

The background of the image is a soft-focus photograph of an urban environment. It features a large, modern-looking bridge with a dark, textured surface and a steel frame. Below the bridge, there's a street with some greenery and what appears to be a bus stop or a small building. The sky is a clear, pale blue.

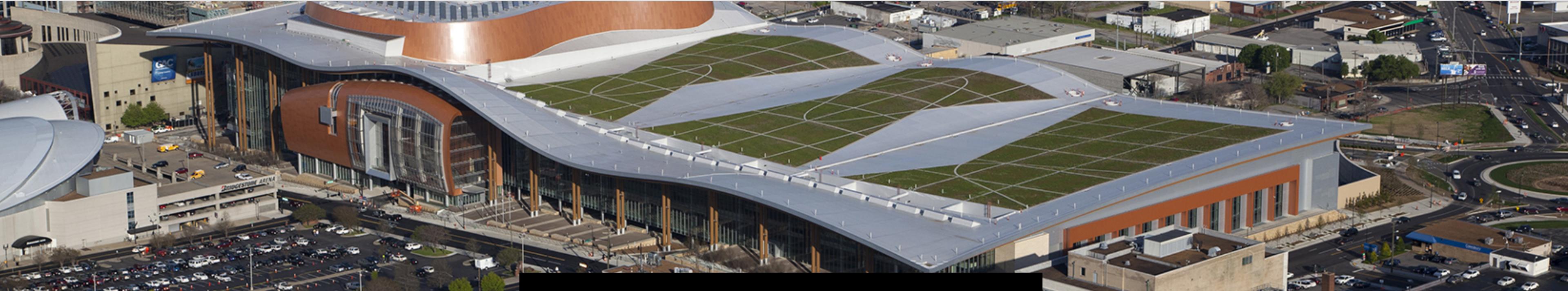
Springboard—DSC

Capstone Project 2

# PREDICTING PUBLIC TRANSPORTATION SAFETY RISK

[H T T P S : / / G I T H U B . C O M / T A M A R A H O R N E / S P R I N G B O A R D / T R E E / M A I N / C A P S T O N E % 2 0 P R O J E C T % 2 0 2](https://github.com/tamarahorne/springboard/tree/main/capstone%20project%202)

# INTRODUCTION



WHAT'S THE SAFETY RISK?

**Increasing Public Transportation Options  
Between Nashville's Convention Center and Airport**

# THE DATA

## National Transit Database

- Monthly modal time series data
- 133,196 rows; 65 columns
- One row per month per agency per mode

## Wrangling

- Filled NaNs in 4000+ with data found in the dataframe
- Filled NaNs in 27 rows with data from FTA

