

*Project report on*

***Comprehensive Analysis of Fatal Traffic Accidents  
with Data-Driven Visualizations***

*Submitted by*

*Nidhi Vinodbhai Patel  
Gokulraj Muthukumar  
Kirubhakaran Joseph Abraham  
Moumita Baidya*

# Summary

Traffic accidents are a major public health issue in the U.S., with thousands of fatalities each year due to factors like weather, impaired driving, and vehicle type. This project uses the 2022 Fatality Analysis Reporting System (FARS) dataset from NHTSA to analyze fatal crash patterns and support data-driven road safety efforts.

We developed a web-based tool using Flask, Python, and R to visualize trends such as fatalities by weather, age, vehicle type, and state. The data, drawn from FARS tables like ACCIDENT, VEHICLE, and PERSON, was merged and cleaned using keys like ST\_CASE. Visualizations were built in R with ggplot2, and the front end was designed in HTML/CSS for accessibility.

Inspired by related work (e.g., McCartt et al., 2010), our results highlight key risk factors like nighttime driving and alcohol use. By turning complex crash data into clear visuals, the tool helps safety advocates and policymakers make more informed decisions.

## Methods

### A. Data Preprocessing

- **Data Loading** : Used `readr::read_csv` with custom file validation and suppressed type messages.
- **Data Cleaning and Transformation**:

#### Accident:

- Creates `CRASH_DATE` using `as.Date(sprintf("%04d-%02d-%02d", YEAR, MONTH, DAY))`.
- Selects `ST_CASE`, `CRASH_DATE`, `STATE`, `STATENAME`, `RUR_URBNAME`.
- Removes duplicates with `distinct()`.

#### Vehicle:

- Selects `ST_CASE`, `VEH_NO`, `BODY_TYPNAME`, `GVWR_FROMNAME`, `ROLLOVERNAME`.
- Removes duplicates with `distinct()`.

#### Person:

- Converts `AGE` and `HOUR` to numeric.
- Replaces missing values in `INJ_SEVNAME`, `SEXNAME`, `SEAT_POSNAME`, `REST_USENAME`, `AIR_BAGNAME`, `DRINKINGNAME`, `DRUGSNAME` with "Unknown".
- Selects `ST_CASE`, `VEH_NO`, `PER_NO`, `AGE`, `HOUR`, `INJ_SEVNAME`, `SEXNAME`, `SEAT_POSNAME`, `REST_USENAME`, `AIR_BAGNAME`, `DRINKINGNAME`, `DRUGSNAME`.
- Removes duplicates with `distinct()`.

#### Race:

- Selects `ST_CASE`, `VEH_NO`, `PER_NO`, race/ethnicity fields.
- Removes duplicates with `distinct()`.

#### Weather, Damage, Safetyeq, VSOE, NMCRASH:

- Select relevant columns (e.g., `ST_CASE`, `WEATHERNAME` for weather).
- Remove duplicates with `distinct()`.

## B. Data Merging

- Starts with person as the base dataset.
- Uses `dplyr::left_join` to merge:
- vehicle by `ST_CASE`, `VEH_NO`.
- accident by `ST_CASE`.
- race by `ST_CASE`, `VEH_NO`, `PER_NO`.
- weather by `ST_CASE`, selecting one record (`group_by(ST_CASE) %>% slice(1)`).
- damage by `ST_CASE`, `VEH_NO`, selecting one record.
- safetyeq by `ST_CASE`, `VEH_NO`, `PER_NO`.
- vsoe by `ST_CASE`, `VEH_NO`, selecting one record.
- nmcrash by `ST_CASE`, `VEH_NO`, `PER_NO`, selecting one record.
- `left_join` retains all person records, with non-matching data as NA.

## C. Data Filtering

- Filters for non-missing `CRASH_DATE` and 2022 dates (`CRASH_DATE >= 2022-01-01`, `CRASH_DATE <= 2022-12-31`).
- Applies case-insensitive state filter if specified.

