

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

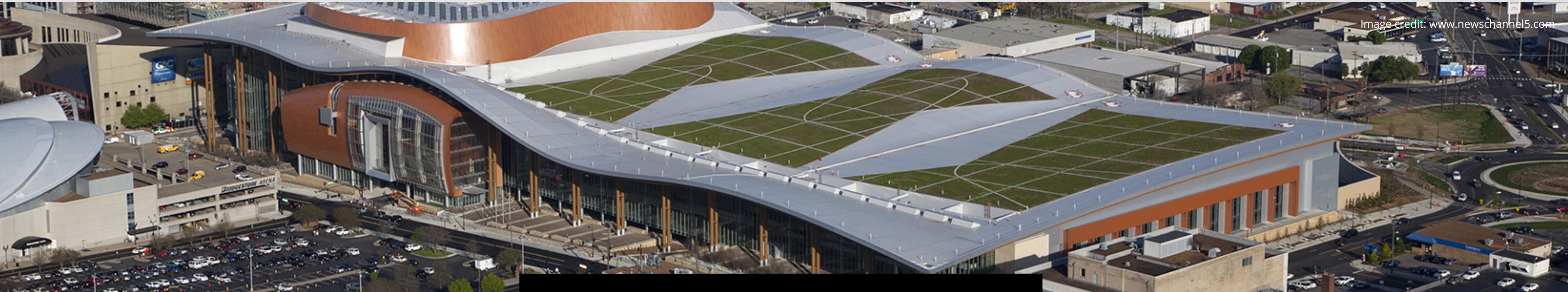


Image credit: www.newschannel5.com

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)

