

# US RoadSafe Analytics

---

## Project Introduction

---

Road accidents represent a significant public safety challenge in the United States, causing substantial loss of life, injuries, and economic burden annually. Understanding the underlying patterns and contributing factors to these accidents is crucial for developing effective prevention strategies and improving transportation infrastructure. The **RoadSafe Analytics** project addresses this critical need through comprehensive data-driven analysis.

This project leverages advanced exploratory data analysis (EDA) techniques to examine a large-scale dataset of over 4 million road accident records from across the United States. By systematically analyzing accident patterns across temporal, spatial, environmental, and severity dimensions, the project aims to extract actionable insights that can inform policy decisions, guide infrastructure investments, and support public safety initiatives.

Using Python-based data science tools including Pandas, Matplotlib, and Seaborn, this analysis transforms raw accident data into meaningful intelligence about when, where, and under what conditions accidents occur, and what factors contribute to their severity. The findings serve as an evidence base for stakeholders including transportation authorities, urban planners, policy makers, and safety advocates.

---

## Project Statement

---

The primary goal of this project is to analyze a large dataset of road accidents to uncover trends, patterns, and key factors contributing to accident severity. The project involves performing in-depth exploratory data analysis (EDA) using Python libraries such as Pandas for data manipulation, Matplotlib and Seaborn for statistical visualization, and Streamlit for developing an interactive dashboard to extract meaningful insights that can help improve road safety across the United States.

---

## Project Objectives

---

The RoadSafe Analytics project is structured around the following core objectives:

### 1. Data Acquisition and Understanding

- Obtain and load the US Accidents dataset containing 4+ million records
- Explore the dataset structure, schema, and dimensions
- Understand the relationships between variables and data characteristics

### 2. Data Quality and Preprocessing

- Handle and preprocess real-world accident data effectively
- Address missing values, outliers, and data inconsistencies
- Transform and engineer features for enhanced analytical capability

### 3. Temporal Pattern Analysis

- Explore accident frequency based on time-related factors
- Identify patterns across hours of the day, days of the week, and months
- Detect seasonal and temporal trends in accident occurrence

### 4. Spatial and Geographic Analysis

- Analyze accident distribution across locations (states and cities)
- Identify accident hotspots and high-risk geographic areas
- Understand regional variations in accident patterns

### 5. Environmental Factor Investigation

- Examine the impact of weather conditions on accident occurrence and severity
- Assess how road surface conditions affect accident outcomes
- Evaluate the relationship between visibility and accident characteristics

### 6. Severity Assessment

- Visualize patterns in accident severity using various plots and graphs
- Identify key determinants of accident severity
- Understand correlations between environmental factors and severity levels

### 7. Insight Generation and Hypothesis Testing

- Answer critical questions about accident patterns
- Test hypotheses regarding contributing factors
- Derive evidence-based insights to support road safety awareness and policy recommendations

## 8. Documentation and Knowledge Transfer

- Prepare detailed documentation of methodology, findings, and results
- Create comprehensive visualizations that communicate insights effectively
- Develop presentation materials for stakeholder communication

---

## Dataset Overview

- **Dataset Name** : US Accidents (2016 - 2023)
- **Geographic Coverage** : 49 states of the USA
- **Time Period** : February 2016 to March 2023
- **Total Records** : Approximately **7.7 million** accident records
- **Source** : Kaggle

## Description

This is a **countrywide car accident dataset** that covers **49 states of the USA**. The accident data were collected from **February 2016 to March 2023**, using multiple APIs that provide streaming traffic incident (or event) data.

