

Springboard—DSC

Capstone Project 2

# PREDICTING PUBLIC TRANSPORTATION SAFETY RISK

[HTTPS://GITHUB.COM/TAMARAHORNE/SPRINGBOARD/TREE/MAIN/CAPSTONE%20PROJECT%202](https://github.com/tamarahorne/springboard/tree/main/capstone%20project%202)

# INTRODUCTION

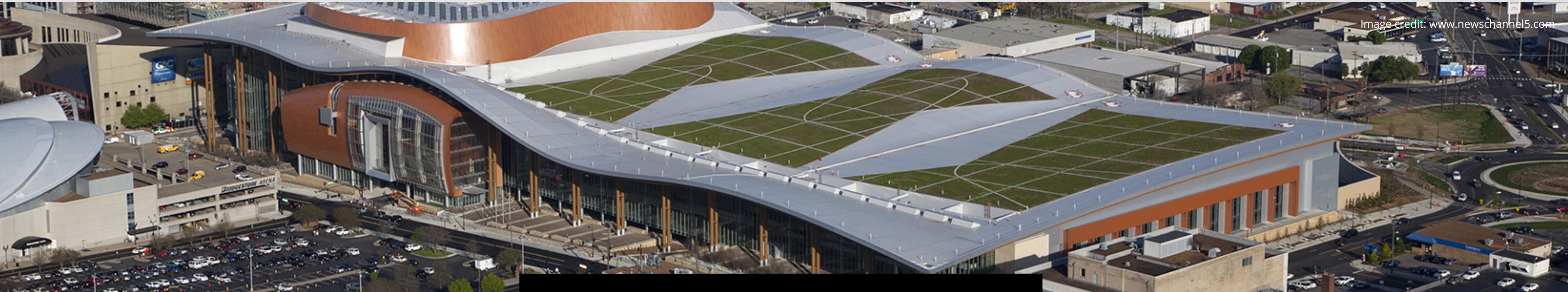


Image credit: www.newschannel5.com

WHAT'S THE SAFETY RISK?

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

TAMARA HORNE  
2023 | March

# 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 (Federal Transit Administration)

# 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

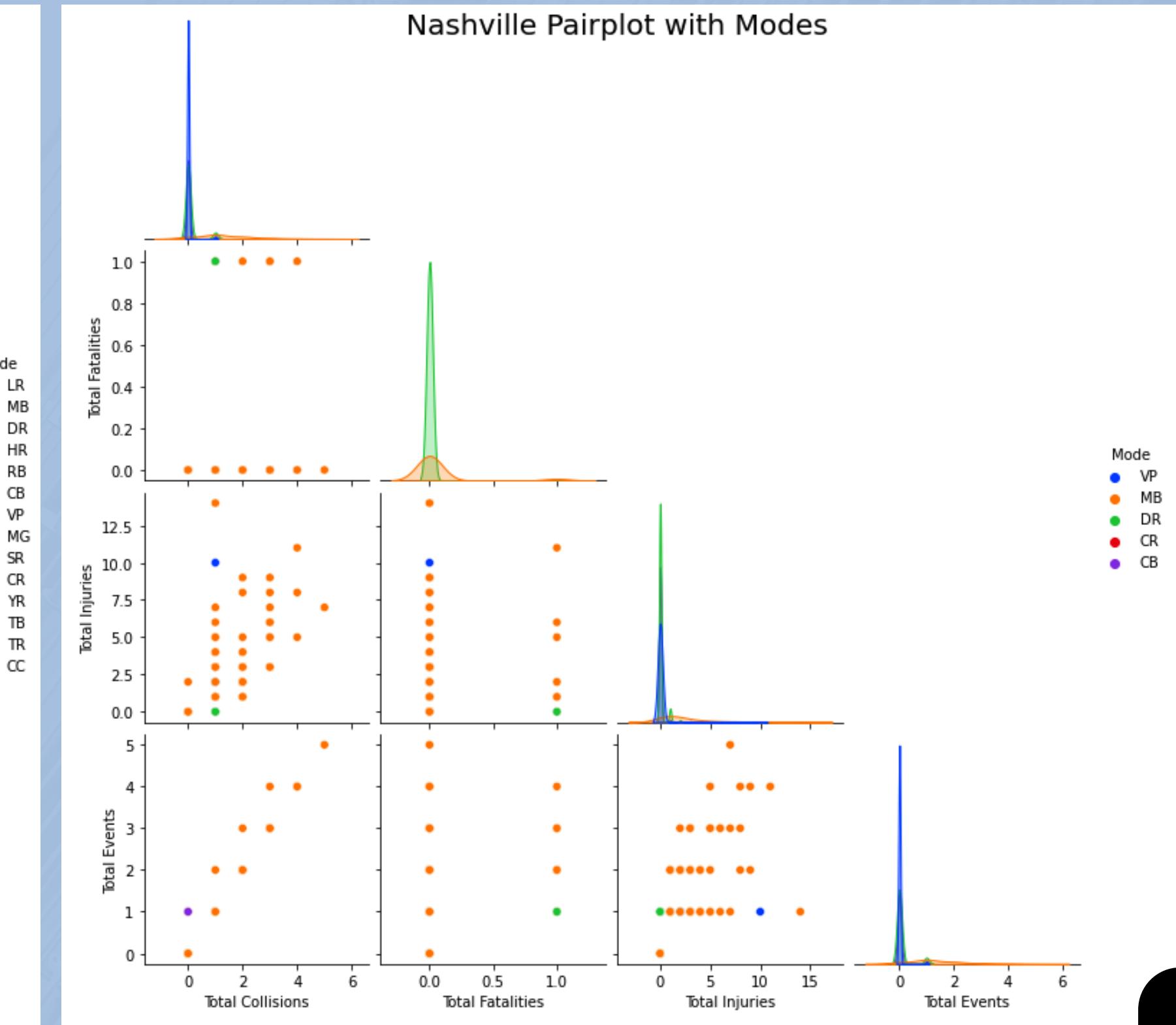
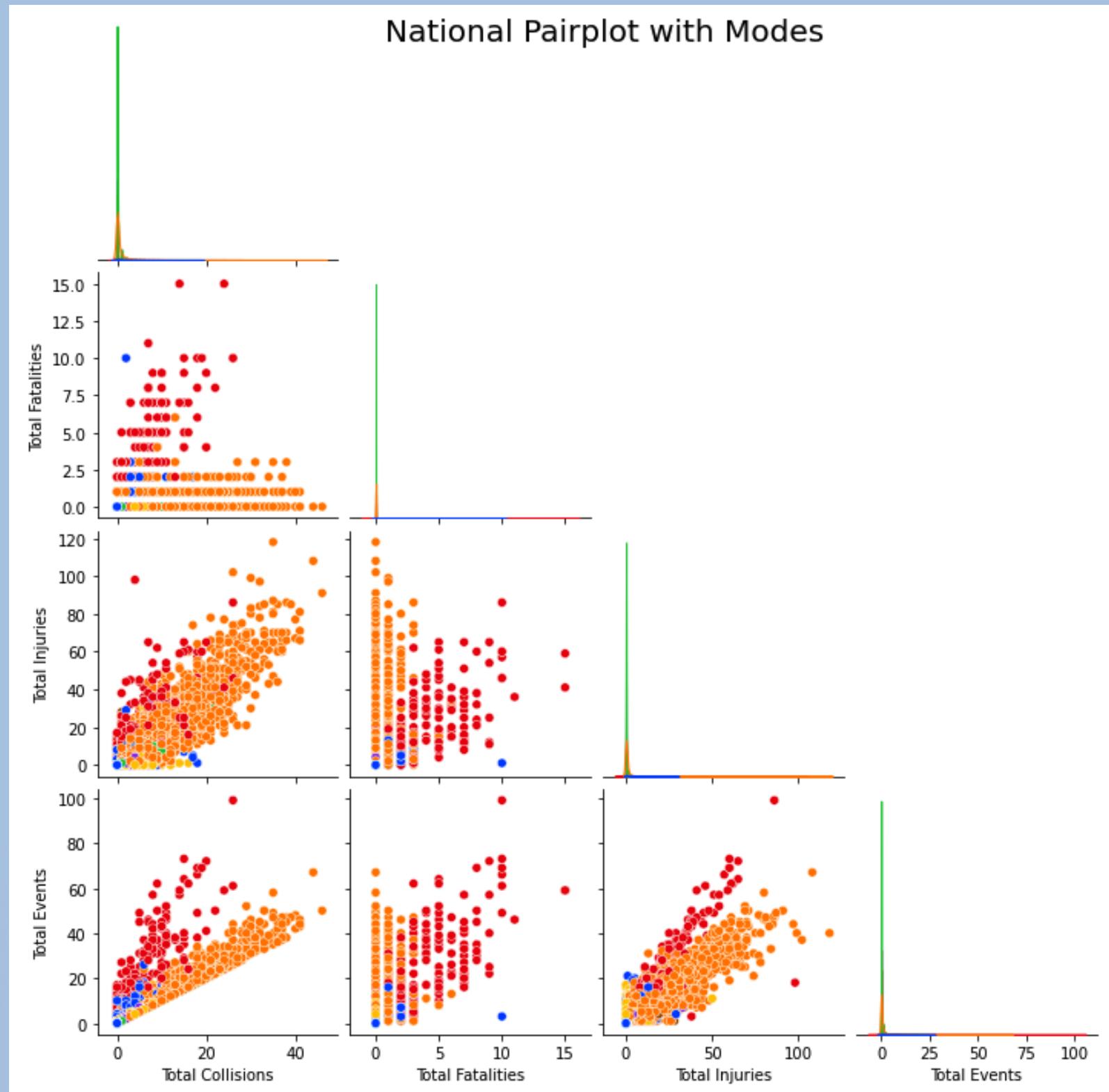