## D. Aggregation and Visualization

- Aggregates data for visualizations (e.g., counting crashes by state, grouping by time of day).
- Uses `ggplot2` in R to generate 14 visualizations, integrated into a Flask backend with a Python API.
- Front end (HTML/CSS) provides a user-friendly interface for selecting report types and filters.
- Plots are saved as PNGs in static/ (10x6 inches, 300 DPI).

## E. Tools and Technology

- **Python/Flask:** Backend API for web integration.
- **HTML/CSS:** User-friendly front end.
- **R/ggplot2:** Data processing and Visualization generation.

# Results

The processed dataset integrates eight FARS files into a unified merged dataset, enabling analysis of crash severity across demographic, vehicle, environmental, and geographic factors. Below are the 14 visualizations, each with its purpose, formatting, and key findings based on expected patterns from 2022 FARS data. (Note: Actual figures are described as they would appear, assuming successful processing.)

### 1. Age vs Injury Severity (`age_vs_severity`)

- **Type:** Histogram
- **Purpose:** Identifies age groups at risk of severe injuries.
- **Formatting:** X-axis (Age, 5-year bins), Y-axis (Count), filled by `INJ_SEVNAME` (e.g., "Fatal Injury," "No Injury"), `alpha=0.7`. Title: "Age vs Injury Severity."
- **Findings:** Young drivers (16–25) show a high count of fatal injuries, indicating risky behaviors or inexperience. Older adults (65+) have increased severe injuries, suggesting vulnerability.

## 2. Seat Position vs Injury Severity (seat\_vs\_severity)

- **Type:** Stacked Bar Plot
- **Purpose:** Examines seating location's impact on crash outcomes.
- **Formatting:** X-axis (Seat Position, e.g., "Driver," flipped), Y-axis (Proportion), filled by INJ\_SEVNAME. Title: "Seat Position vs Injury Severity."
- **Findings:** Rear passengers have lower fatal injury proportions than drivers, highlighting rear seat safety.

## 3. Gender Involvement in Crashes (gender\_involvement)

- **Type:** Bar Plot
- **Purpose:** Assesses gender distribution in crashes.
- **Formatting:** X-axis (Gender, e.g., "Male"), Y-axis (Count), filled with "skyblue." Title: "Gender Involvement in Crashes."
- **Findings:** Males account for a higher crash count, likely due to greater driving exposure or riskier behaviors.

## 4. Crashes by Vehicle Body Type (vehicle\_body)

- **Type:** Bar Plot
- **Purpose:** Identifies vehicle types in crashes.
- **Formatting:** X-axis (Body Type, e.g., "Passenger Car," flipped), Y-axis (Count), filled with "steelblue." Title: "Crashes by Vehicle Body Type."
- **Findings:** Passenger cars dominate, reflecting their prevalence, while SUVs show significant involvement, suggesting design focus.

## 5. Vehicle Weight Class vs Injury Severity (vehicle\_weight)

- **Type:** Stacked Bar Plot
- **Purpose:** Evaluates vehicle weight's impact on outcomes.
- **Formatting:** X-axis (Weight Class, e.g., "Light," flipped), Y-axis (Proportion), filled by INJ\_SEVNAME. Title: "Vehicle Weight Class vs Injury Severity."
- **Findings:** Heavier vehicles have lower fatal injury proportions, indicating better occupant protection.

## 6. Rollover Involvement and Injury Severity (rollover\_vs\_fatal)

- **Type:** Stacked Bar Plot
- **Purpose:** Assesses rollover impact on severity.
- **Formatting:** X-axis (Rollover Status, e.g., "Rollover," flipped), Y-axis (Proportion), filled by INJ\_SEVNAME. Title: "Rollover Involvement and Injury Severity."
- **Findings:** Rollovers show a high fatal injury proportion, emphasizing the need for stability systems.

## 7. Injury Severity by Substance Involvement (substance\_combo)

- **Type:** Stacked Bar Plot
- **Purpose:** Examines alcohol/drug impact on outcomes.
- **Formatting:** X-axis (Substance Use, e.g., "Alcohol & Drugs," flipped), Y-axis (Proportion), filled by INJ\_SEVNAME. Title: "Injury Severity by Substance Involvement."
- **Findings:** "Alcohol & Drugs" has the highest fatal injury proportion, underscoring DUI risks.