## Key Features

- **Never** expected update frequency (static historical dataset)
- Real-time streaming data collection from multiple sources
- Comprehensive coverage across the United States
- Multi-year temporal span (7+ years)
- Rich feature set including temporal, spatial, environmental, and severity attributes

## Dataset Utility

This dataset is valuable for:

- Traffic accident pattern analysis
- Road safety research
- Predictive modeling for accident risk
- Policy-making and infrastructure planning
- Geospatial analysis and hotspot identification
- Time-series analysis of accident trends
- Machine learning applications in transportation safety

## Stepwise Pipeline Structure

### Initial Data

- Loaded from CSV file (example: `US_Accidents_March23.csv`)
- Original data shape: millions of rows, approximately 47 columns
- Notable missing data mainly in weather features and geographic details

### Duplicate Records Removal

- Duplicates eliminated based on unique `"ID"` field
- Several thousand duplicate entries dropped to ensure uniqueness

### Columns with High Missing Data Dropped

- Only `End_Lat` and `End_Lng` columns were dropped due to >44% missing values
- Other columns with missing data were imputed rather than dropped

### Non-Analytical Text/ID Columns Removed

- Columns such as `"ID"`, `"Source"`, `"Description"`, `"Street"`, `"Country"`, `"Zipcode"`, `"Timezone"`, `"Airport_Code"`, and `"Amenity"` were removed to reduce noise and focus on analytical features

### Temporal Columns Cleaning

- Converted `"Start_Time"` and `"End_Time"` to datetime objects with error coercion
- Rows with invalid timestamps were removed for data quality

### Geographic Coordinates Cleaning

- Latitude and Longitude columns ( `"Start_Lat"` , `"Start_Lng"` ) converted to numeric types

- Renamed to "Latitude" and "Longitude"
- Rows with missing or invalid coordinates were dropped

### Severity Filtering

- Kept only records where "Severity" is 1, 2, 3, or 4
- Outliers or invalid severity levels removed

### Minor Missing Data Rows Removed

- Columns with low missingness (0–3%) were identified
- Rows missing data in those columns were dropped for increased data purity

### Weather Feature Imputation

- **Wind\_Speed(mph)** : Median imputation for missing values to maintain robustness
- **Precipitation(in)** : Missing values filled with 0.0, assuming no precipitation
- **Wind\_Chill(F)** : Missing values predicted by linear regression model trained on weather-related features (Wind\_Speed(mph) , Temperature(F) , Humidity(%))

### Numeric Imputation

- Remaining numeric columns with missing values filled using median imputation

### Temporal Feature Engineering

- Created the following new features:
  - Duration\_Minutes : Difference between "End\_Time" and "Start\_Time" in minutes
  - Year : Extracted from "Start\_Time"
  - Hour : Hour of day from "Start\_Time"
  - DayOfWeek : Weekday extracted from "Start\_Time"
  - Month : Month extracted from "Start\_Time"
  - IsWeekend : Binary indicator if day is Saturday or Sunday

### Categorical Encoding

- Boolean columns ( "Roundabout" , "Station" , "Stop" , "Traffic\_Calming" , "Traffic\_Signal" , "Turning\_Loop" ) converted to integer (0/1)
- "Sunrise\_Sunset" column converted to binary "IsDay" feature (Day=1, Night=0)

### Redundant Temporal/Weather Columns Dropped

- Columns removed after feature engineering to avoid redundancy:
  - "Start\_Time" , "End\_Time"
  - "Weather\_Timestamp"
  - Twilight indicators: "Civil\_Twilight" , "Nautical\_Twilight" , "Astronomical\_Twilight"
  - Original "Sunrise\_Sunset"

### Final Cleanup

- Any remaining rows with NaN values across any feature removed to ensure a complete dataset

### Data Export

- Final, cleaned dataframe saved as CSV at user-specified path
- Typical final dataset shape: ~ 6.9 millions of rows, 35 columns with no missing values

### Imputation Summary

Column	Imputation Technique	Notes
Wind_Speed(mph)	Median imputation	Robust central tendency for skew-affected data
Precipitation(in)	Zero-fill	Assumes missing = no precipitation
Wind_Chill(F)	Linear regression imputation	Predicted from Wind_Speed(mph) , Temperature(F) , and Humidity(%)

# Streamlit Dashboard Features Explanation

## 1. Interactive User Interface

### Page Header

- **Title:** "🚀 Data Preprocessing Pipeline"
- **Subtitle:** "Automated Data Cleaning & Feature Engineering"
- Professional emoji-enhanced UI for visual appeal and clarity