All contributing columns

# MANY SAFETY INCIDENTS FOR MB AND HR

## MB = MODE, BUS      HR = MODE, HEAVY RAIL

TAMARA HORNE  
2023 | March



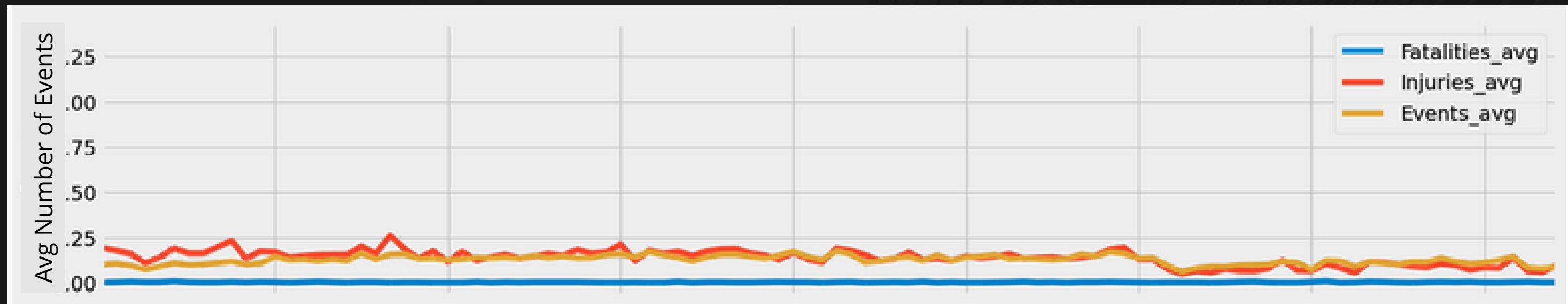
# DEFINING THE TARGET

<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

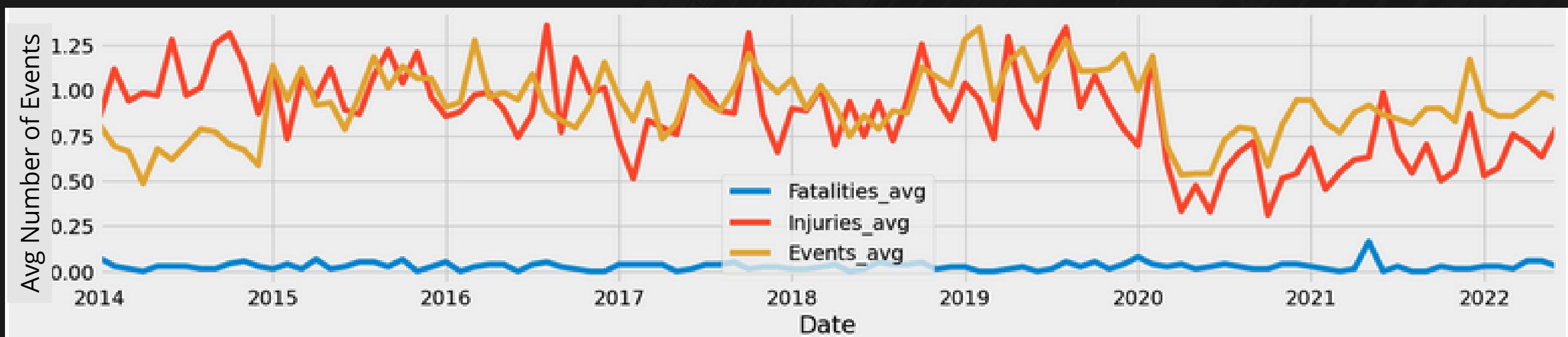
**(TOTAL FATALITIES + TOTAL INJURIES + TOTAL EVENTS) / VEHICLE REVENUE MILES**

# EXPLORATORY AGGREGATION

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



*Locations with Modes up to and Including Nashville's  
Plus Light Rail*

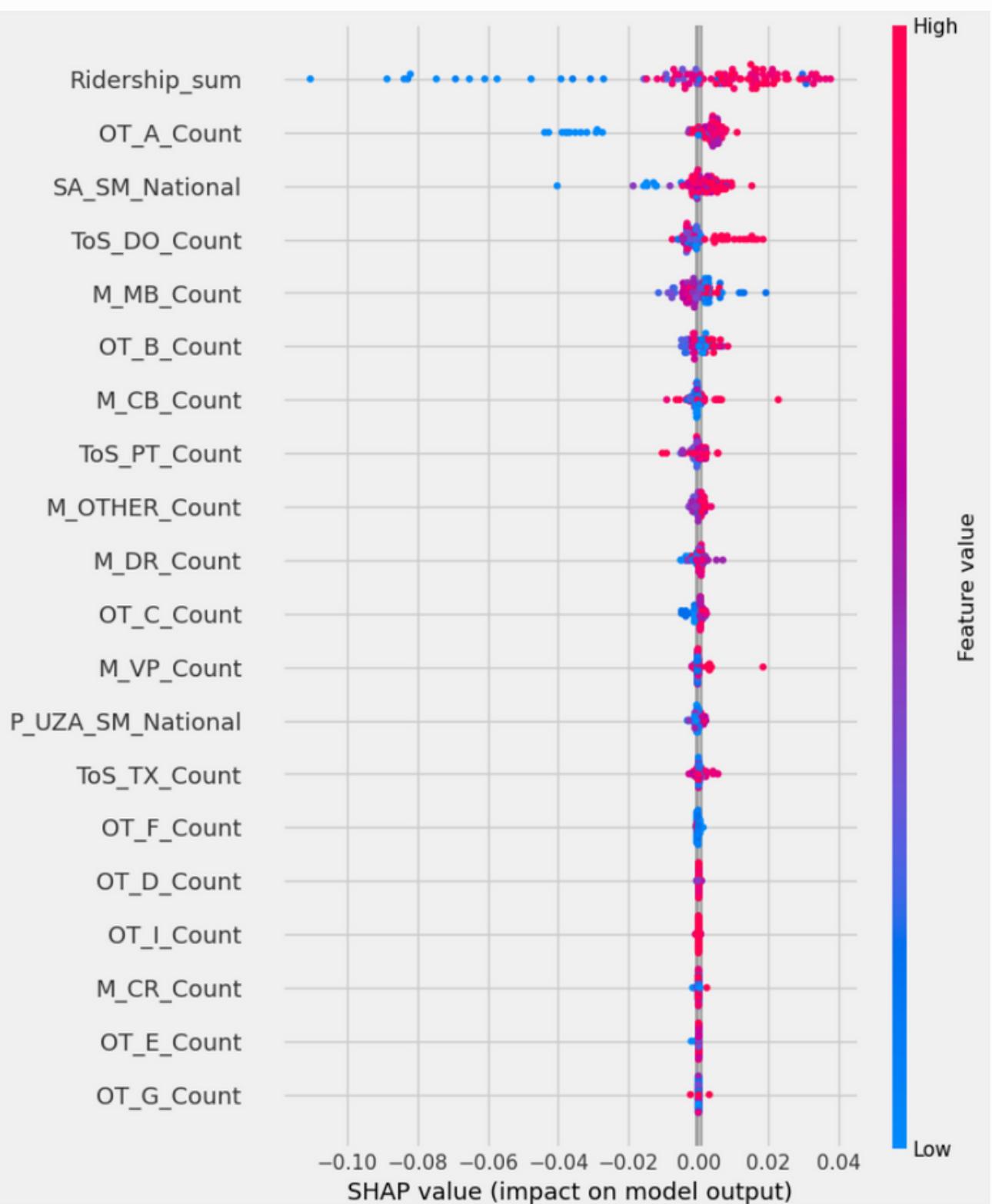


<i>~modeling_df~</i>	MAPE	RMSE	MAE	R Squared
Linear Regression	12.7%	0.088	0.058	0.645
Random Forest Regressor	7.1%	0.041	0.033	0.902
K Neighbors Regressor	7.5%	0.043	0.035	0.365
XGBoost Regressor	4.8%	0.027	0.022	0.999

# MODELING

*First Aggregation*   Date, nation

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



# FEATURE IMPORTANCE

<b>Ridership_sum</b>	Lower values decrease safety risk
<b>OT_A_Count</b>	'Independent Public Agency or Authority of Transit Service' organization type. Most frequently occurring organization type in the national data. Lower values decrease safety risk
<b>SA_SM_National</b>	Average service area square miles for the nation during each given month. Lower values decrease safety risk