TAMARA HORNE

2023 | March

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

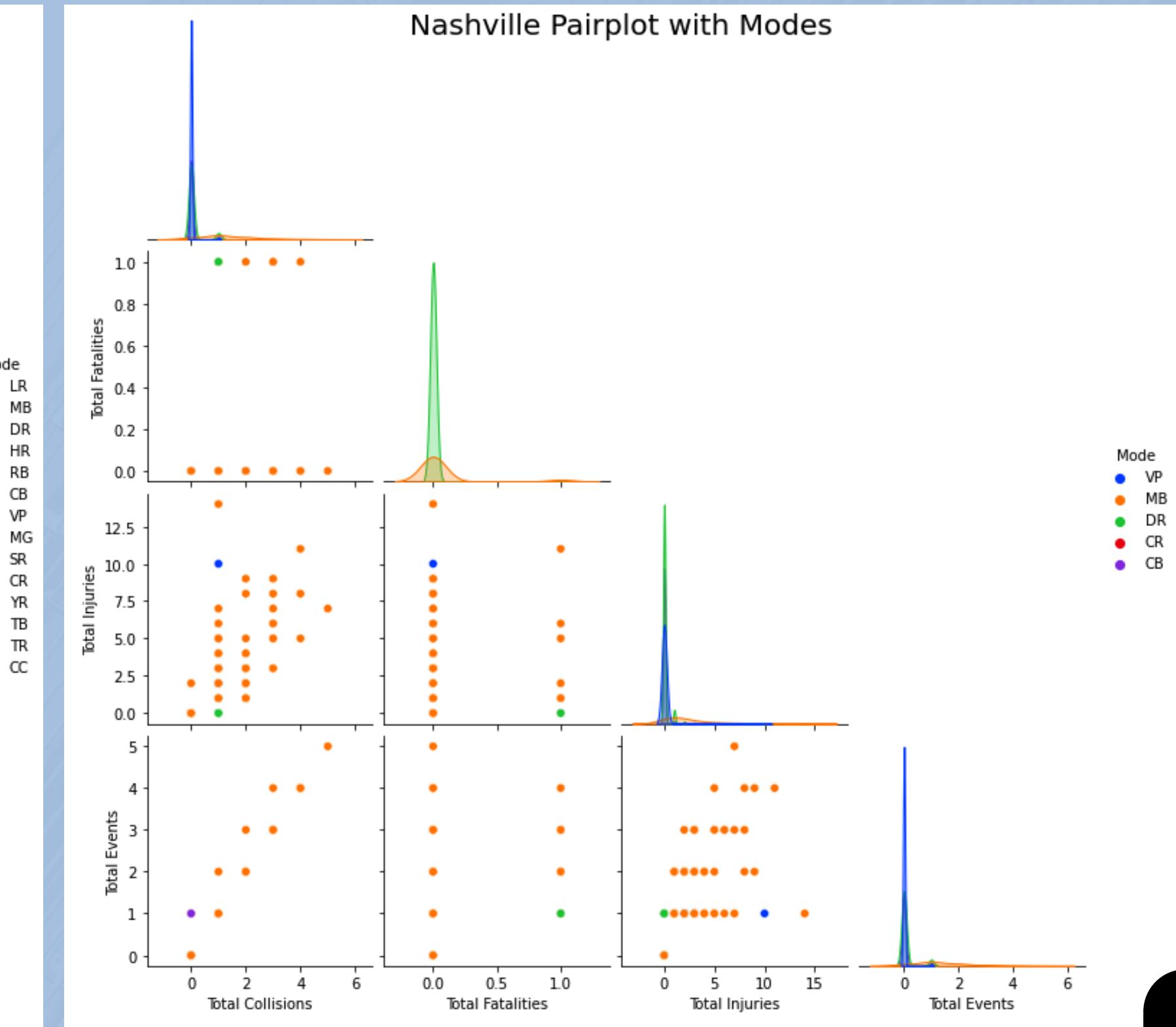
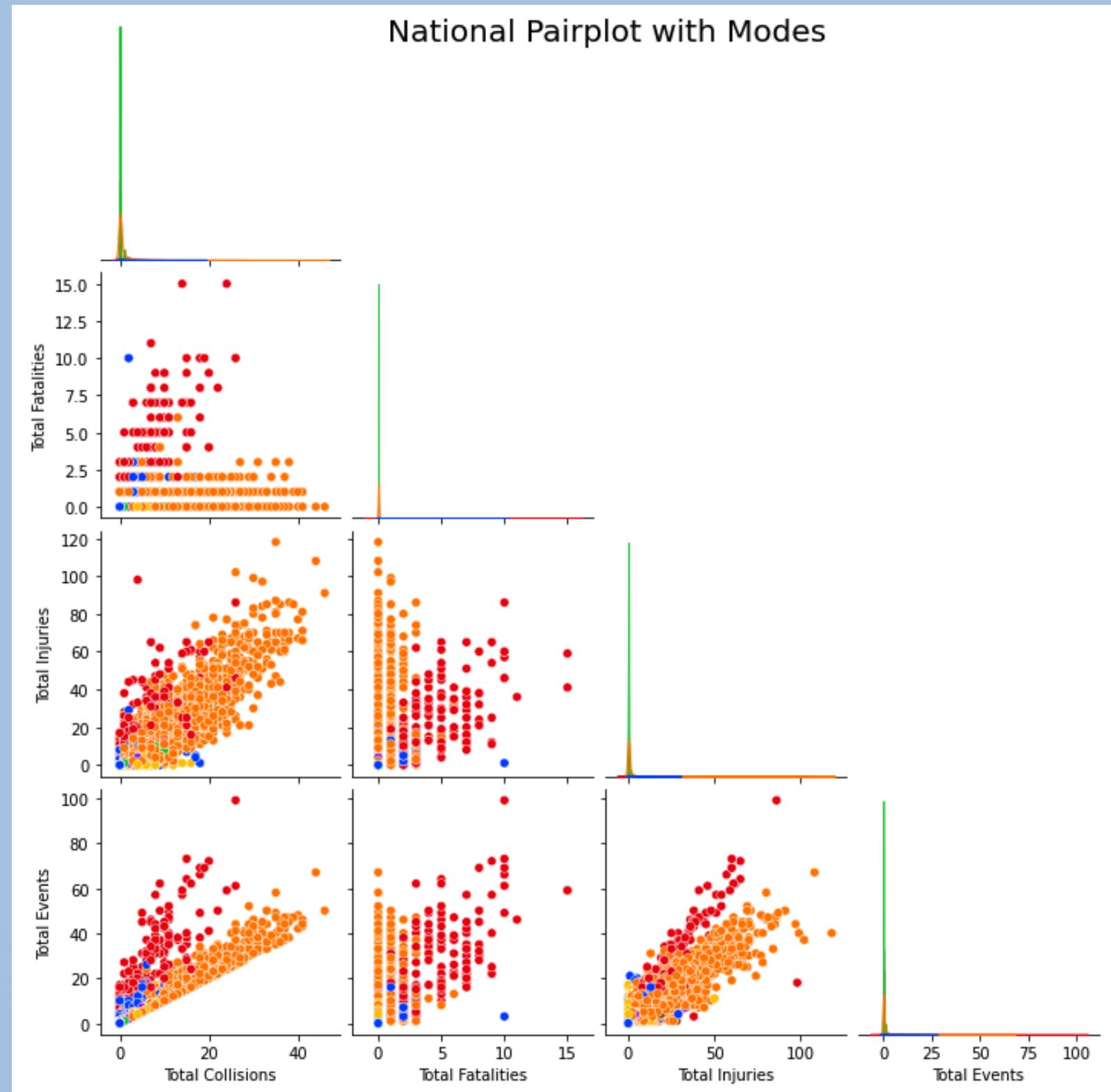