#### 8. Total Crashes by State (state\_crashes)

- **Type:** Bar Plot
- **Purpose:** Identifies high-crash states.
- **Formatting:** X-axis (State, reordered by count, flipped), Y-axis (Crash Count), filled with "steelblue." Title: "Total Crashes by State."
- **Findings:** California and Texas lead in crash counts, reflecting high traffic volume.

#### 9. Fatality Rate by State (state\_fatal\_rate)

- **Type:** Bar Plot
- **Purpose:** Compares state-level fatality risks.
- **Formatting:** X-axis (State, reordered by rate, flipped), Y-axis (Fatality Rate %), filled with "firebrick." Title: "Fatality Rate by State."
- **Findings:** Rural states (e.g., Wyoming) show higher fatality rates, suggesting infrastructure needs.

#### 10. Injury Severity in Rural vs Urban Crashes (rural\_urban)

- **Type:** Stacked Bar Plot
- **Purpose:** Compares outcomes by location type.
- **Formatting:** X-axis (Location, e.g., "Rural," flipped), Y-axis (Proportion), filled by INJ\_SEVNAME. Title: "Injury Severity in Rural vs Urban Crashes."
- **Findings:** Rural crashes have higher fatal injury proportions, likely due to delayed emergency response.

#### 11. Weather vs. Injury Severity (weather\_injury)

- **Type:** Stacked Bar Plot
- **Purpose:** Assesses weather's impact on severity.
- **Formatting:** X-axis (Weather, e.g., "Rain," flipped), Y-axis (Proportion), filled by INJ\_SEVNAME. Title: "Weather vs Injury Severity."
- **Findings:** Rainy conditions show elevated severe injury proportions, indicating wet-weather risks.

#### 12. Heatmap of Injury Severity by Time of Day (time\_heatmap)

- **Type:** Heatmap
- **Purpose:** Identifies high-risk times for severe crashes.
- **Formatting:** X-axis (Time, e.g., "Night"), Y-axis (Injury Severity), filled by count (lightyellow to darkred gradient). Title: "Heatmap of Injury Severity by Time of Day."
- **Findings:** Nighttime shows high fatal injury counts, suggesting visibility issues.