~modeling_df2~	MAE	Mean y-test	Mean y-pred	RMSE	R Squared
Linear Regression	0.425	0.299	0.302	0.810	0.251
Random Forest Regressor	0.522	0.299	0.420	0.937	0.833
K Neighbors Regressor	0.395	0.299	0.270	0.938	0.280
XGBoost Regressor	0.454	0.299	0.340	0.899	0.548

\*MAPE not available due to division by zero\*

# MODELING

## *Second Aggregation Date, location*

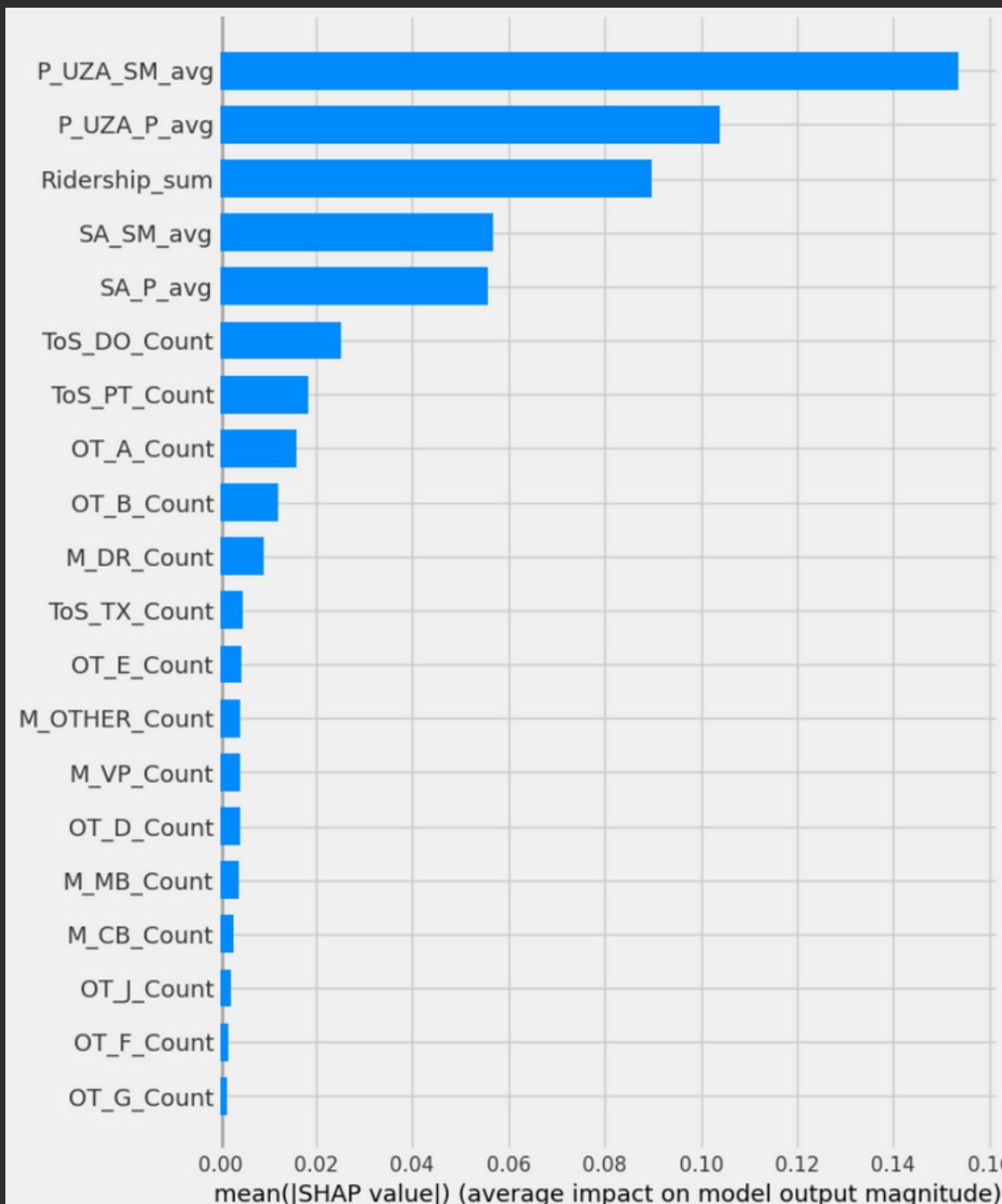
What features influence safety risk when data is aggregated to both date and location (Primary UZA Name)?



