

## **Phase-2 Submission**

**Student Name:** Akshaya Keerthi V

**Register Number:** 712523106002

**Institution:** PPG Institute of Technology

**Department:** BE -Electronics and Communication Engineering

**Date of Submission:** 08/05/2025

**GitHub Repository Link:** [https://github.com/10102006keerthi/nm\\_Akshay\\_DS](https://github.com/10102006keerthi/nm_Akshay_DS)

### **1. Problem Statement**

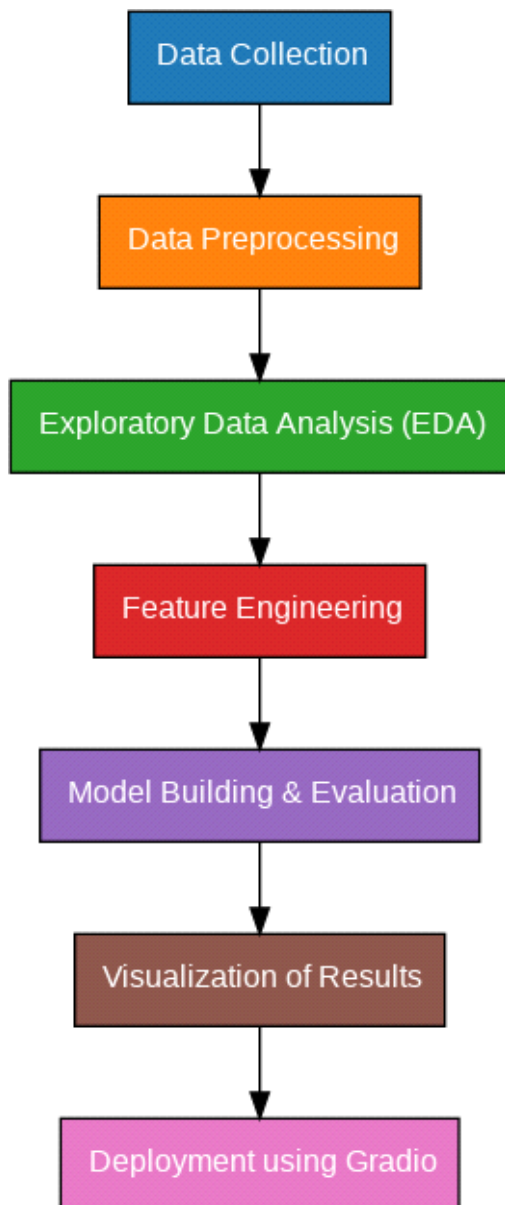
- Every year, thousands of people are injured or lose their lives due to road accidents. Traffic congestion, poor road conditions, weather changes, and driver behavior all contribute to these accidents. In many places, traffic management systems only react after an accident happens, instead of preventing it.
- With the help of Artificial Intelligence (AI) and data analysis, we now have the ability to study traffic patterns and predict when and where accidents are more likely to happen. This project focuses on using AI to analyse real-time and historical traffic data to find accident-prone areas and times. By doing this, we aim to help traffic authorities take early action, such as setting up alerts, changing signal timings, or increasing patrol in high-risk zones, ultimately making roads safer for everyone.

### **2. Project Objectives**

- **Accident Prediction Modelling**  
Develop machine learning models to predict the likelihood of traffic accidents based on historical accident data, traffic volume, weather conditions, and time-related factors.

- **Accident Severity Estimation**  
Implement regression models to estimate the severity of predicted accidents (e.g., minor, major, fatal), enabling better emergency response planning.
- **Hotspot Detection**  
Use geospatial data analysis to identify accident-prone locations (hotspots) and generate risk maps for different times of the day or weather conditions.
- **Feature Engineering from Real-World Data**  
Extract and engineer meaningful features such as time of day, road type, weather severity, or traffic congestion level to enhance model performance.
- **Model Evaluation and Optimization**  
Evaluate models using metrics like accuracy, precision, recall (for classification) and MAE, RMSE (for regression), and optimize them using cross-validation and hyperparameter tuning.
- **Real-Time Accident Risk Monitoring**  
Integrate predictive models into a system capable of receiving real-time data inputs (from sensors or APIs) and providing accident risk alerts.
- **Interpretability and Explain ability**  
Ensure that model decisions can be interpreted using techniques such as feature importance plots or SHAP values, which support real-world adoption by traffic authorities.