# REDUCING THE COLUMNS

## ID Columns



5 Digit NTD ID



4 Digit NTD ID

## Location Columns



Primary UZA Population



Primary UZA Code

## Totals Columns



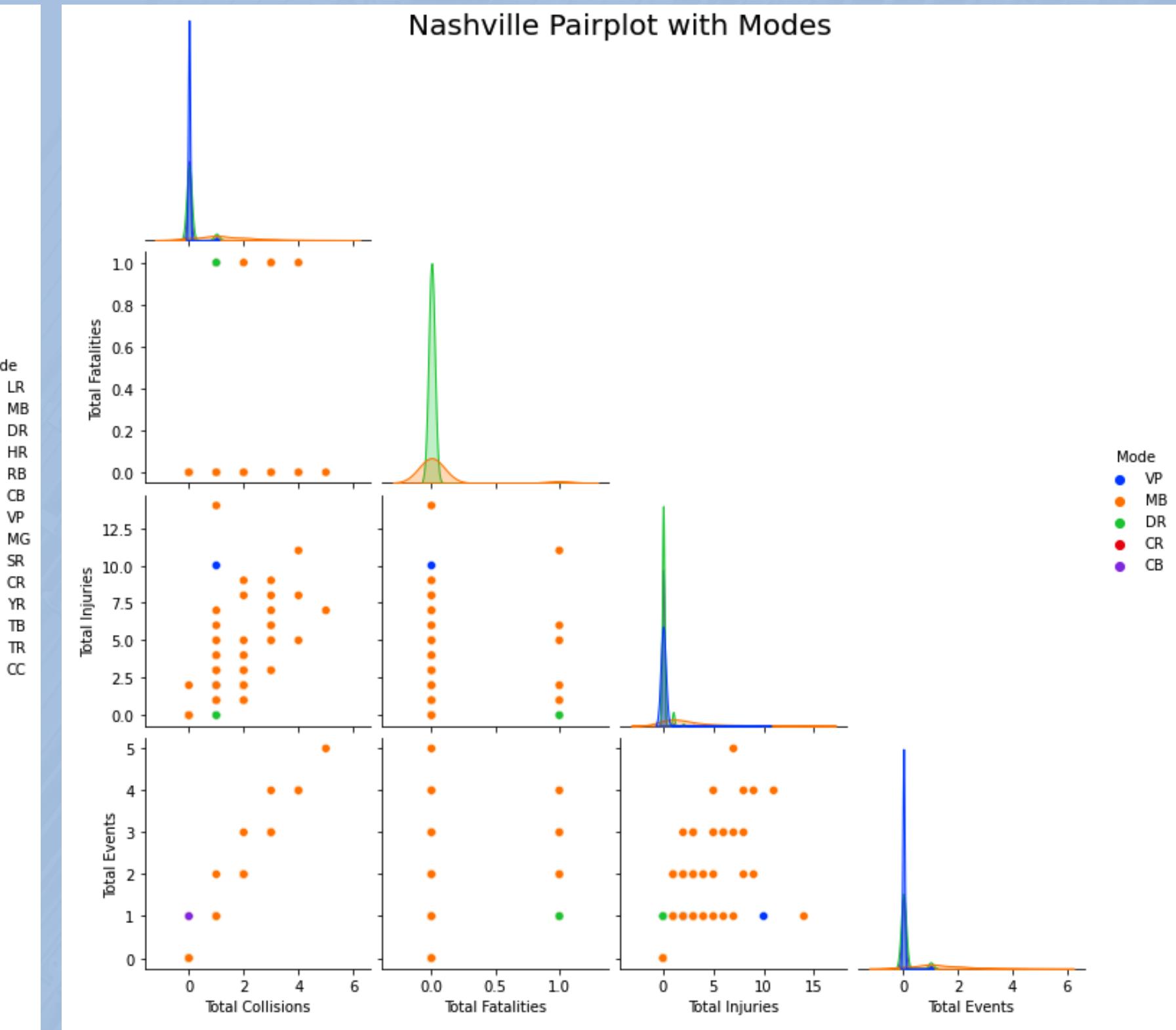
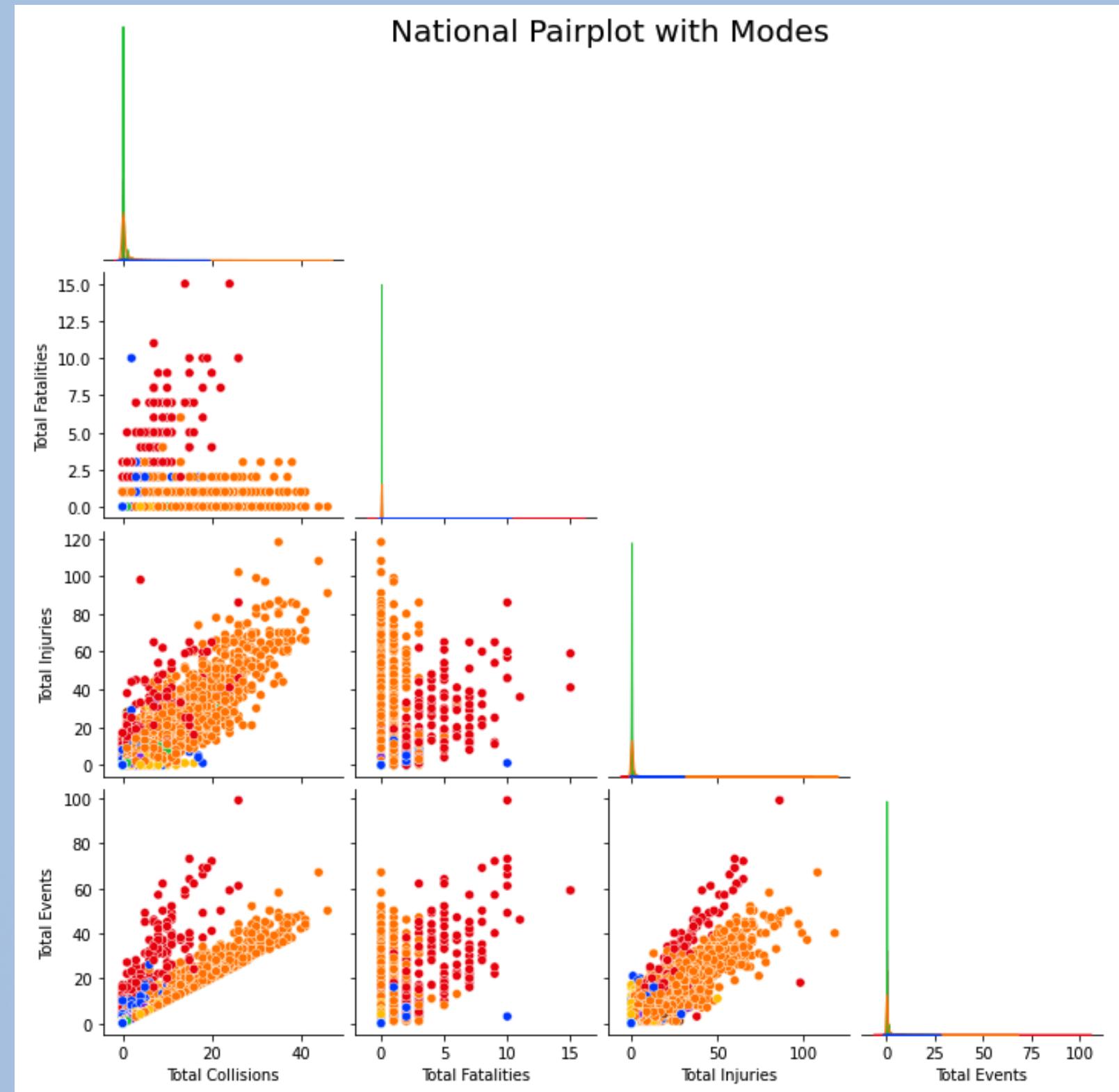
Total Fatalities  
Total Injuries  
Total Events