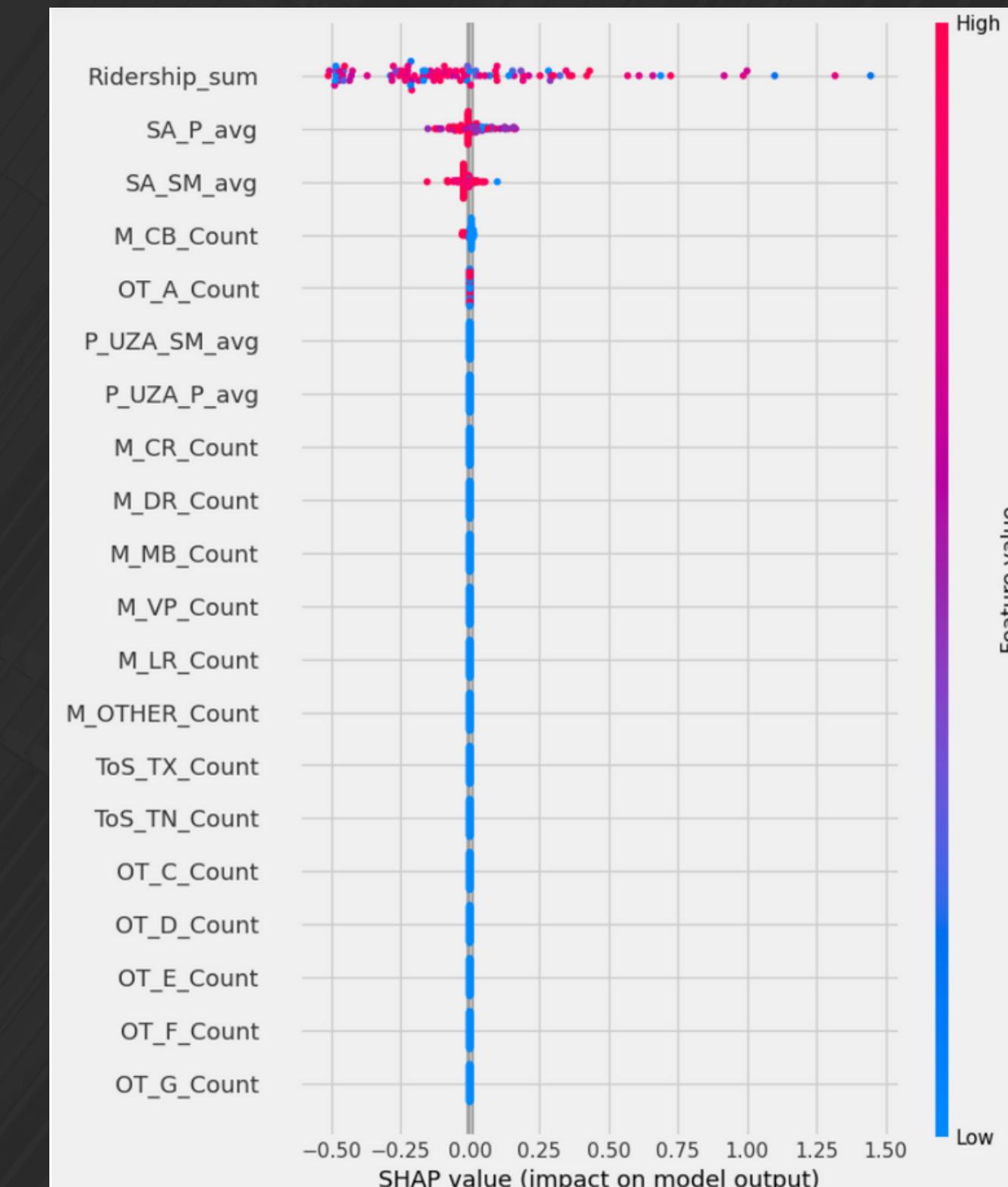
**SHIFTED RESIDUALS MAKE THE RESULTS LESS TRUSTWORTHY**

# FEATURE IMPORTANCE COMPARISON

## National



## Nashville



<i>~modeling_df3~</i>	MAE	Mean y-test	Mean y-pred	RMSE	R Squared
Linear Regression	1.403	1.149	1.594	1.893	0.210
Random Forest Regressor	0.680	1.149	0.774	1.503	0.883
K Neighbors Regressor	0.591	1.149	0.767	1.490	0.595
XGBoost Regressor	0.744	1.149	0.830	1.608	0.945