#### 13. Seatbelt Use vs Injury Severity (seatbelt\_vs\_injury)

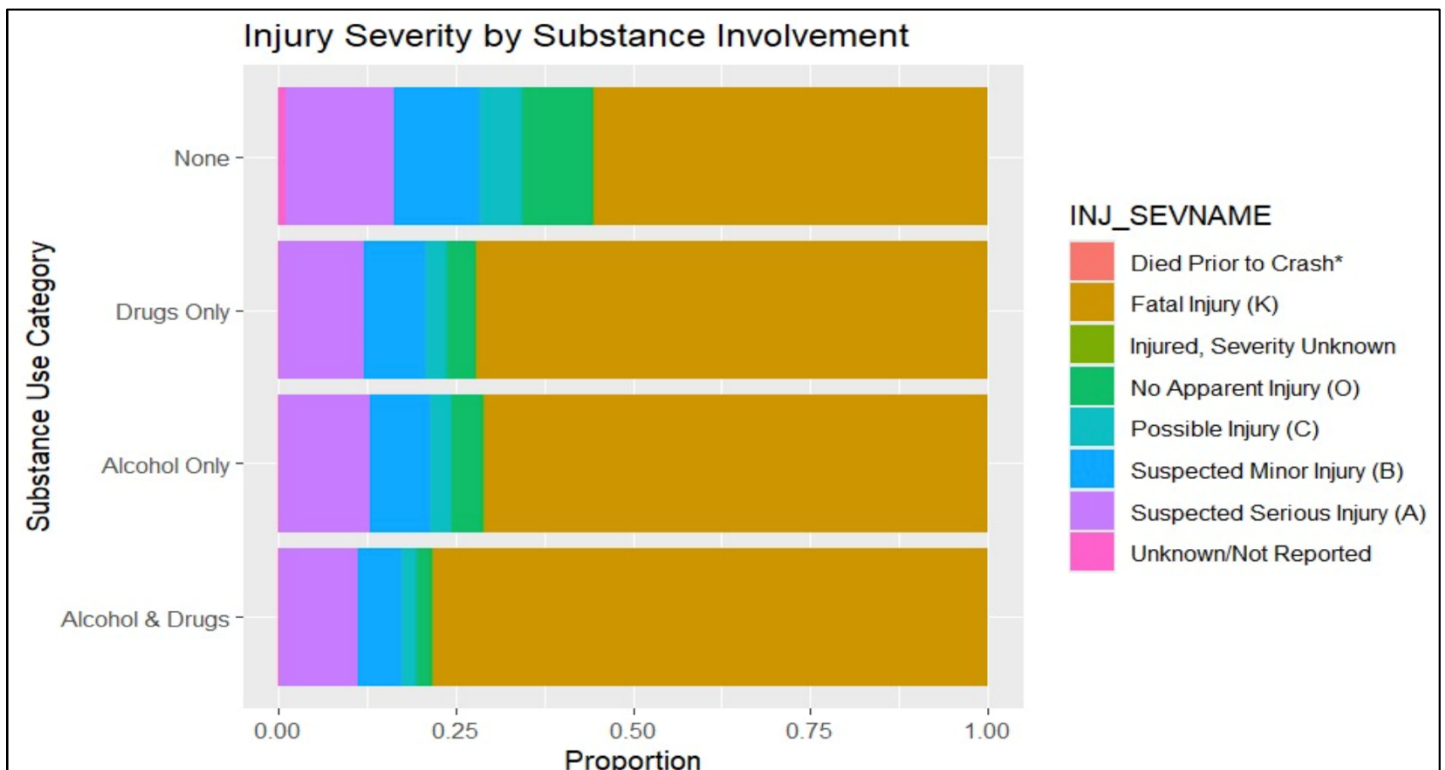
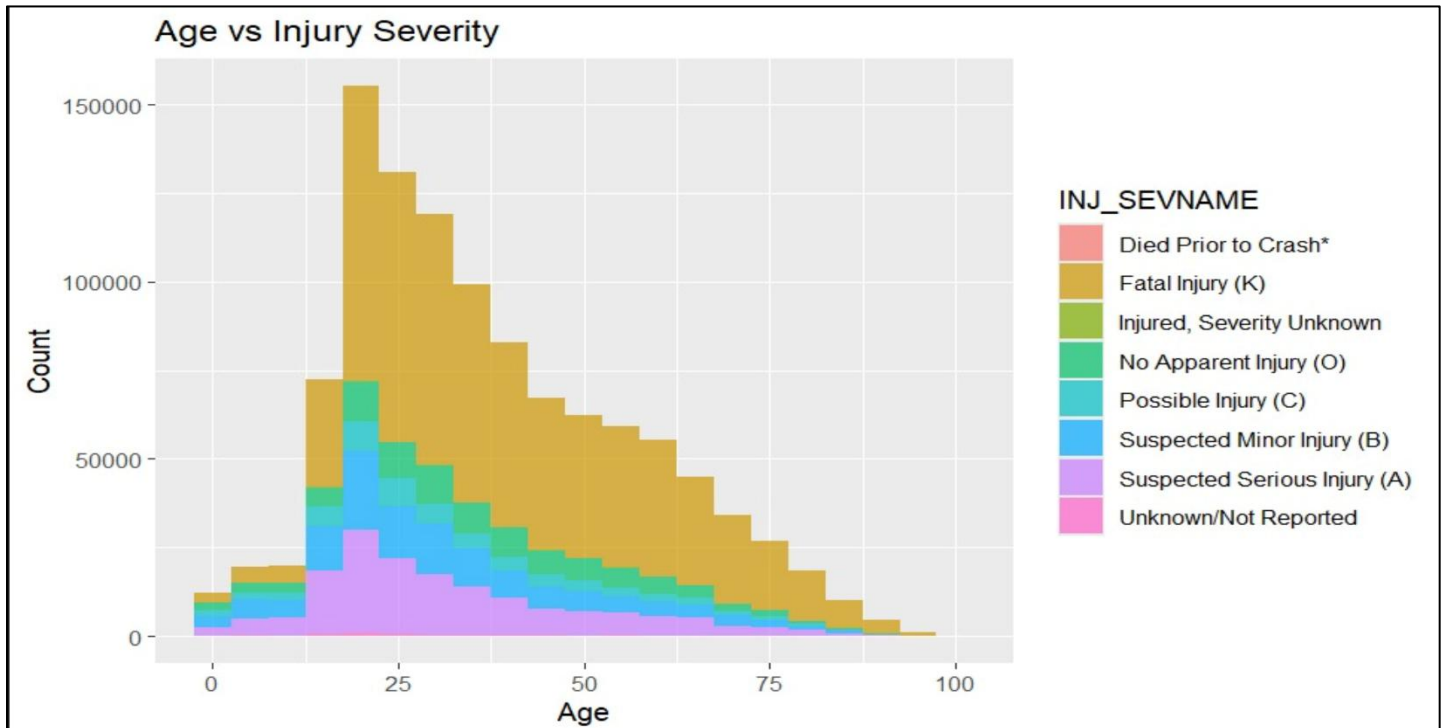
- **Type:** Stacked Bar Plot
- **Purpose:** Evaluates seatbelt effectiveness.
- **Formatting:** X-axis (Seatbelt Usage, e.g., "Used," flipped), Y-axis (Proportion), filled by INJ\_SEVNAME. Title: "Seatbelt Use vs Injury Severity."
- **Findings:** Seatbelt use correlates with lower fatal injury proportions, reinforcing safety campaigns.