### Configuration Inputs

- 📄 **Input Data Path:** Text field to specify raw dataset location
  - Default: data/US\_Accidents\_March23.csv
  - Allows users to load any custom CSV file
- 📄 **Output Data Path:** Text field to specify where cleaned data will be saved
  - Default: data/US\_Accidents\_preprocessed.csv
  - Enables flexible output destination configuration

## 2. Pipeline Overview Table

Interactive markdown table displaying all 15 preprocessing steps:

- Step number
- Operation name
- Purpose of each step

**Benefit:** Users can understand the complete pipeline flow before execution

## 3. Start Preprocessing Button

- **Primary action button** with rocket emoji (🚀)
- Full-width responsive design
- Triggers the complete 15-step preprocessing pipeline
- Provides clear visual affordance for pipeline execution

## 4. Real-Time Progress Tracking

### Progress Bar

- Visual indicator showing completion percentage
- Updates after each step (1/15 → 15/15)
- Reaches 100% on completion

### Status Text Display

- Shows current step number and total steps
- Provides brief description of operation in progress
- **Format:** "Step X/15: [Operation Description]"

### Animation Delays

- 0.3-second pause between steps for readable visual feedback
- Prevents overwhelming rapid updates

## 5. Live Metrics Dashboard

Real-time metric cards displaying key dataset statistics:

1. **Rows Count**
  - Updates after each operation
  - Shows impact on data volume
  - Formatted with thousand separators (e.g., "2,450,123")
2. **Columns Count**

- Tracks column changes through pipeline
  - Shows feature optimization progress
3. **Missing Values Count**
- Monitors data completeness
  - Decreases as imputation progresses
  - Target: 0 by final step
4. **Progress Status**
- Displays "X/15 steps" completed
  - Visual checkpoint indicator

**Layout:** 4-column responsive grid for side-by-side viewing

---

## 6. Step-by-Step Logging

### Scrollable Log Container

- Accumulates all processing events
- Height-limited with vertical scroll
- **Features:**
  - Each step logged with checkmark (✓)
  - Step number and description
  - Data shape after each operation
  - HTML-styled container with monospace font for clarity

**Benefit:** Complete audit trail of all transformations for debugging and verification

---

## 7. Error Handling & Recovery

### FileNotFoundException Handling

- Displays error message with file path
- Provides helpful advisory to check file path
- Graceful failure without crashing

### General Exception Handling

- Catches all preprocessing errors
- Shows error message to user
- Expandable error details section
- Full stack trace available for debugging

---

## 8. Final Summary Statistics

Comprehensive 3-column dashboard displaying before/after metrics:

### Column 1: Row Analysis

- **Initial Rows:** Original data volume
- **Final Rows:** Cleaned data volume
- **Rows Removed:** Count with percentage change indicator

### Column 2: Column Analysis

- **Initial Columns:** Original feature count (~47)
- **Final Columns:** Optimized feature count (~35-40)
- **Columns Changed:** Net change with positive/negative indicator

### Column 3: Data Quality

- **Missing Values:** Always "0" on completion
- **Data Quality:** "✓ Validated"
- **File Saved:** "✓ Success"

---

## 9. Sample Data Preview

### Dataframe Display

- Shows first 10 rows of cleaned dataset
  - Full-width responsive table
  - Interactive sorting and column inspection
  - Allows stakeholders to verify data quality
- 

## 10. Final Column List

---

### Expandable Section

- **Total Columns Count:** Numeric summary
  - **Complete Column Names:** Comma-separated list in code block
  - **Display:** Monospace font for clarity
  - **Purpose:** Reference all final features in cleaned dataset
- 

## 11. Download Button

---

### One-Click CSV Export

- **Label:** "📄 Download Preprocessed Data (CSV)"
- **Format:** CSV (Comma-Separated Values)
- **File Name:** Derived from output path (user-configurable)
- **MIME Type:** Properly set for browser download
- **Styling:** Primary button type, full-width
- **Purpose:** Direct download without file system navigation

