salari, F. (2020). [Exploring the relationship between alcohol and driver characteristics in motor vehicle accidents](https://doi.org/10.1186/s12889-020-08643-6). *BMC Public Health*, 20, 654. <https://doi.org/10.1186/s12889-020-08643-6>
4. Li, Y., Sun, T., Qin, C., & Chen, Y. (2020). [Highway crash detection and risk estimation using deep learning](https://doi.org/10.1016/j.trc.2020.02.020). *Transportation Research Part C: Emerging Technologies*, 112, 39-56. <https://doi.org/10.1016/j.trc.2020.02.020>
5. Chen, J., Zhang, Y., Liu, Y., & Liu, X. (2019). [Traffic accident analysis using machine learning paradigms](https://doi.org/10.1016/j.aap.2019.105282). *Accident Analysis & Prevention*, 135, 105282. <https://doi.org/10.1016/j.aap.2019.105282>
6. Ahmad, M., Khalid, M. S., Hussain, T., & Habib, U. (2023). Improving traffic accident severity prediction using MobileNet transfer learning model and SHAP XAI technique. *Applied Sciences*, 13(6), 3334. <https://doi.org/10.3390/app13063334>
7. Singh, H., Gupta, R., & Kumar, M. (2022). [Evaluating the impact of weather conditions on traffic accidents in urban areas](https://doi.org/10.1109/TITS.2021.3132703). *IEEE Transactions on Intelligent Transportation Systems*, 23(3), 1809-1819. <https://doi.org/10.1109/TITS.2021.3132703>
8. GitHub Repository : <https://github.com/Pranavsharma13/RoadSense--Advanced-Predictive-Modelling-for-Traffic-Safety>
9. Road Sense App Link: <https://roadsense-trafficaftey.streamlit.app/>
10. Kaggle Dataset Source: [US Accidents \(2016 - 2023\) \(kaggle.com\)](https://www.kaggle.com/datasets/usaccidents/us-accidents-2016-2023)