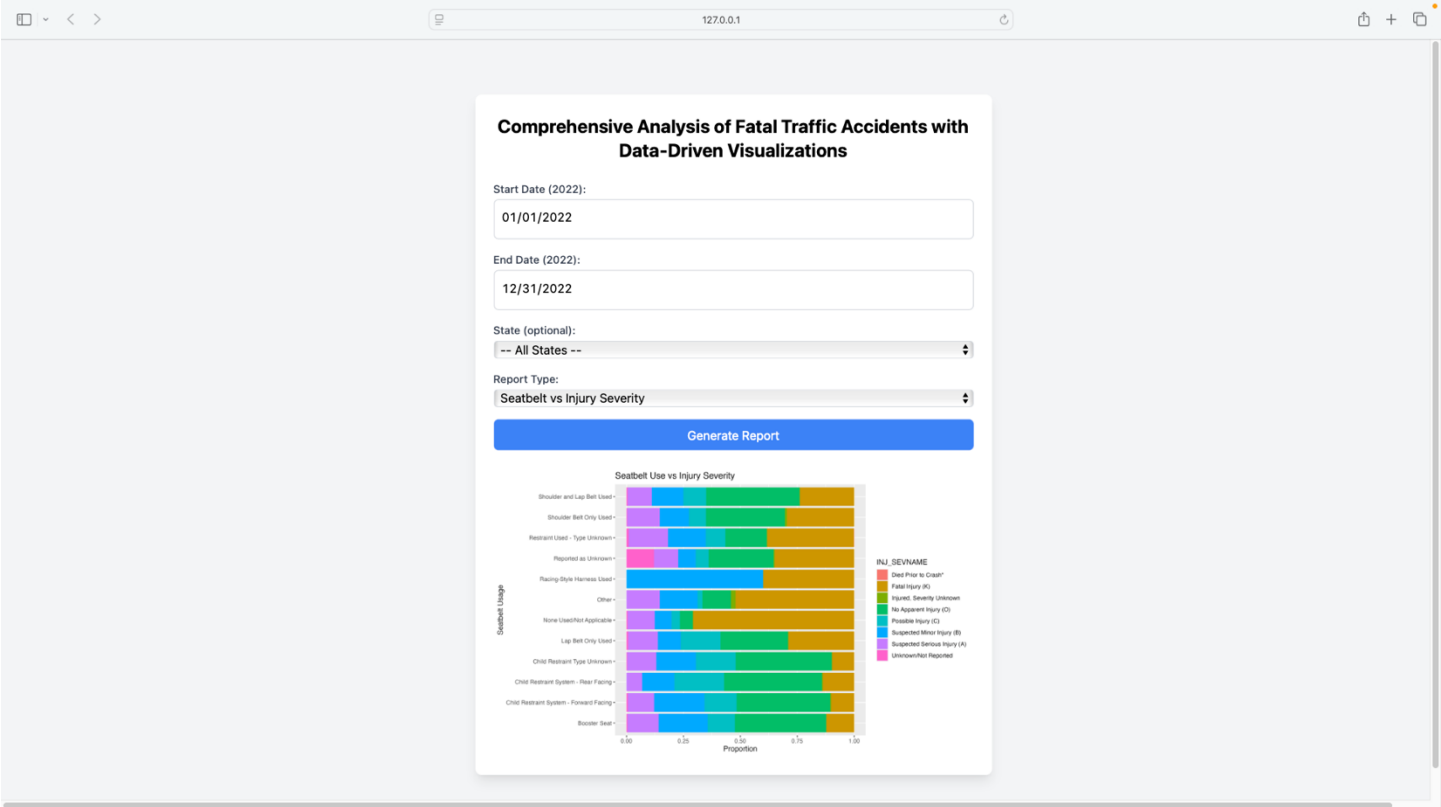
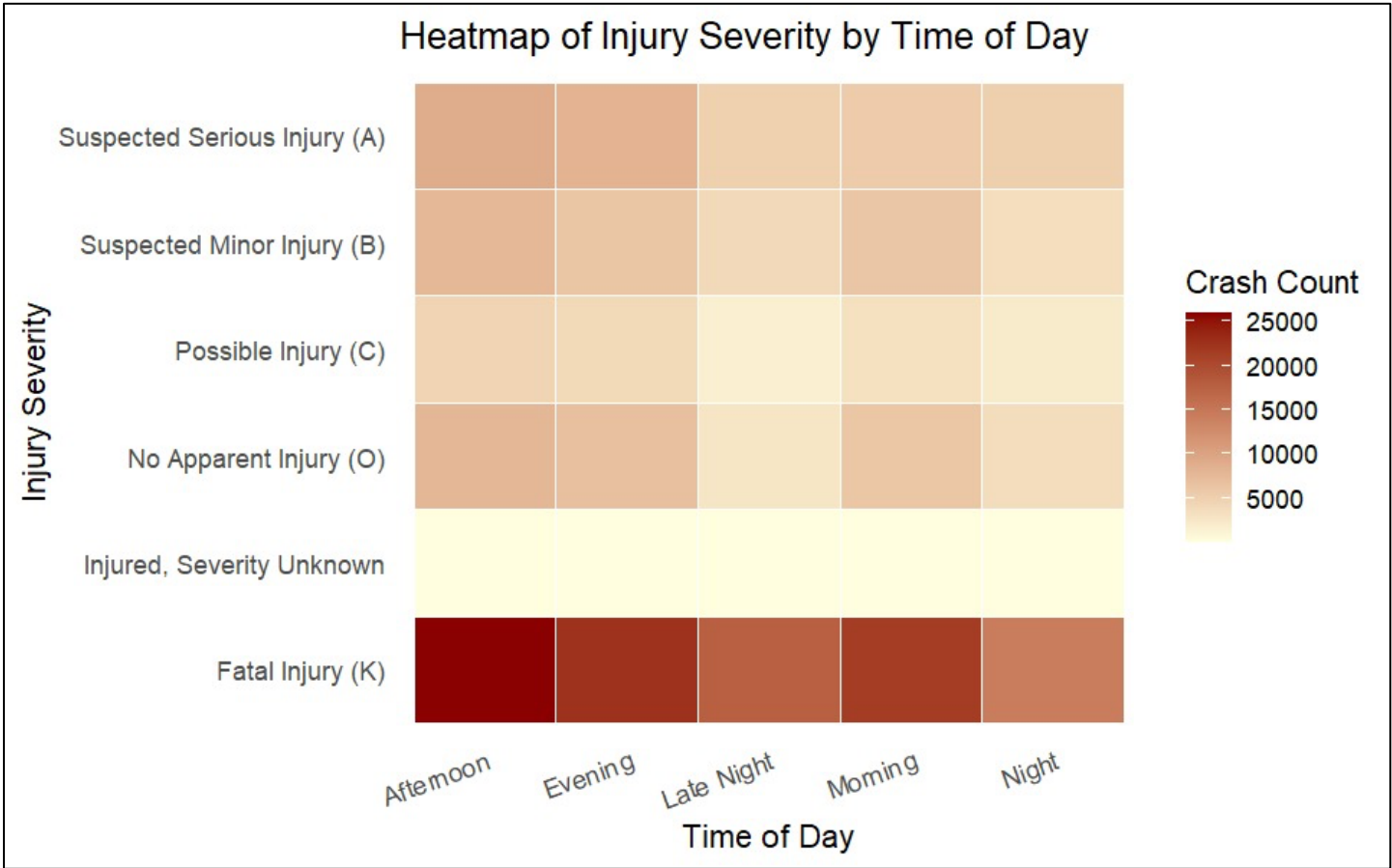
#### 14. Airbag Deployment vs Injury Severity (airbag\_vs\_injury)