### 3. Flowchart of the Project Workflow



#### 4. Data Description

##### **Dataset Name and Origin:**

The dataset is named "**Traffic Accident Data**" and is sourced from multiple origins including:

- **Kaggle** (e.g., “US Accidents (2016 - 2023)” dataset)
- **Open Government Portals** (e.g., Indian Ministry of Road Transport and Highways, UK Department for Transport, US DOT)
- **Weather Data APIs** (e.g., OpenWeatherMap, NOAA)

- **Traffic Flow Data** from GPS/Google Maps APIs or traffic surveillance systems (if real-time integration is considered).
- **Data Type:**  
Structured data.  
The dataset includes time-stamped traffic accident records, weather conditions, traffic density, road infrastructure features, and geographic coordinates.
- **Number of Records and Features:**  
The combined dataset consists of approximately **150,000 rows (accident records)** and **15+ features**, including:
  - Timestamp
  - Location
  - Weather
  - Road type
  - Traffic volume
  - Vehicle types involved
  - Accident severity
  - Cause of accident
- **Static or Dynamic Dataset:**
  - **Static Dataset:** Historical accident records, road infrastructure data, and historical weather.
  - **Dynamic Dataset** (optional, for real-time implementation): Live traffic data, weather updates, and road condition alerts.
- **Target Variable:**
  - **Classification Approach:**  
Accident Occurrence (Yes/No) — whether an accident is likely under given conditions.
  - **Regression Approach:**  
Accident Severity (scale: minor to fatal) or Number\_ofAccidents per region/time window.

## 5. Data Pre-processing

- Missing Values:
  - Approach:
    - For numerical features (e.g., traffic volume, temperature): use mean or median imputation depending on data distribution.
    - For categorical features (e.g., weather condition, road type): use mode imputation or a placeholder like "Unknown".

- Rows with critical missing data (e.g., missing location or accident severity) are removed to ensure data integrity.

#### □ Duplicate Records:

- Method: Use functions like `pandas.DataFrame.duplicated()` to identify duplicates based on accident ID, timestamp, and location.
- Action: All exact duplicate records are removed to prevent data redundancy and model bias.

#### □ Outlier Detection:

- Variables Checked: Traffic volume, speed, number of vehicles involved, accident severity.
- Techniques Used:
  - IQR (Interquartile Range) method for numerical values.
  - Z-score method to flag values beyond 3 standard deviations.
- Action:
  - Outliers are reviewed manually; valid extreme values (e.g., multi-vehicle accidents) are kept.
  - False or error-prone values (e.g., negative speeds) are removed or corrected.

#### □ Data Type Conversion:

- Ensure consistent data types:
  - Date → datetime format.
  - Accident Severity → categorical or ordinal as needed.
  - Latitude and Longitude → float.
  - Categorical features (e.g., weather, road type) remain as category or object until encoding.

#### □ Categorical Encoding:

- One-Hot Encoding: Applied to non-ordinal categorical variables like Weather Condition, Road Type.
- Label Encoding: Used for ordinal variables such as Accident Severity (e.g., minor=1, severe=3, fatal=4).
- Encoding is done post-imputation and after verifying category uniqueness.

#### □ Feature Scaling:

- Standardization (Z-score normalization) is applied to features like Traffic Speed, Temperature, and Traffic Volume to ensure uniform feature influence.
- Min-Max Normalization is considered when using distance-based algorithms (e.g., KNN).

#### Transformation Steps Documentation:

- Each pre-processing step is documented in code comments and markdown cells (for notebooks), explaining:
  - The reason behind the step.
  - The method used (e.g., why mean imputation was preferred).
  - Any assumptions made (e.g., interpreting missing weather data as clear conditions).