All contributing columns

# MANY SAFETY INCIDENTS FOR MB AND HR

TAMARA HORNE  
2023 | March



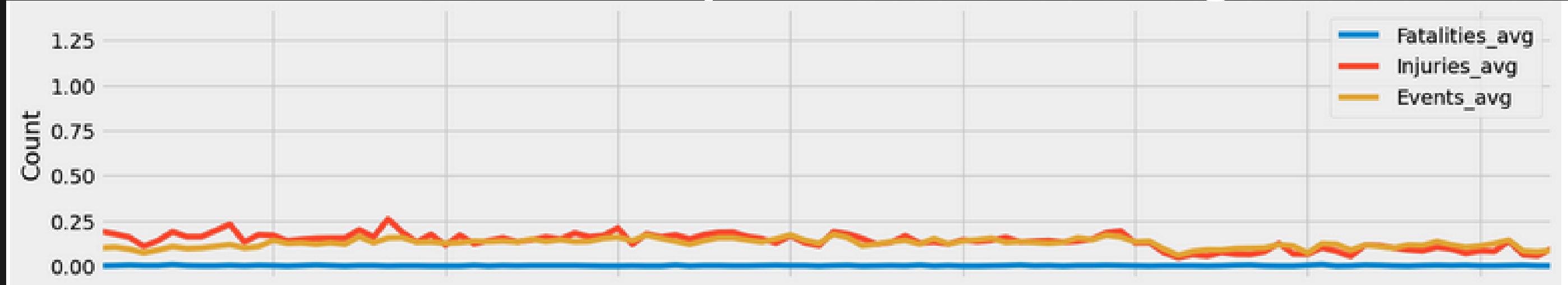
# THE TARGET

<b>Safety Performance Targets as Reported to the National Transit Database (NTD)</b>							
The targets listed below are based on reviews of the previous five years of MTA dba WeGo Public Transit's safety performance data.							
<b>Mode of Transit Service</b>	<b>Fatalities (total)</b>	<b>Fatalities (per 100 thousand VRM)</b>	<b>Injuries (total)</b>	<b>Injuries (per 100 thousand VRM)</b>	<b>Safety Events (total)</b>	<b>Safety Events (per 100 thousand VRM)</b>	<b>System Reliability (VRM / failures)</b>
<b>Fixed Route Bus</b>	0	0	35	.55	24	.45	5,500
<b>Demand Response Bus</b>	0	0	6	.27	6	.26	24,800
<b>Demand Response Taxi</b>	0	0	0	0	0	0	0