All contributing columns

MANY SAFETY INCIDENTS FOR MB AND HR

MB = MODE, BUS HR = MODE, HEAVY RAIL

TAMARA HORNE
2023 | March



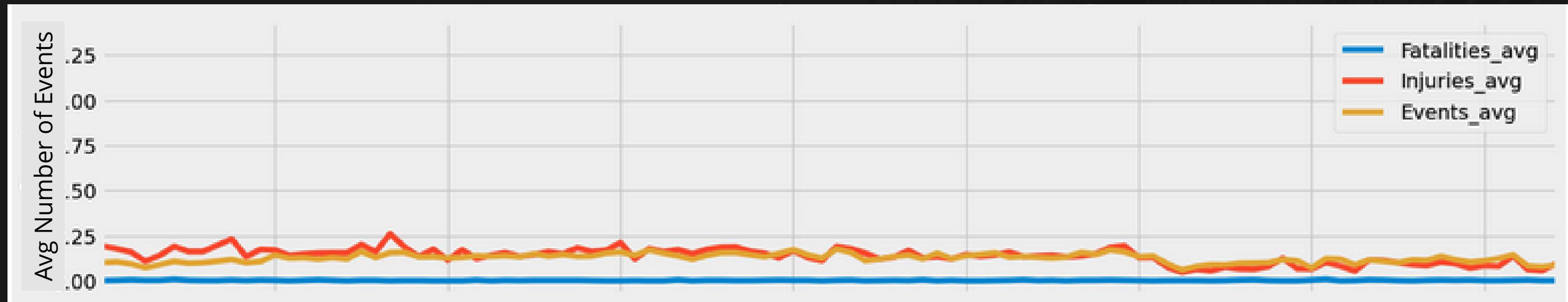
DEFINING THE TARGET

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

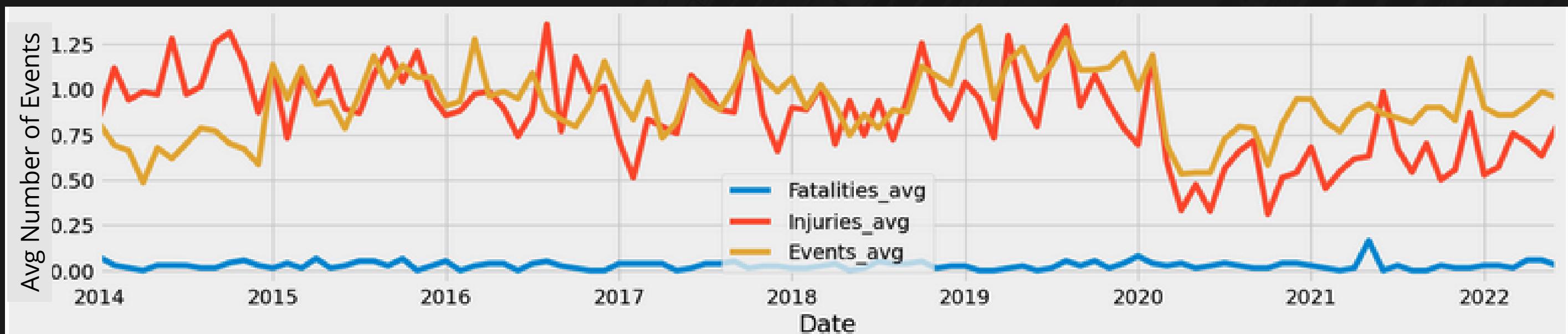
(TOTAL FATALITIES + TOTAL INJURIES + TOTAL EVENTS) / VEHICLE REVENUE MILES

EXPLORATORY AGGREGATION