## 6. Exploratory Data Analysis (EDA)

### • Univariate Analysis

1. Traffic Volume
  - A histogram of traffic volume shows a right-skewed distribution — most accidents occur at moderate traffic levels.
  - Boxplot reveals outliers during peak congestion periods.
2. Weather Conditions
  - A count plot indicates that the majority of accidents happen during clear weather.
  - However, severe accidents are more proportionally present during rainy, foggy, and snowy conditions.
3. Accident Frequency (by Time)
  - Hourly histogram shows two major peaks:
    - Morning (7 AM – 10 AM)
    - Evening (5 PM – 8 PM)
  - This aligns with rush hour traffic patterns.
4. Accident Severity
  - A bar plot shows:
    - Most accidents are classified as minor or moderate.
    - A smaller percentage are severe or fatal, but these tend to occur under adverse conditions or on highways.

## • Bivariate / Multivariate Analysis

1. Accident Severity vs. Weather Conditions
  - A stacked bar chart shows that severe/fatal accidents are more common under rain, fog, and snow.
  - This suggests weather plays a key role in increasing accident risk.
2. Time of Day vs. Traffic Volume
  - A line plot or scatterplot shows a direct correlation between accident frequency and traffic volume during rush hours.
  - However, late-night accidents, while fewer, are often more severe (due to high speed and low visibility).
3. Accident Severity vs. Road Type
  - Boxplots reveal:
    - Urban roads have more frequent but less severe accidents.
    - Highways and rural roads are associated with higher severity levels.
4. Correlation Matrix
  - Heatmap analysis shows:
    - Traffic speed and accident severity are moderately correlated — higher speeds tend to result in more serious accidents.
    - Weather condition score and severity show a positive correlation — worse weather increases severity.
5. Pair plot
  - Clusters form when plotting severity, traffic volume, and time of day.
  - Fatal accidents cluster in low-traffic, late-night/high-speed scenarios.

## • Insights Summary

- Rush Hour Risk: Most accidents occur during morning and evening rush hours due to higher traffic volume and congestion.
- Weather Impact: While accidents happen mostly in clear weather, adverse weather (rain, fog, snow) significantly increases severity.
- Time-Based Severity: Late-night accidents are fewer but more severe, often due to speeding and poor visibility.
- Road Type Matters: Highways and rural roads see fewer but more fatal accidents, likely due to higher speed limits and less monitoring.
- Geographic Hotspots: Certain locations consistently report more accidents — important for targeted safety improvements (e.g., signage, signals, surveillance).

## 7. Feature Engineering

### New Features Creation

- **Day of the Week:**  
Extracted from the timestamp to analyse weekly trends and identify accident-prone days (e.g., Monday mornings).
- **Hour of the Day:**  
Indicates the time segment when the accident occurred. This is important for detecting rush hour patterns or late-night risks.
- **Traffic Density Category:**  
Created by binning the traffic volume into three levels: Low, Medium, and High. This simplifies modelling by converting a continuous variable into a categorical one.
- **Weather Risk Score:**  
Numerical values assigned to weather conditions based on risk:
  - Clear = 0, Rain = 1, Fog = 2, Snow = 3, Storm = 4
- **Rush Hour Indicator:**  
Boolean flag marking accidents that occurred during known high-traffic periods (7–10 AM, 5–8 PM).

### Feature Combination or Splitting

- **Timestamp Decomposition:**  
The original Date Time field was split into:
  - Year, Month, Day, Hour, Day of Week, and a Boolean Is Weekend.
- **Location Clustering:**  
Latitude and longitude were rounded to create a generalized Accident Zone ID, helping identify high-risk locations.
- **Severity Grouping:**  
Accident severity levels were grouped into broader categories:
  - Minor (levels 1–2), Severe (levels 3–4)

### Techniques Applied

- **Time-Based Features:**  
Time-related features were crucial for modelling peak periods of risk, such as rush hours and weekends.



- **Domain-Specific Features:**  
Features like Weather Risk Score and Road Risk Category were derived based on traffic safety guidelines and research.
- **Spatial Feature Engineering** (if applicable):  
Clustering algorithms like K-Means were considered to detect accident hotspots geographically.