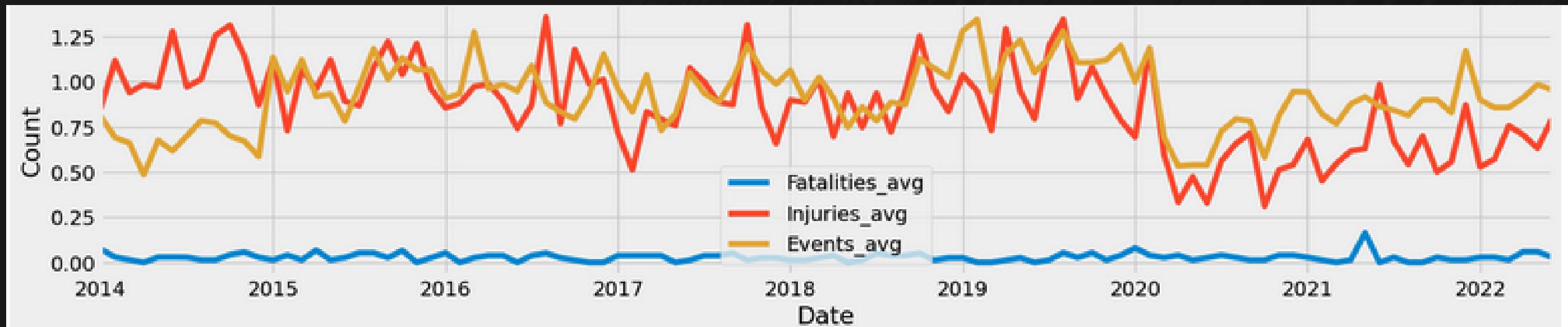
**(TOTAL FATALITIES + TOTAL INJURIES + TOTAL EVENTS) / VEHICLE REVENUE MILES**

# EXPLORATORY AGGREGATION

*Locations with Modes up to and Including Nashville's*



*Locations with Nashville's Some or All of Modes Plus Light Rail*

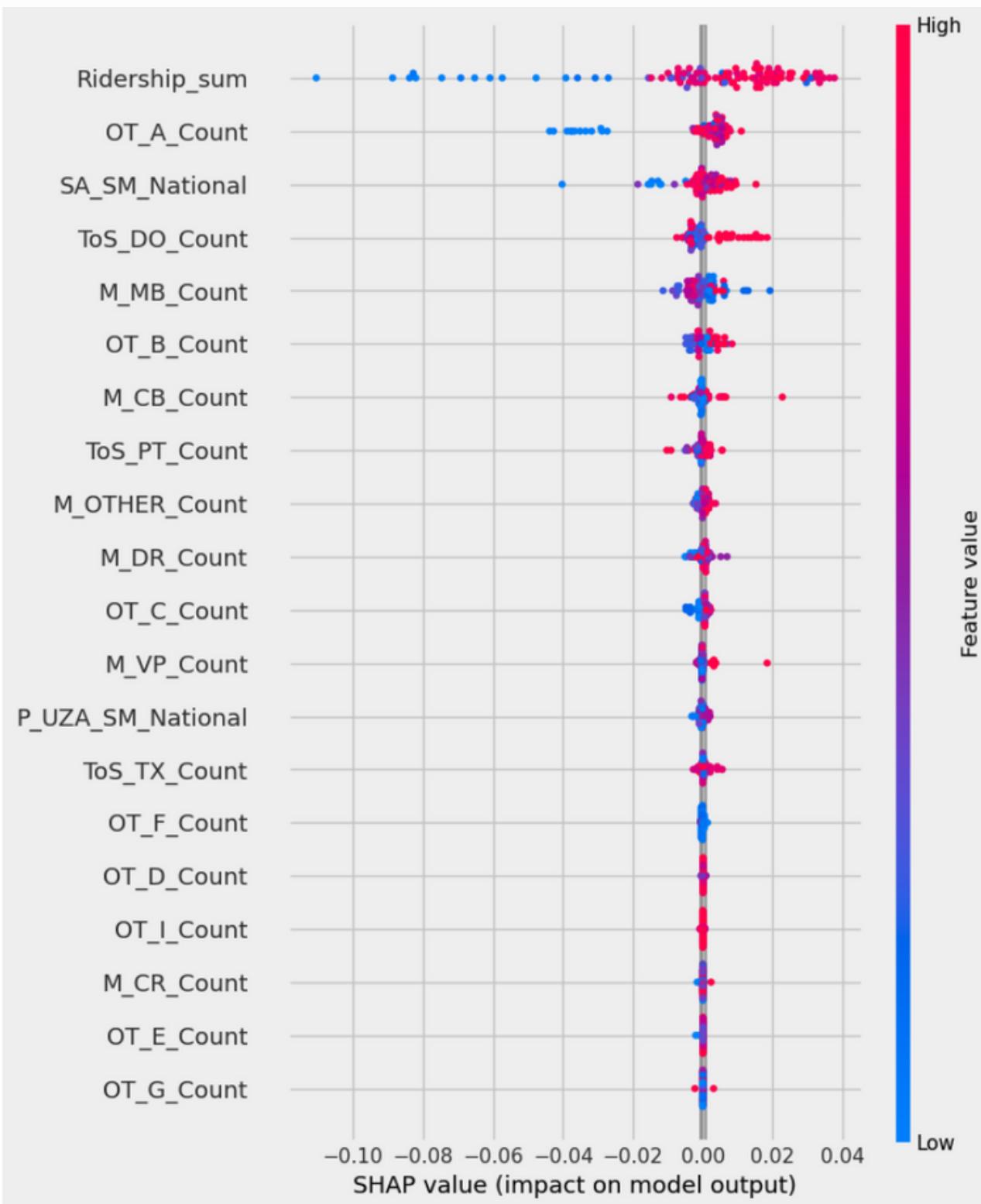


<i>~modeling_df~</i>	R Squared	MAPE	MAE	MSE
Linear Regression	0.645	0.127	0.058	0.007
Random Forest Regressor	0.902	0.071	0.033	0.002
K Neighbors Regressor	0.365	0.075	0.035	0.002
XGBoost Regressor	0.999	0.048	0.022	0.001

# MODELING

*First Aggregation*    Date, nation

What features influence safety risk for the nation as a whole?