Locations with Modes up to and Including Nashville's



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

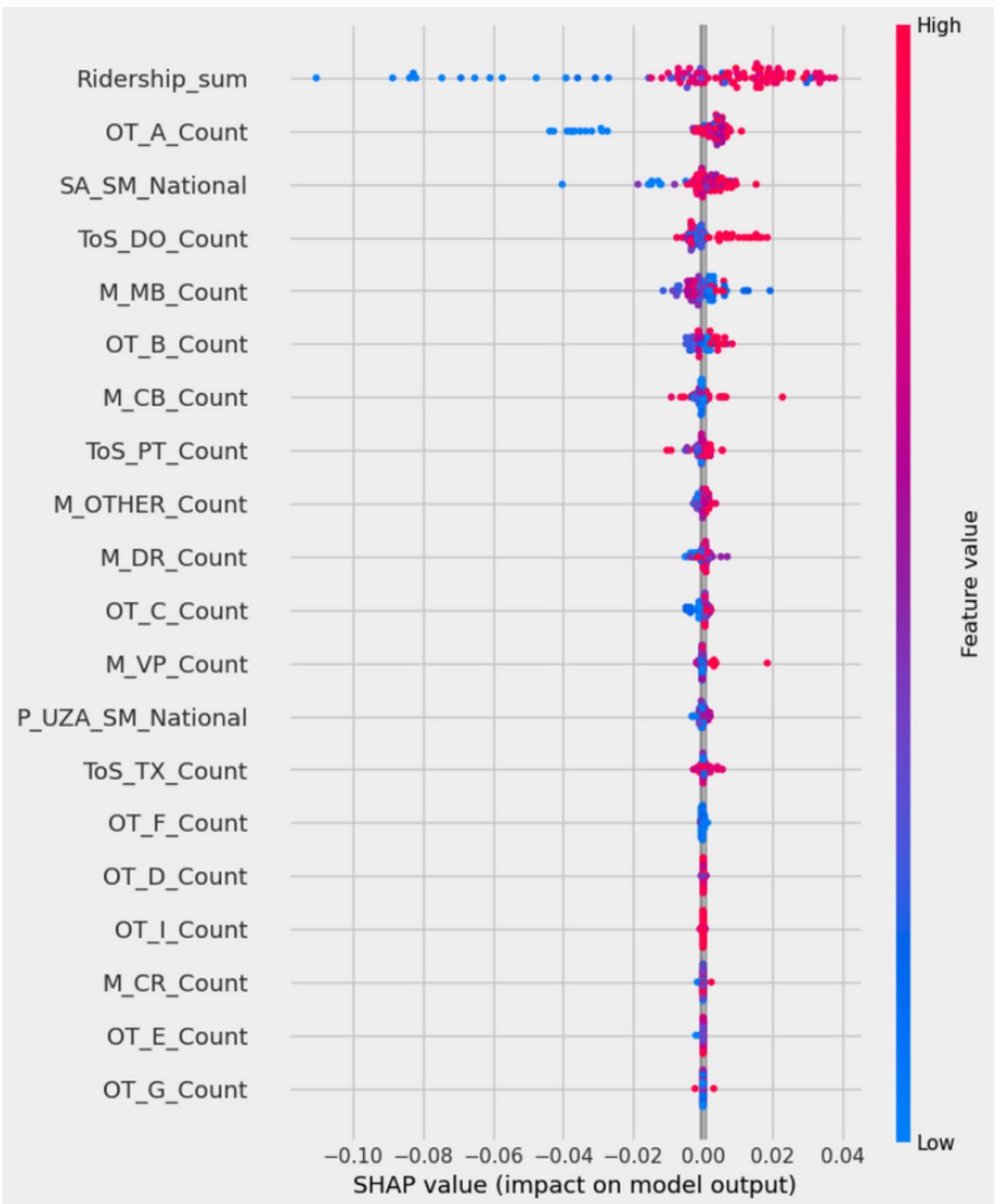


<i>~modeling_df~</i>	R Squared	MAPE	MAE	MSE
Linear Regression	0.645	12.7%	0.058	0.007
Random Forest Regressor	0.902	7.1%	0.033	0.002
K Neighbors Regressor	0.365	7.5%	0.035	0.002
XGBoost Regressor	0.999	4.8%	0.022	0.001