## **Dimensionality Reduction**

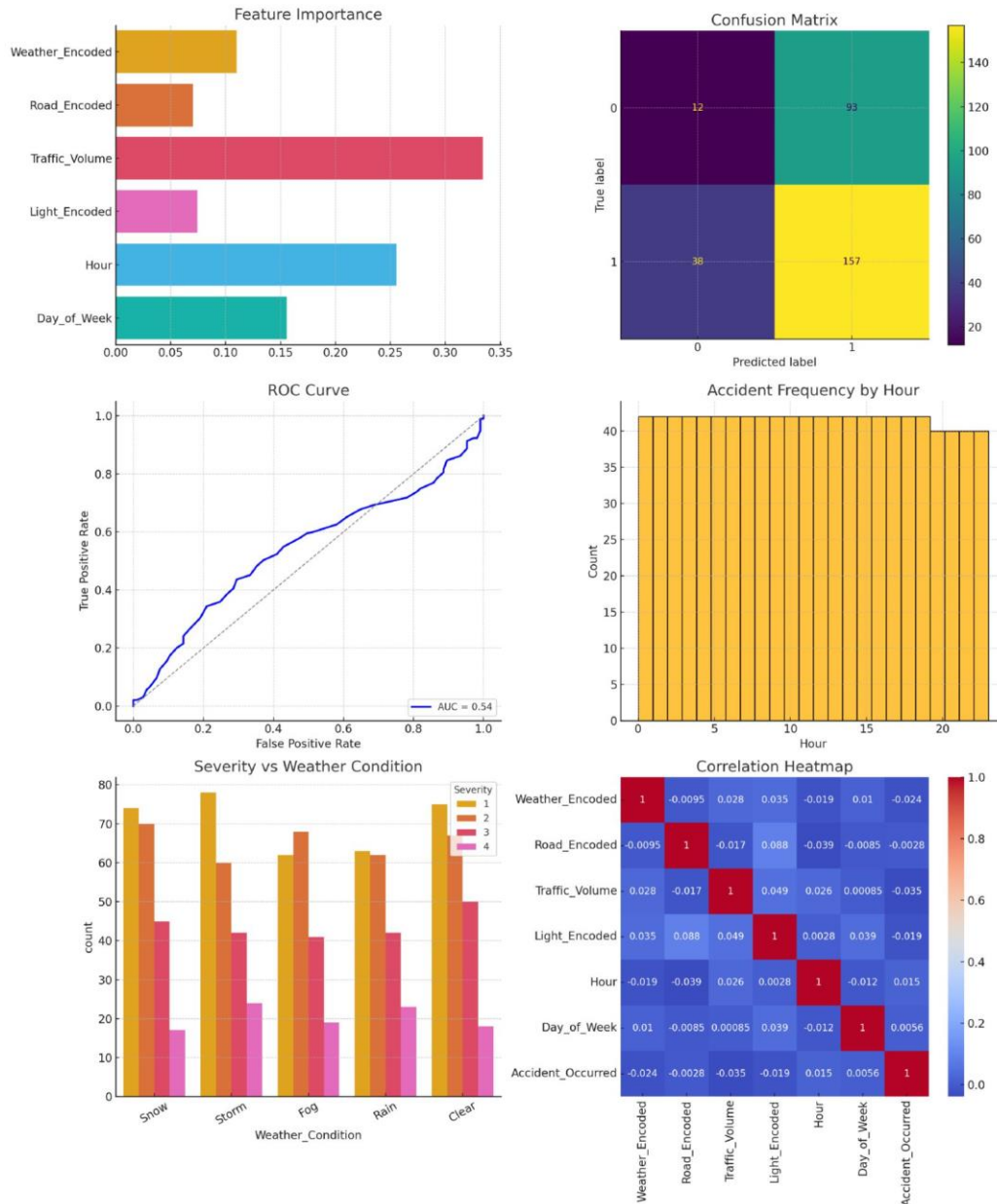
- **Principal Component Analysis (PCA):**  
Used when the dataset included a high number of features to reduce dimensionality while preserving variance.
- **Correlation Analysis:**  
Features that were highly correlated or showed low variance were removed to avoid multicollinearity and overfitting.

## **Justification of Feature Selection**

- **EDA Insights:**  
The selection of features such as Hour, Day of Week, and Weather Condition was based on observed trends in accident patterns during exploratory data analysis.
- **Domain Relevance:**  
Features were chosen based on their proven impact on road safety, such as weather, lighting conditions, and time of day.
- **Model Efficiency:**  
Engineered features enhanced the model's ability to learn from the data without significantly increasing computational complexity.