# FEATURE IMPORTANCE

## Ridership\_sum

Lower values decrease safety risk

## OT\_A\_Count

'Independent Public Agency or Authority of Transit Service' organization type.  
Most frequently occurring organization type in the national data

## SA\_SM\_National

Lower values decrease safety risk

Average service area square miles for the nation during each given month

Lower values decrease safety risk

<i>~modeling_df2~</i>	R Squared	MAPE	MAE	MSE
Linear Regression	0.251		0.425	0.656
Random Forest Regressor	0.833		0.522	0.871
K Neighbors Regressor	0.280		0.395	0.880
XGBoost Regressor	0.548		0.454	0.808

# MODELING

## *Second Aggregation*

Date, location

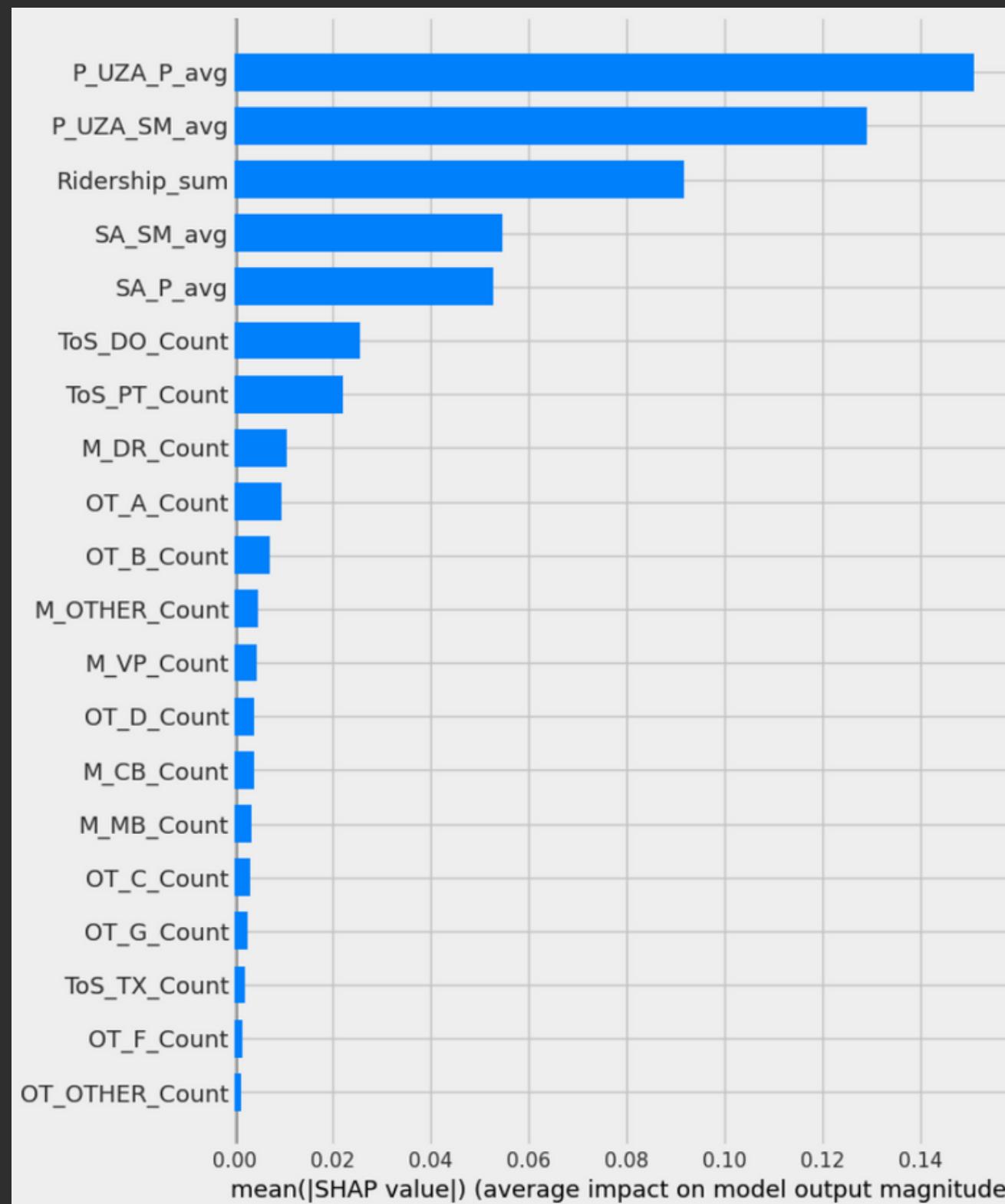
What features influence safety risk when data is aggregated to both date and location (Primary UZA Name)?