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

Most frequently occurring organization type in the national data

Lower values decrease safety risk

SA_SM_National

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

Lower values decrease safety risk

<i>~modeling_df2~</i>	R Squared	Mean y-test	Mean y-pred	MAE	MSE
Linear Regression	0.251	0.299	0.302	0.425	0.656
Random Forest Regressor	0.833	0.299	0.420	0.522	0.871
K Neighbors Regressor	0.280	0.299	0.270	0.395	0.880
XGBoost Regressor	0.548	0.299	0.340	0.454	0.808

MAPE not available due to division by zero

MODELING

Second Aggregation Date, location

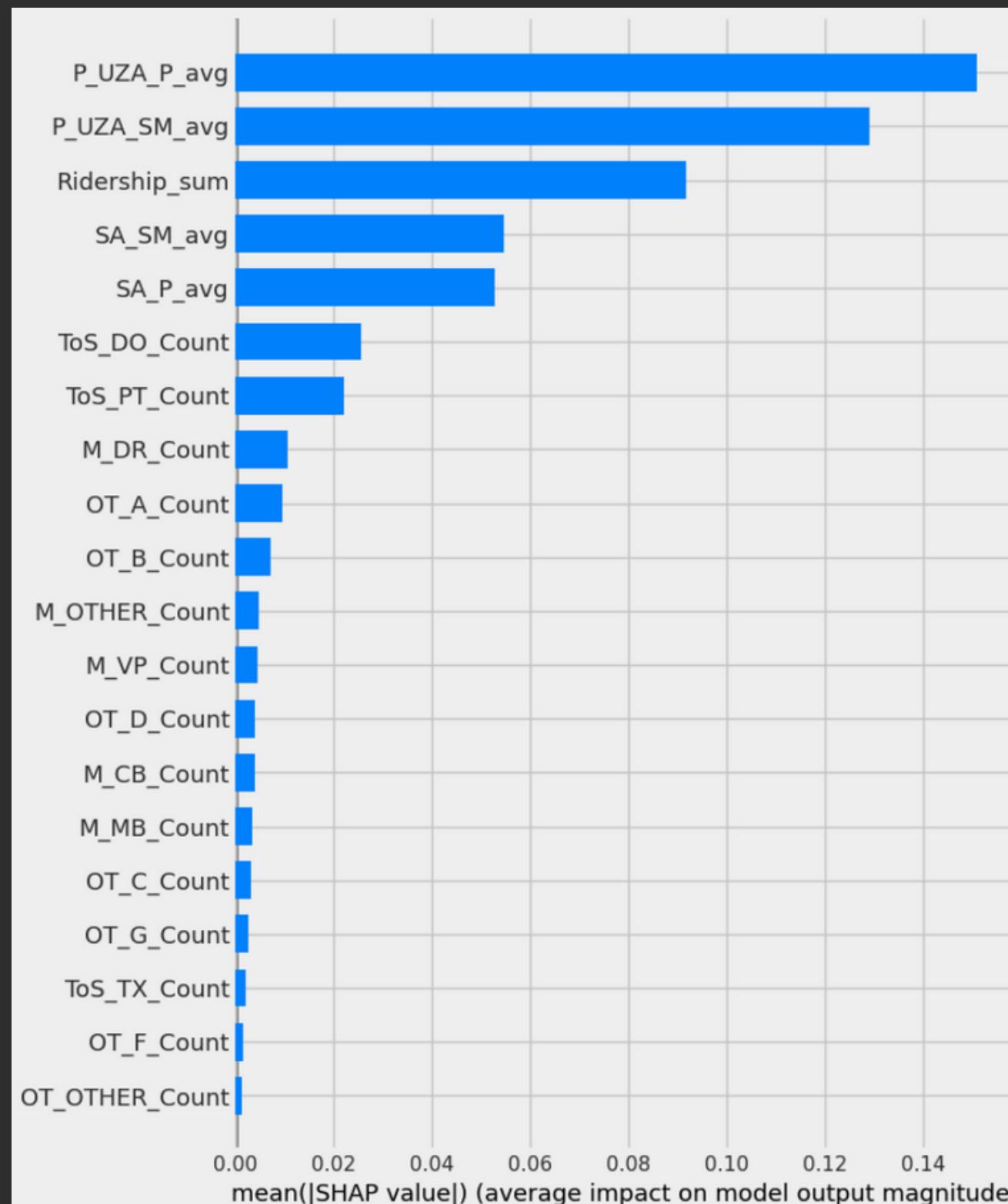
What features influence safety risk when data is aggregated to both date and location (Primary UZA Name)?