\*MAPE not available due to division by zero\*

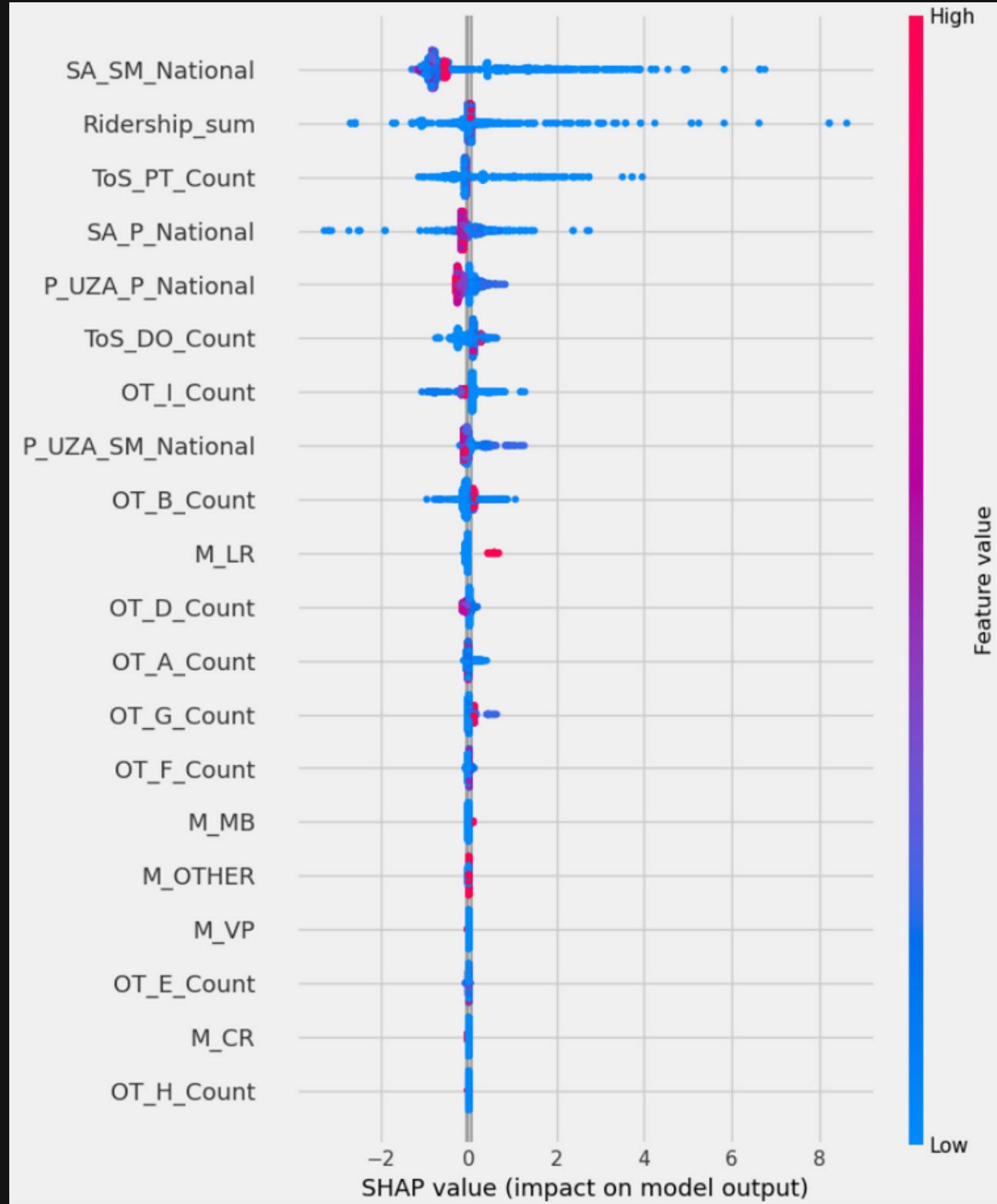
# MODELING

## *Third Aggregation Date, mode*

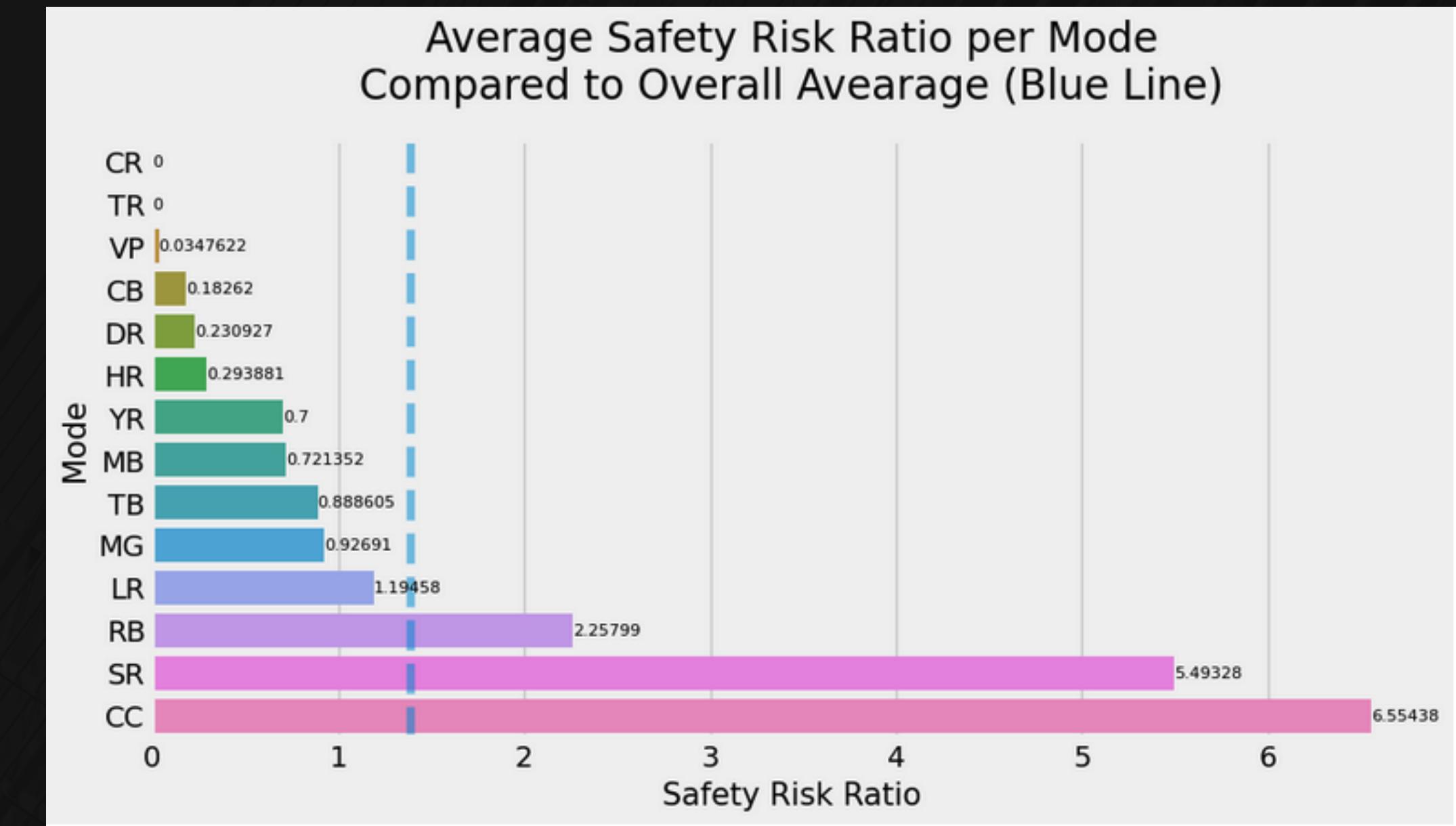
What features influence safety risk for the individual modes?

Which modes have the highest and lowest safety risk?