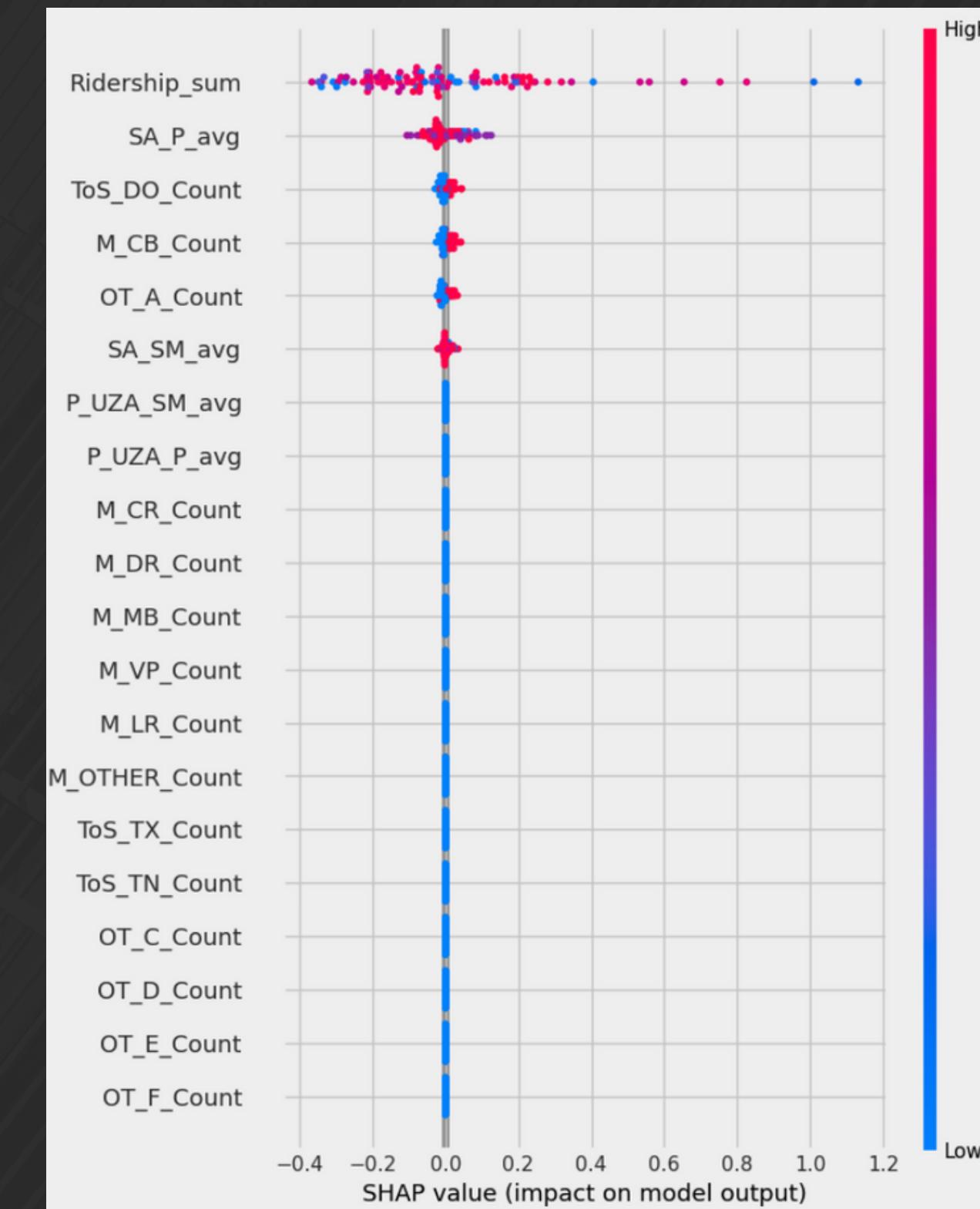
**SHIFTED RESIDUALS MAKE THE RESULTS LESS TRUSTWORTHY**

# FEATURE IMPORTANCE COMPARISON

## National



## Nashville



<b>~modeling_df3~</b>	<b>R Squared</b>	<b>MAPE</b>	<b>MAE</b>	<b>MSE</b>
Linear Regression	0.210		1.403	3.587
Random Forest Regressor	0.883		0.680	2.335
K Neighbors Regressor	0.595		0.591	2.221
XGBoost Regressor	0.945		0.744	2.587