**Benefit:** Instant access to cleaned data for downstream analysis or modeling

---

## 12. Dynamic State Management

---

### Before Execution

- Displays pipeline overview table
- Shows informational message
- Waits for user action

### During Execution

- Updates all metrics in real-time
- Fills progress bar incrementally
- Accumulates log entries
- Blocks new executions

### After Completion

- Shows success message (✅)
  - Displays final summary statistics
  - Shows sample data preview
  - Enables CSV download
  - Lists all final columns
- 

## 13. User Experience Features

---

### Responsive Design

- Multi-column layouts adapt to screen size
- Mobile-friendly interface
- Touch-accessible buttons

### Visual Hierarchy

- Clear section separation with markdown separators (---)
- Emoji icons for quick scanning
- Color-coded success/error messages
- Expandable sections for detailed info

Accessibility

- Descriptive button labels
  - Help text in input fields
  - Error messages are informative
  - Sequential step-by-step process
- 

14. Data Quality Assurance

Validation Checkpoints

- File existence verification
- Data type checking
- Row count tracking
- Missing value monitoring
- Column consistency validation

Quality Metrics

- Initial vs. final row comparison
  - Column count optimization
  - Missing value elimination (0% target)
  - Data completeness percentage
- 

15. Performance Metrics Capture

The dashboard captures and displays:

- Rows removed at each step
  - Columns dropped at each step
  - Values imputed per feature
  - Processing time indicators
  - Data shape transitions
  - Quality score improvements
- 

Key Advantages of This Streamlit Dashboard

👤 **User-Friendly:** No coding required for data preprocessing 🔄 **Transparent:** Complete visibility into all transformations ⚡ **Fast Feedback:** Real-time progress updates 🛡️ **Error-Resilient:** Robust error handling 📄 **Reproducible:** Consistent, documented pipeline 📁 **Exportable:** Direct CSV download capability 📈 **Scalable:** Handles millions of rows efficiently 📖 **Educational:** Clear logging for learning purposes 🏢 **Professional:** Publication-ready data quality 🤝 **Production-Ready:** Suitable for team collaboration

---

Univariate Analysis Module

1. Column Selector

- **Selectbox Dropdown:**
    - Allows users to pick any column from the dataset for univariate analysis.
    - Prevents accidental action by setting default to “–Choose a column–”.
    - If no column is selected, an informative message prompts user action.
  - **Benefit:** Ensures analysis is intentional and input-driven.
- 

2. Data Type Detection and Branching

- Automatically detects whether the selected column is numeric or categorical.
  - Sets the analytic workflow (visualization type) accordingly, ensuring context-appropriate visual feedback.
- 

3. Numerical Feature Analysis

- **Histogram Visualization (Plotly)**

- Plots the distribution of the selected numerical feature.
    - Uses 30 bins for clarity and interpretability.
  - KDE Option**
    - Checkbox toggle: "Include KDE plot in histogram"
    - If selected, overlays a Kernel Density Estimate (KDE) line using `scipy.stats.gaussian_kde`.
    - KDE is scaled to match the histogram's bin counts for direct comparison.
    - KDE curve is rendered in red for immediate distinction.
  - Benefit:** Enables rich exploration of numeric data with both frequency and probability density views.
- 

## 4. Categorical Feature Analysis

---

- For non-numeric columns, displays the top 10 most frequent categories.
  - Bar Chart Visualization (Plotly)**
    - Plots categories (sorted by frequency) on the X axis and counts on the Y axis.
    - Simple, intuitive format for spotting dominant category values.
  - Benefit:** Helps users quickly identify common classes, outliers, and skewness in categorical data.
- 

## 5. Modern Visualization

---

- All charts are rendered with Plotly Express, providing:
    - Interactivity: tooltips, zoom, pan, and export options included by default
    - Consistency: container-wide layout ensures full-width readability
    - Responsiveness: adapts well to browser resizing
- 