# FEATURE IMPORTANCE



# MODAL SAFETY RISK



# CONCLUSIONS

## Best Models

- XGBoost Regressor
- K Neighbors
- K Neighbors

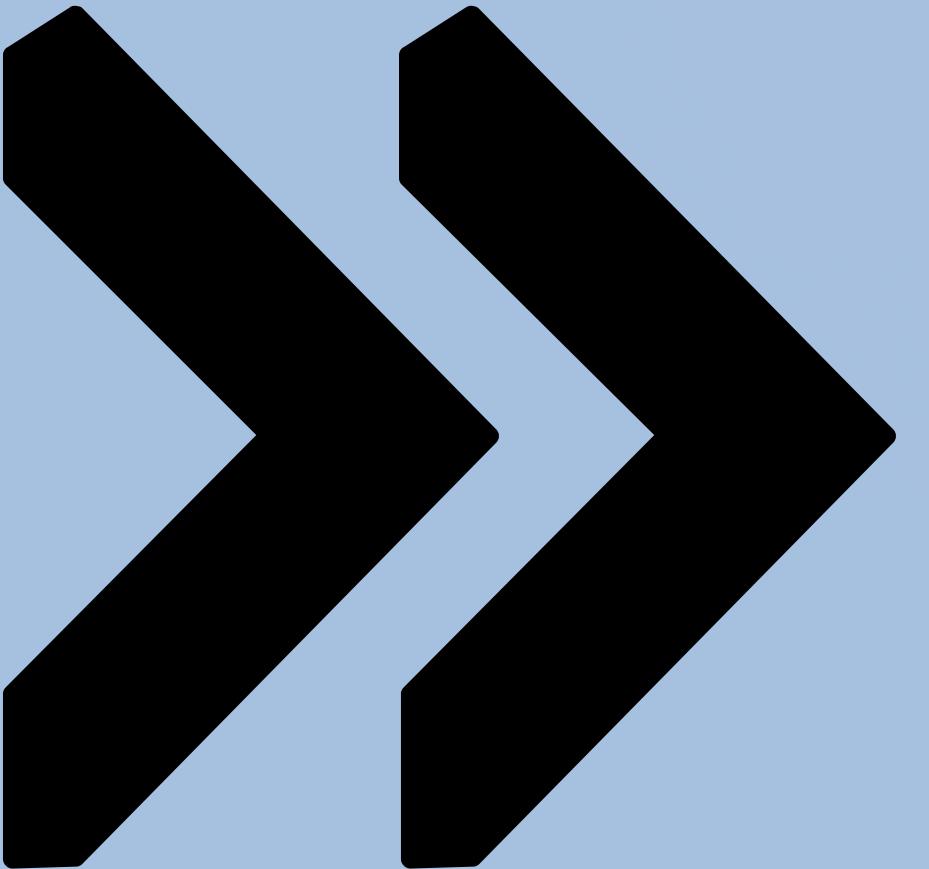
## Aggregations

- 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
- Use K Neighbors to further explore modeling\_df2 and modeling\_df3.