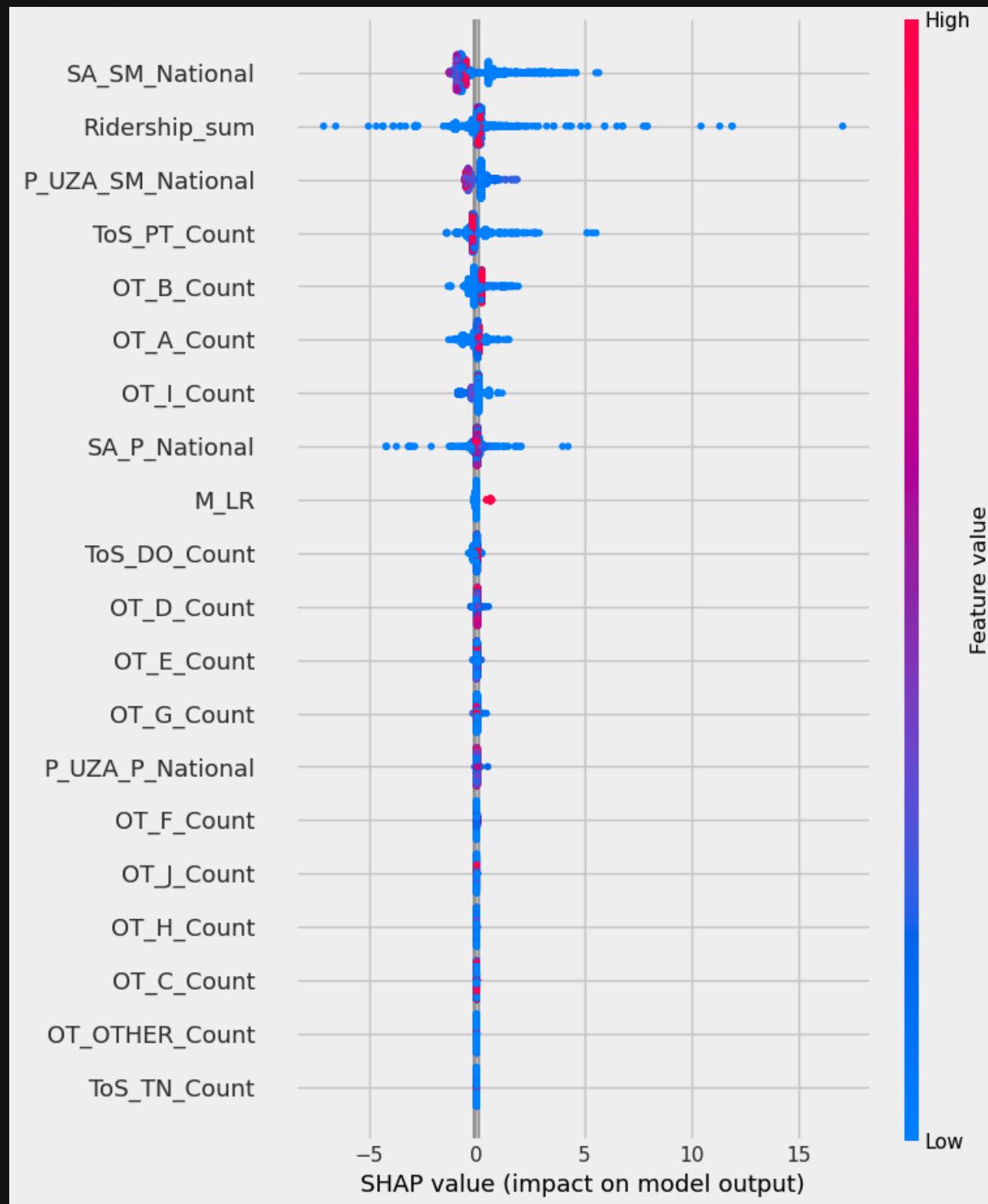
# MODELING

*Third Aggregation* Date, mode

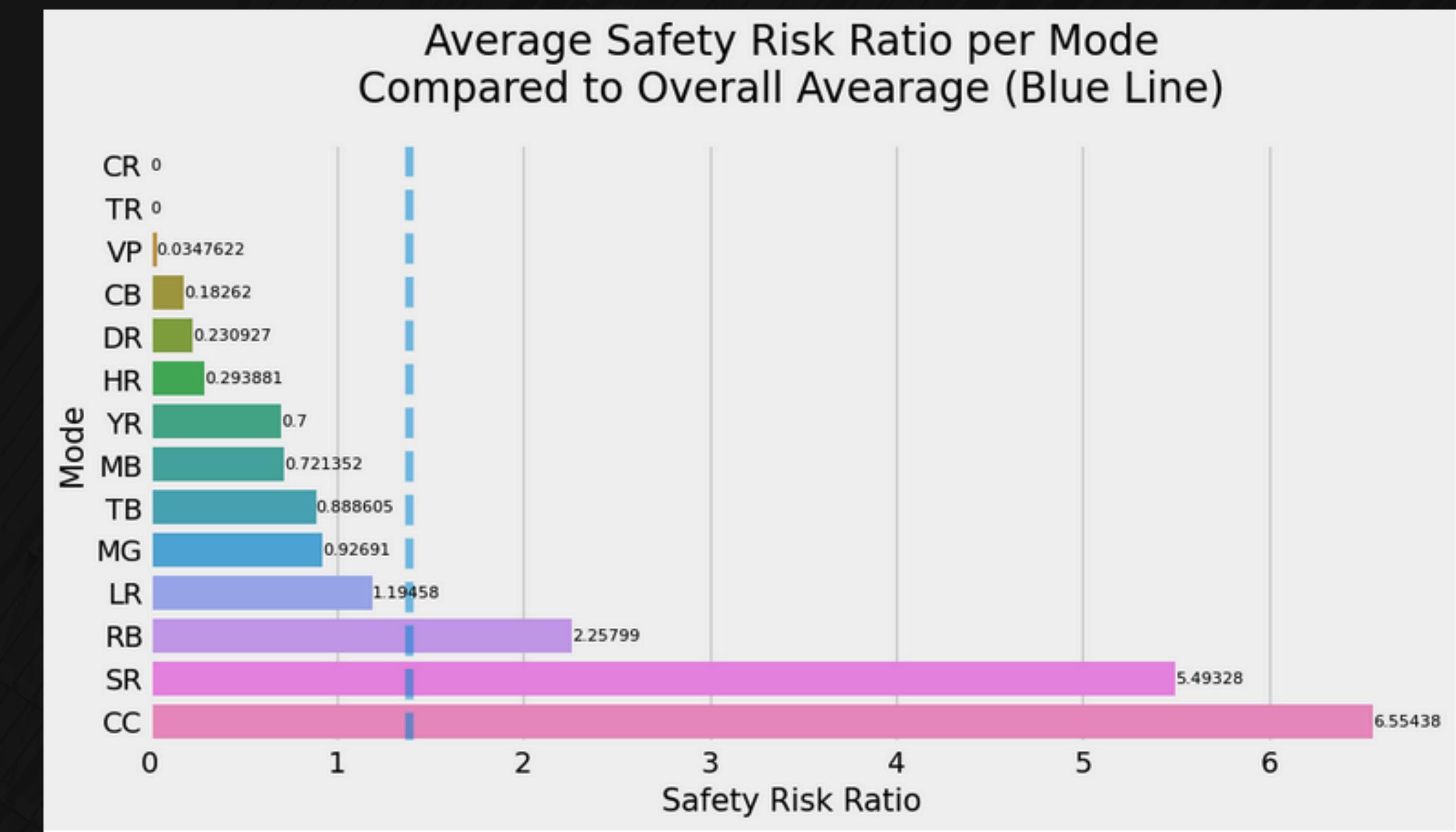
What features influence safety risk for the individual modes?