## Comparative Analysis Module

---

### 1. Feature Type Extraction and Classification

---

- Automatic separation of features into:
    - Numerical Features:** All float and integer columns
    - Categorical Features:** All object, boolean, and category columns
- 

### 2. Chart Type Selection

---

- Selectbox Dropdown:**
    - Lets the user choose among "Scatterplot", "Box Plot", and "Heatmap"
    - Updates the interface dynamically depending on the selected plot type
- 

### 3. Scatterplot Mode

---

- Purpose:** Visualizes numeric-to-numeric variable associations, with class-based coloring
  - Inputs:** User selects two numeric columns (X and Y axes) from dropdowns
  - Validation:** If either axis is unselected, informs user to choose both
  - Plot:**
    - Rendered using Plotly Express
    - Points colored by "Severity" (accident severity)
    - Output: Interactive, publication-ready scatterplot
  - Benefit:** Instantly reveals linear, non-linear, and outlier relationships and class separations
- 

### 4. Box Plot Mode

---

- Purpose:** Explores numeric feature distributions grouped by accident severity classes
  - Inputs:** Single numeric column selection for box plot (Y axis)
  - Validation:** If not selected, info message prompts for input
  - Plot:**
    - Drawn using Plotly Express
    - Box plot of the selected numeric variable, colored/grouped by "Severity"
    - Output: Interactive distribution summary per class
  - Benefit:** Quickly detects medians, variability, and class-wise outliers in numeric features
-



## 5. Heatmap Mode

---

### a. Numerical Correlation Heatmap

- **Purpose:** Displays pairwise Pearson correlations for all numeric columns (including "Severity")
- **Feature Set:** All numerical columns except 'Severity', which is appended for correlation reference
- **Validation:** Checks for at least two numeric features; otherwise informs the user
- **Plot:**
  - Annotated correlation matrix drawn with Plotly Figure Factory
  - Uses color intensity (Viridis scale) and numeric annotations for clarity
  - Title: "Numerical Correlation Heatmap with Severity"
- **Benefit:** Reveals multicollinearity, direct and inverse relationships

### b. Categorical Correlation Heatmap (Cramér's V)

- **Purpose:** Measures and visualizes the association between all categorical variables and "Severity"
- **Metric:** Cramér's V, calculated via chi-squared contingency analysis
- **Feature Set:** All categorical columns, forcibly including "Severity" even if stored as numeric
- **Validation:** At least two categorical features required
- **Plot:**
  - Annotated Cramér's V matrix, heatmapped and labeled via Plotly Figure Factory
  - Color-coded strength of association (0: none, 1: perfect)
  - Title: "Categorical Correlation Heatmap with Severity (Cramér's V)"
- **Benefit:** Identifies strong categorical-categorical associations, aiding feature engineering and EDA

---

## Geospatial Accident Analysis Module

---

### 1. Geographic Level Selection

---

- **Radio buttons:** lets users choose among three geography granularities:
  - Country-wide view
  - State-level filtering
  - City-level filtering
- UI adapts based on selection, enabling specific city or state pickers for zoomed-in views

---

### 2. Visualization Type Selection

---

- Radio buttons to choose plot type:
  - **Point Map:** Individual accident locations shown on map
  - **Hotspot Density:** Clusters of accident hotspots detected and sized by occurrence count

---

### 3. Accident Severity Filtering

---

- Dropdown box contains accident severity options (dynamically derived from data)
- User must select a severity level for filtering data shown on maps
- Info prompt if no severity selected to prevent empty display

---

### 4. Dynamic Geographic Filtering

---

- For **Country:**
  - No further filtering; all data visible for selected severity
- For **State:**
  - Dropdown to pick a US state by full name (using state abbreviation mapping)
  - Zoom and map center reposition to state's accident mean coordinates
  - Data filtered by state and severity
- For **City:**
  - Dropdown to pick city from selected data subset
  - Zoom and center repositioned to city coordinate mean
  - Data filtered by city and severity