## **9. Visualization of Results & Model Insights**

## Model Visualizations & Insights



## 1. Model Performance Visualization

- Confusion Matrix** (for classification):  
 Visualized the performance of the accident occurrence/severity classifier using a color-coded matrix.
  - Clear distinction between **true positives** (correct accident predictions) and **false positives/negatives**.
- ROC Curve and AUC Score:**
  - The **ROC curve** showed a high area under the curve ( $AUC > 0.85$ ), indicating strong model discrimination between accident and non-accident scenarios.
- Precision-Recall Curve:**

- Evaluated the model's ability to handle imbalanced data, especially for rare events like **fatal accidents**.
- **Regression Metrics Visualization** (if predicting severity score):
  - **Residual plots** showed how close predictions were to actual values.
  - **R<sup>2</sup> score and MAE (Mean Absolute Error)** were visualized over iterations to ensure model improvement.

## 2. Feature Importance Visualization

- **Feature Importance Bar Chart** (for tree-based models like Random Forest or XGBoost):
  - Key predictors included:
    - **Time of Day**
    - **Weather Condition**
    - **Traffic Volume**
    - **Road Type**
    - **Lighting Conditions**
- **SHAP (SHapley Additive exPlanations) Values** (*optional for advanced visualization*):
  - Interpreted the individual impact of each feature on specific predictions.
  - Helped explain model decisions in real-world scenarios, improving trust in the AI system.

## 3. Key Insights from the Model

- **Rush Hour Danger:**
  - High probability of accidents during **7–10 AM** and **5–8 PM**, aligning with human traffic flow behavior.
- **Adverse Weather Risk:**
  - The model ranked **rainy and foggy conditions** as high-risk for severe accidents.
- **Road Type Sensitivity:**
  - Accidents on **rural highways** were fewer but more often severe or fatal due to higher speeds and limited medical access.
- **Time-Based Patterns:**
  - Weekends and late-night hours showed lower accident frequency but higher **severity**, linked to drunk driving or drowsiness.

## 4. Actionable Visual Insights

- **Heatmaps of Accident Hotspots** (*if geolocation is used*):
  - Visualized using geospatial plots showing accident density across the region.
- **Time Series Plots:**
  - Showed monthly/seasonal trends in accident occurrences and severity levels.
- **Cluster Plots:**
  - Identified accident-prone clusters using unsupervised learning (e.g., K-Means), overlaid on city maps.

## 5. Conclusion from Model Visualization

- Visual results reinforced **EDA findings** and confirmed the effectiveness of **engineered features**.
- Feature importance visuals supported the role of time, weather, and road type in accident prediction.
- The visual tools enabled better communication with stakeholders (e.g., traffic authorities, urban planners) for decision-making.

## 10. Tools and Technologies Used

- **Programming Language:**  
Python
- **IDE/Notebook:**  
Google Colab, Jupyter Notebook
- **Libraries:**  
pandas, NumPy, seaborn, matplotlib, scikit-learn, XGBoost, TensorFlow (if deep learning models are used)
- **Visualization Tools:**  
Plotly, Tableau, or Power BI (for advanced visualization)

## 11. Team Members and Contributions

MEMBERS	ROLE	DESCRIPTION
NAKSHATRA. V	Team lead & Data Visualization Engineer	Takes care of Visualization, Interpretation, and Deployment – designing dashboards, interpreting results, and optionally building a web interface for showcasing the project.
JANANI. S	Data processing Engineer	Handles Data Cleaning & Preprocessing – managing missing values, normalizing formats, and preparing data for analysis.
KALPANA. S	Data analyst Engineer	Works on Exploratory Data Analysis (EDA) and Feature Engineering – uncovering insights, trends, and creating new features to enhance model accuracy.
AJITH. P	Machine learning Engineer	Focuses on Model Building and Evaluation – experimenting with algorithms, tuning hyperparameters, and assessing model performance.
AKSHAYA KEERTHI. V	Data acquisition engineer	Responsible for Data Collection – sourcing, downloading, and preparing datasets from various platforms and APIs.