## *Big*

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

# RECOMMENDATIONS

## For the Client

YES *Mode* MATTERS

TAMARA HORNE  
2023 | March

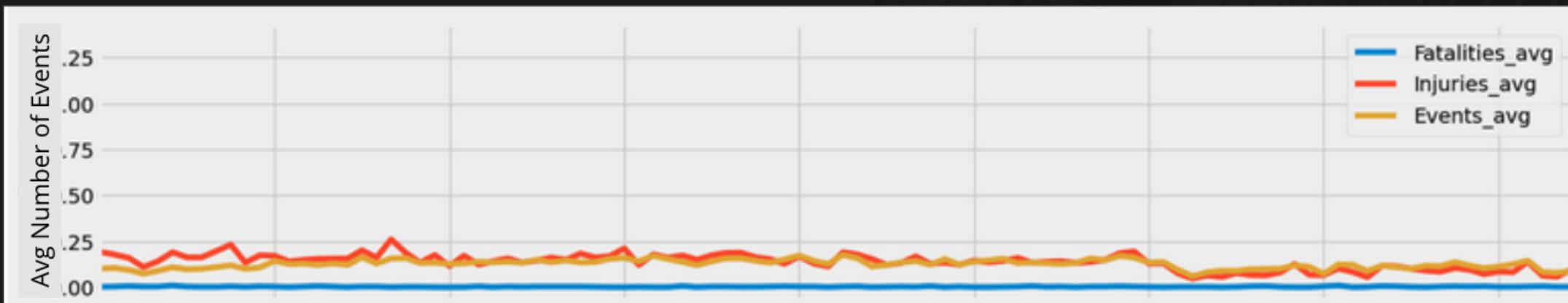
\*Correlation Does Not Mean Causation



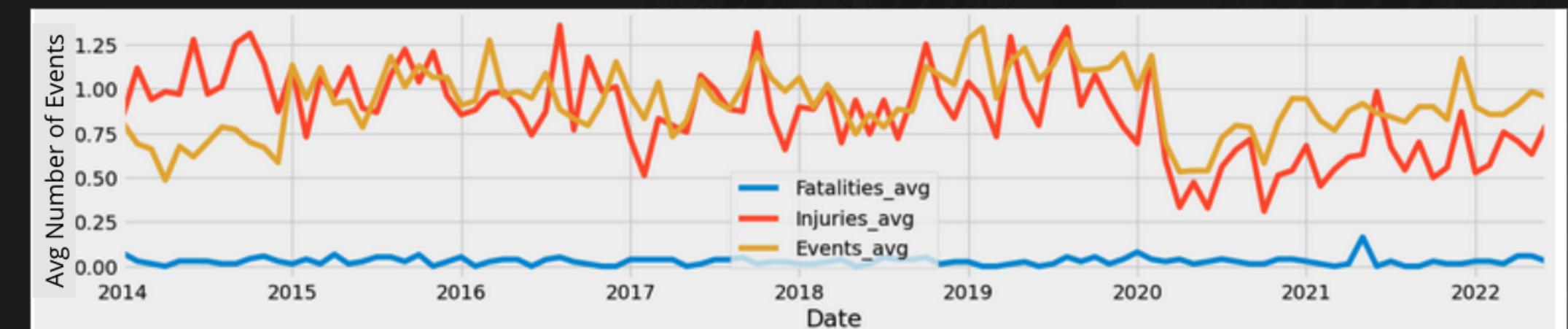
## DON'T ADD LIGHT RAIL

Adding light rail would increase safety risk

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



*Locations with Modes up to and Including Nashville's Plus Light Rail*



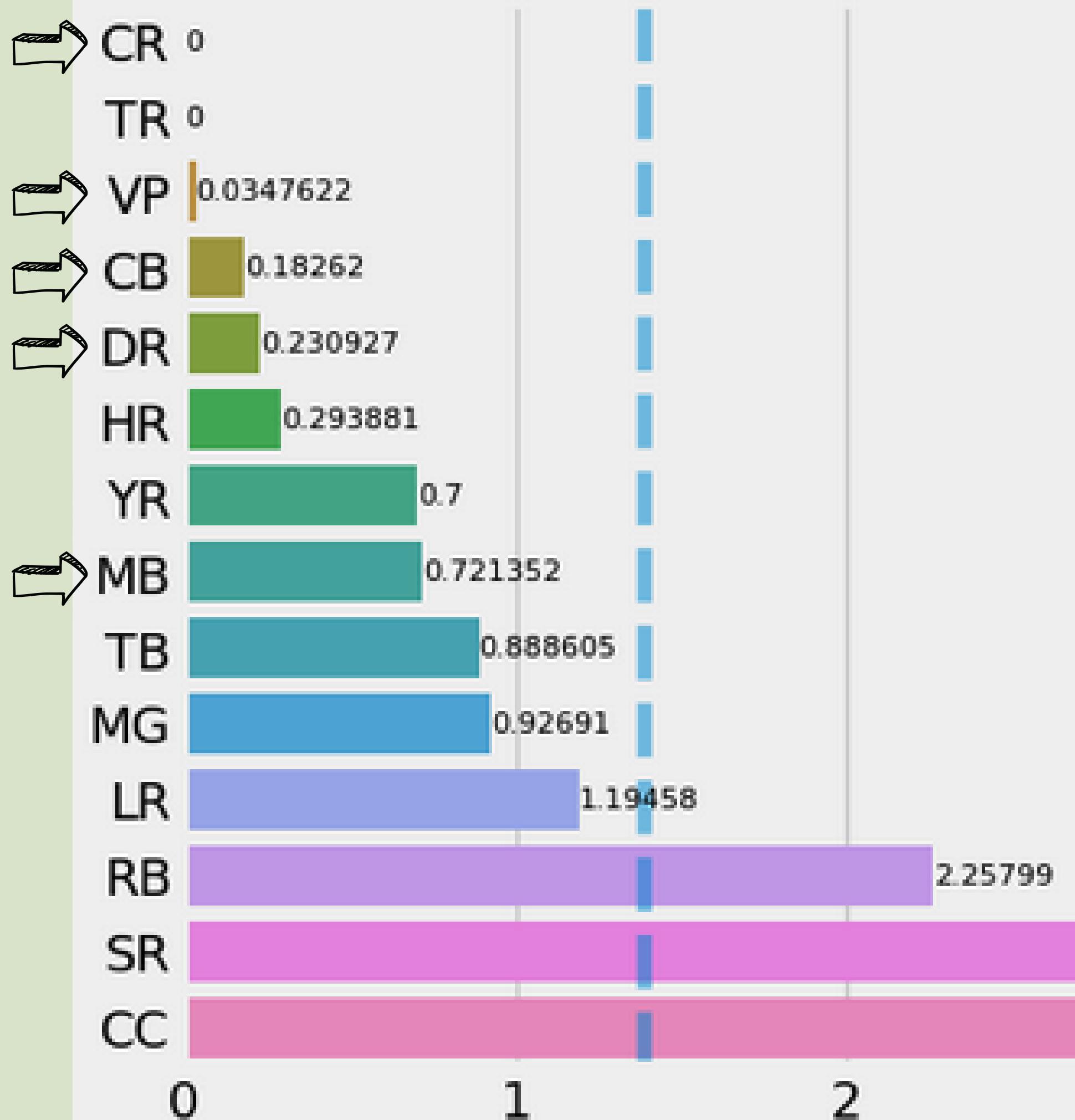