---

## 5. Map Visualizations

---

## Point Map

- Scatter mapbox showing accident points colored by severity
- Color mapping: Green (1) → Red (4)
- Hover tool shows region details (state, city, or country label)
- Map styled with `carto-positron` for clean backgrounds
- Zoom and center dynamically adjusted by geography selection

## Hotspot Density Map

- Uses DBSCAN clustering algorithm on geospatial coordinates (Haversine distance)
- Epsilon set for ~1km radius clusters, min samples=5
- Clusters labeled; noise (-1) points ignored
- Creates summarized cluster dataframe aggregating accident counts and mean lat/lon
- Cluster size proportional to accident count at hotspots, colored by severity
- Interactive hover showing accident counts and cluster location coordinates
- Map style and zoom follow selections for great UX

---

# Insight Extraction & Hypothesis Testing

---

## 1. Interactive User Interface

---

### Page Header

- **Header:** "Insight Extraction & Hypothesis Testing with Statistical Validation"
- Presents hypotheses and data-driven conclusions with real-time visuals

---

## 2. Data Loading

---

- Loads a sample of 40,000 rows from the preprocessed dataset
- Ensures fast interaction and manageable compute for statistical tests

---

## 3. Insight 1: Effect of Weather Conditions on Accident Severity

---

- Mean severity by different weather conditions shown via bar chart
- Hypothesis: Weather impacts average severity
- Result: **Theory Proven FALSE** ( $p=0.1493$ ); no significant difference between "Clear" and "Rain" severities

---

## 4. Insight 2: Accident Frequency by Hour of Day

---

- Line chart of hourly accident counts
- Compares accidents during morning rush (7–9 am) vs. late night (12–3 am)
- Result: **Theory Proven TRUE**; more accidents occur during morning rush hours

---

## 5. Insight 3: Correlation Between Temperature and Severity

---

- Binned temperature with mean severity per bin shown in bar chart
- Pearson correlation near zero (0.002),  $p=0.676$  (non-significant)
- Result: Weak relationship, theory on temperature extremes affecting severity not supported

---

## 6. Insight 4: Accident Counts by Visibility Range

---

- Binned visibility ranges with accident counts, bar chart displayed
  - Hypothesis: Low visibility (<2mi) leads to more accidents
  - Chi-square test shows **Theory Proven TRUE** with significant association ( $p=0.0001$ )
-

## 7. Insight 5: Rain vs No Rain Accident Counts

- Binary classification of rain vs no rain from weather condition text
- Bar chart counts displayed
- Chi-square test indicates **Theory Proven TRUE**; rain significantly affects accident severity/frequency ( $p \approx 0$ )

## 8. Insight 6: Humidity vs Severity Correlation

- Pearson correlation shown (-0.001), p-value not significant (0.780)
- Result: **Theory Proven FALSE**; no significant correlation

## 9. Insight 7: Atmospheric Pressure vs Severity

- Pearson correlation -0.013,  $p=0.0096$  (statistically significant)
- Result: **Theory Proven TRUE**; pressure shows significant but weak correlation with severity

## 10. Insight 8: Road Features Effects on Severity

- Road features tested: Bump, Crossing, Give\_Way, Junction, No\_Exit, Railway, Roundabout, Station, Stop, Traffic\_Calming, Traffic\_Signal, Turning\_Loop
- For each, t-test comparing severity for presence/absence of feature
- Significant results ( $p < 0.05$ ) for most features showing impact on severity, e.g.:
  - Crossing, Give\_Way, Junction, No\_Exit, Railway, Station, Stop, Traffic\_Calming, Traffic\_Signal
- Features with insufficient data or no significance flagged
- Example: Bump showed no significant impact ( $p=0.4584$ )

# Streamlit Dashboard Summary & Key Findings Report

## 1. Total Accidents Analyzed

- Over 6.7 million accidents in the dataset were included in the analysis.
- Provides large-scale, comprehensive national coverage of US road safety.

## 2. Peak Accident Hour

- Most accidents occur around 7 AM.
- Suggests high-risk morning traffic periods requiring focus for interventions.