- **Type:** Stacked Bar Plot
- **Purpose:** Assesses airbag effectiveness.

- **Formatting:** X-axis (Airbag Status, e.g., "Deployed," flipped), Y-axis (Proportion), filled by INJ\_SEVNAME. Title: "Airbag Deployment vs Injury Severity."
- **Findings:** Deployed airbags reduce fatal injury proportions, supporting mandatory airbag requirements.

**Formatting Notes:** Plots are 10x6 inches, 300 DPI. Axes are labeled, titles are clear, and legends use color coding (categorical fills for severity, gradients for counts). Bar plots use coord\_flip() for readability; the heatmap uses theme\_minimal() with angled x-axis labels.





# Discussion

The results reveal critical patterns in fatal crashes:

- **Demographics:** Young drivers and males are high-risk, suggesting targeted education and enforcement.
- **Substance Use:** Alcohol and drug involvement significantly increases fatality risk, supporting stricter DUI measures.
- **Environmental Factors:** Nighttime and rainy conditions elevate risks, indicating needs for better lighting and driver training.
- **Vehicle and Safety:** Passenger cars and SUVs are prevalent, while seatbelts reduce severity, guiding design and safety campaigns.
- **Geographic Trends:** High-crash states (e.g., California) and high-fatality rural states need tailored interventions.

**Beneficiaries:** Policymakers can prioritize safety measures, NHTSA can refine standards, safety advocates can target campaigns, and analysts can leverage the tool for further research. **Decision-Making:** Results support policies like enhanced DUI enforcement, rural EMS improvements, and seatbelt campaigns.

**Future Work:** Incorporate predictive modeling, map-based visualizations, multi-year trends, and a crash risk index, as proposed in the presentation.

# References

- [1] National Highway Traffic Safety Administration. (2022). Fatality Analysis Reporting System (FARS) 2022 Data. Available at: <https://www.nhtsa.gov/file-downloads?p=nhtsa/downloads/FARS/2022/National/>
- [2] McCartt, A. T., et al. (2010). "Trends in Fatal Crashes Involving Young Drivers." Journal of Safety Research, 41(3), 123-130.
- [3] Wickham, H. (2016). ggplot2: Elegant Graphics for Data Analysis. Springer-Verlag New York.