## ADD SERVICE TO EXISTING MODES

- ✓ Safety Risk is lowest for CR
- ✗ Safety Risk is highest for MB

TAMARA HORNE

2023 | March





## ADD SERVICE TO EXISTING MODES



Added vehicle revenue miles will differ according to mode and this will affect the target

# Nashville's Modes

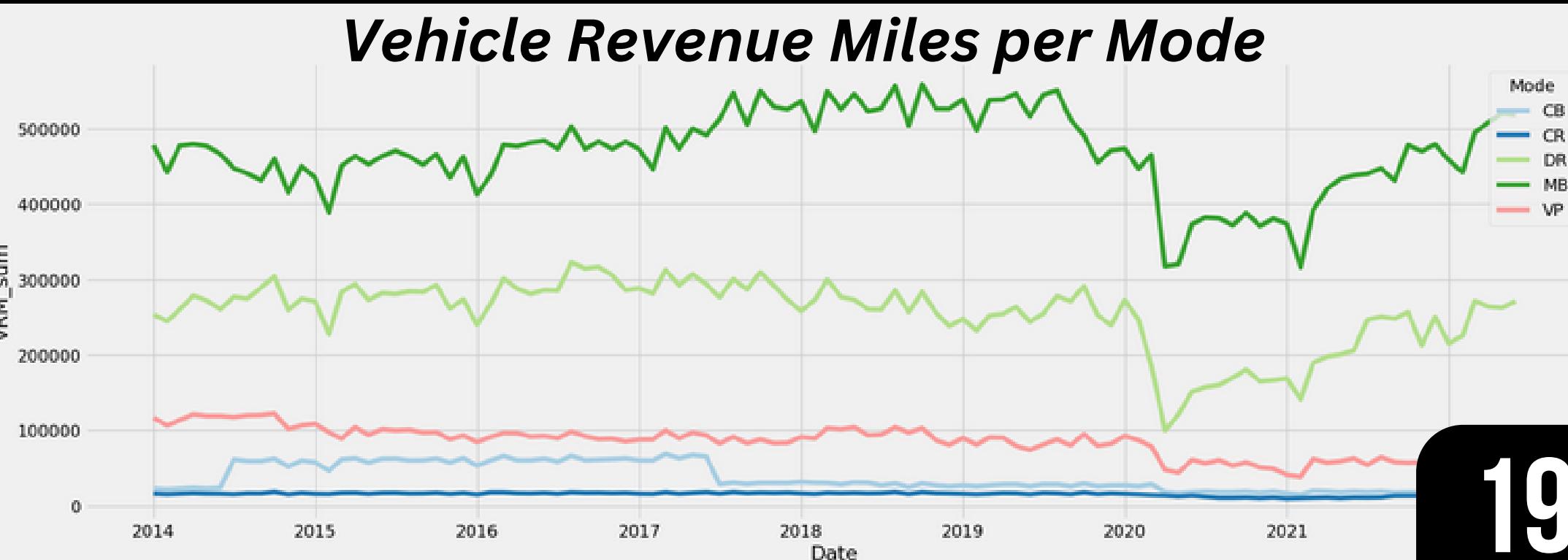
CB: Commuter Bus

CR: Commuter Rail

DR: Demand Response (taxis, etc)

MB: Bus

VP: Vanpool



# MORE INFORMATION

---

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



Image Creator: jmsilva | Credit: Getty Images

TAMARA HORNE  
2023 | March