## 3. Top 5 Accident-Prone States

- The states with the highest number of accidents are:
  1. California (CA) - 1.5 million accidents
  2. Florida (FL)
  3. Texas (TX)
  4. South Carolina (SC)
  5. North Carolina (NC)
- These states require prioritized road safety policies and resources.

## 4. Top 5 Accident-Prone Cities

- The cities with highest accident counts include:
  1. Houston
  2. Miami
  3. Los Angeles
  4. Charlotte
  5. Dallas
- Urban traffic safety and enforcement measures suggested.

## 5. Top 5 Weather Conditions During Accidents

- Weather conditions contributing most to accidents include:
  1. Clear

- 2. Rain
- 3. Cloudy/Partly Cloudy
- 4. Light Drizzle
- 5. Light Rain
- Weather-related risk factors are prominent.

## 6. Top 5 Road Surface / Feature Conditions in Accidents

- Road conditions features most associated with accidents include:
  1. Traffic Signal
  2. Crossing
  3. Junction
  4. Stop
  5. Station
- Indicates complex road infrastructures and intersections as high-risk zones.

## Summary Report Contents

### Executive Overview

- Project scope: 6.7M+ accidents analyzed across 49 US states
- Time period: February 2016 - March 2023
- Key stakeholders and business impact

### Key Project Objectives

1. Data preprocessing with 15-step pipeline
2. Temporal pattern discovery (Peak hour: 7 AM)
3. Spatial hotspot identification (Top states & cities)
4. Environmental factor analysis (Weather impacts)
5. Road infrastructure impact assessment
6. Statistical hypothesis validation (7 of 8 proven)

### Critical Findings Table

- 8 hypothesis tests with p-values and conclusions
- 59% more accidents during morning rush vs. late night
- Visibility impact:  $p=0.0001$  (highly significant)
- Rain impact:  $p<0.001$  (extremely significant)
- Road features: 9 of 12 features significant

### Data Processing Pipeline Summary

- Complete 15-step breakdown with rows/columns affected
- Quality improvement: 60% → 95%+
- Final dataset: 6.9M rows, 35 columns, 0% missing
- Feature engineering: +6 temporal + 7 categorical features

## Technical Stack Details

### Core Technologies

Technology	Purpose
Streamlit	Interactive web dashboard
Pandas	Data manipulation & aggregation
NumPy	Numerical computing
Plotly	Interactive visualizations
SciPy	Statistical testing
Scikit-learn	ML algorithms (DBSCAN, LinearRegression)
Python 3.x	Core language

## Data Pipeline Architecture

- CSV input (6.9M records)
  - Pandas preprocessing & transformation
  - SciPy statistical analysis
  - Scikit-learn geospatial clustering
  - Plotly interactive output
- 

## Deployment Architecture

- Modular design: 5 independent Streamlit modules
  - GitHub repository structure defined
  - Local, Streamlit Cloud, and Docker deployment options
- 

This executive summary is perfect for:

- **Stakeholder presentations**
- **Academic publications**
- **Technical documentation**
- **Project portfolio**
- **Grant proposals**

## Created by

- **Name:** Infosys Springboard Internship
  - **Role:** Infosys Internship Team
  - **Date:** 28th October, 2025
  - **Organization:** Infosys Springboard
  - **Contact:** [springboardmentor942x@gmail.com](mailto:springboardmentor942x@gmail.com)
  - **Github Link:** [RoadSafe Analytics GitHub Repository](#)
- 

## Acknowledgments

This report was prepared as part of the **US RoadSafe Analytics** project, utilizing contributions from the data science and analytics team, with project creation dates spanning from initial data collection to final documentation in October 2025.

---

## Credits & Development

- Data preprocessing pipeline and Analysis report created by **Infosys Springboard Intern Team**.
- Visualization and Dashboard Development by **Infosys Springboard Intern Team**.
- Dataset Source: [Kaggle US Accidents Dataset](#)
- Model and Analysis techniques: Python , Pandas , Plotly , SciPy , scikit-learn