Which modes have the highest and lowest safety risk?

## FEATURE IMPORTANCE



## MODAL SAFETY RISK



# CONCLUSIONS

## Best Models

- XGBoost Regressor
- Random Forest Regressor

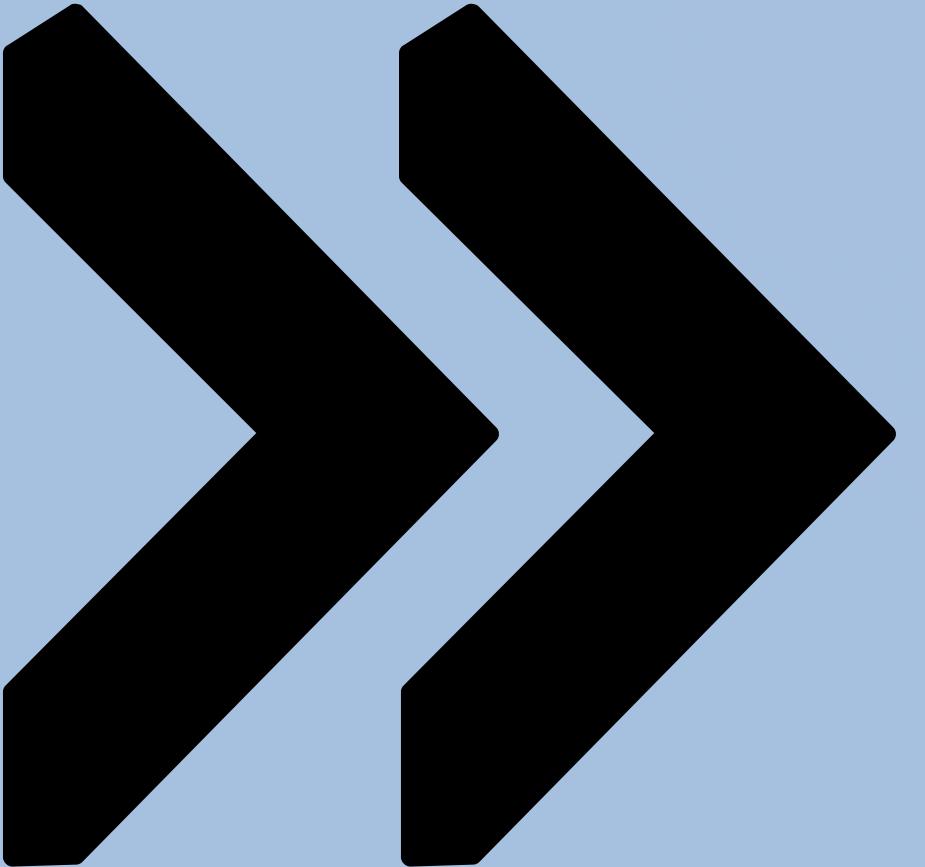
## Aggregations

- Date, Nation
- Date, Location
- Date, Mode

## Prediction Power

- Safety risk overall for the nation
- Safety risk for Nashville
- Safety risk for specific modes

# FUTURE WORK



## *Small*

- Diagnose shifted data in modeling\_df2
- Explore two low ridership/high SHAP values for modeling\_df2

## *Big*

- Import data needed for System Reliability and incorporate into the definition of the target
- Time series modeling

# RECOMMENDATIONS FOR THE CLIENT



- Adding light rail would increase safety risk
- Adding service to existing modes would be a good choice
  - Added vehicle revenue miles will differ according to mode; this will affect the target
  - Safety risk is lowest for CR and highest for MB

# MORE INFORMATION

---

<https://github.com/tamarahorne/Springboard/tree/main/Capstone%20Project%202>