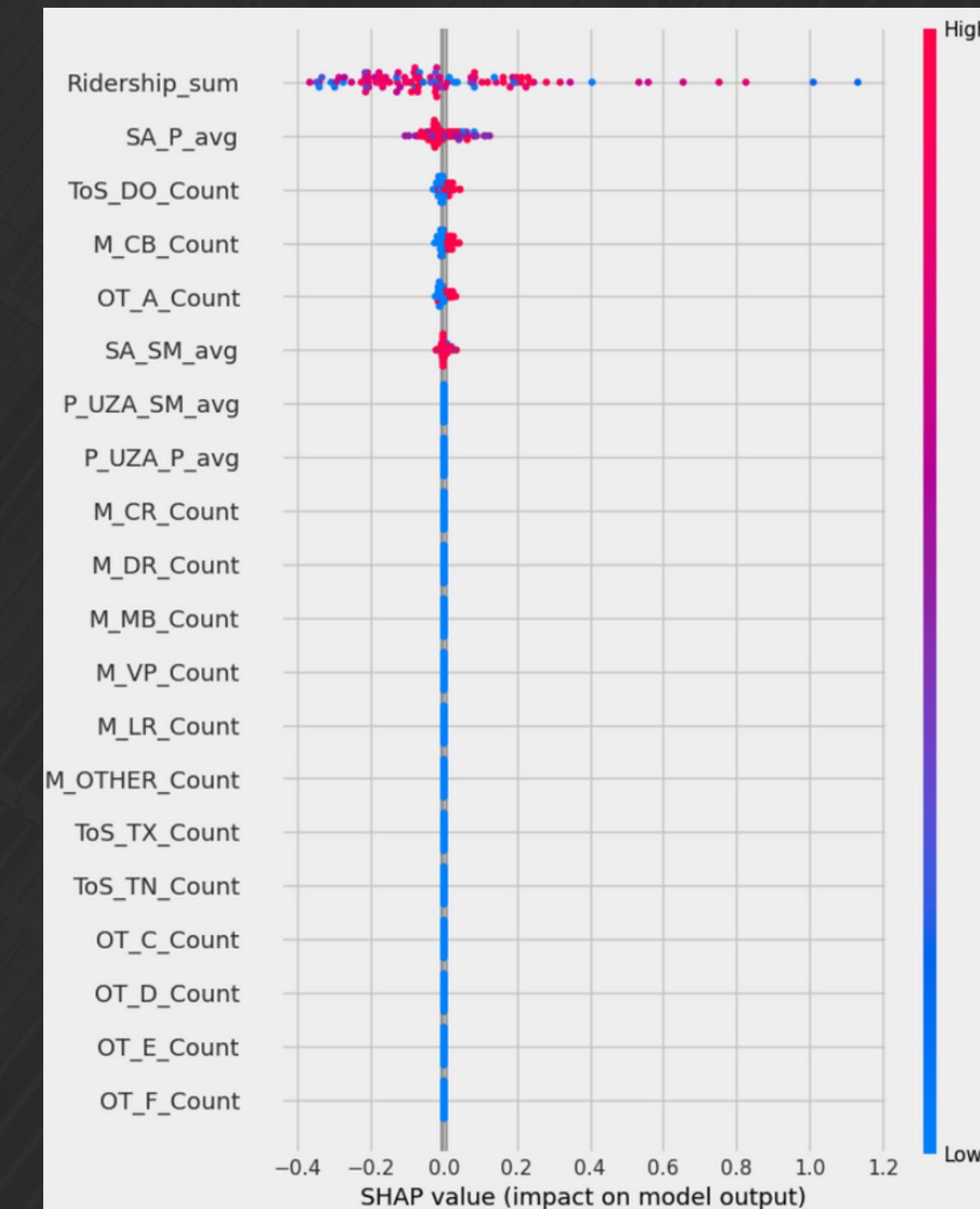
SHIFTED RESIDUALS MAKE THE RESULTS LESS TRUSTWORTHY

FEATURE IMPORTANCE COMPARISON

National



Nashville



<i>~modeling_df3~</i>	R Squared	Mean y-test	Mean y-pred	MAE	MSE
Linear Regression	0.210	1.149	1.594	1.403	3.587
Random Forest Regressor	0.883	1.149	0.774	0.680	2.335
K Neighbors Regressor	0.595	1.149	0.767	0.591	2.221
XGBoost Regressor	0.945	1.149	0.830	0.744	2.587

MAPE not available due to division by zero

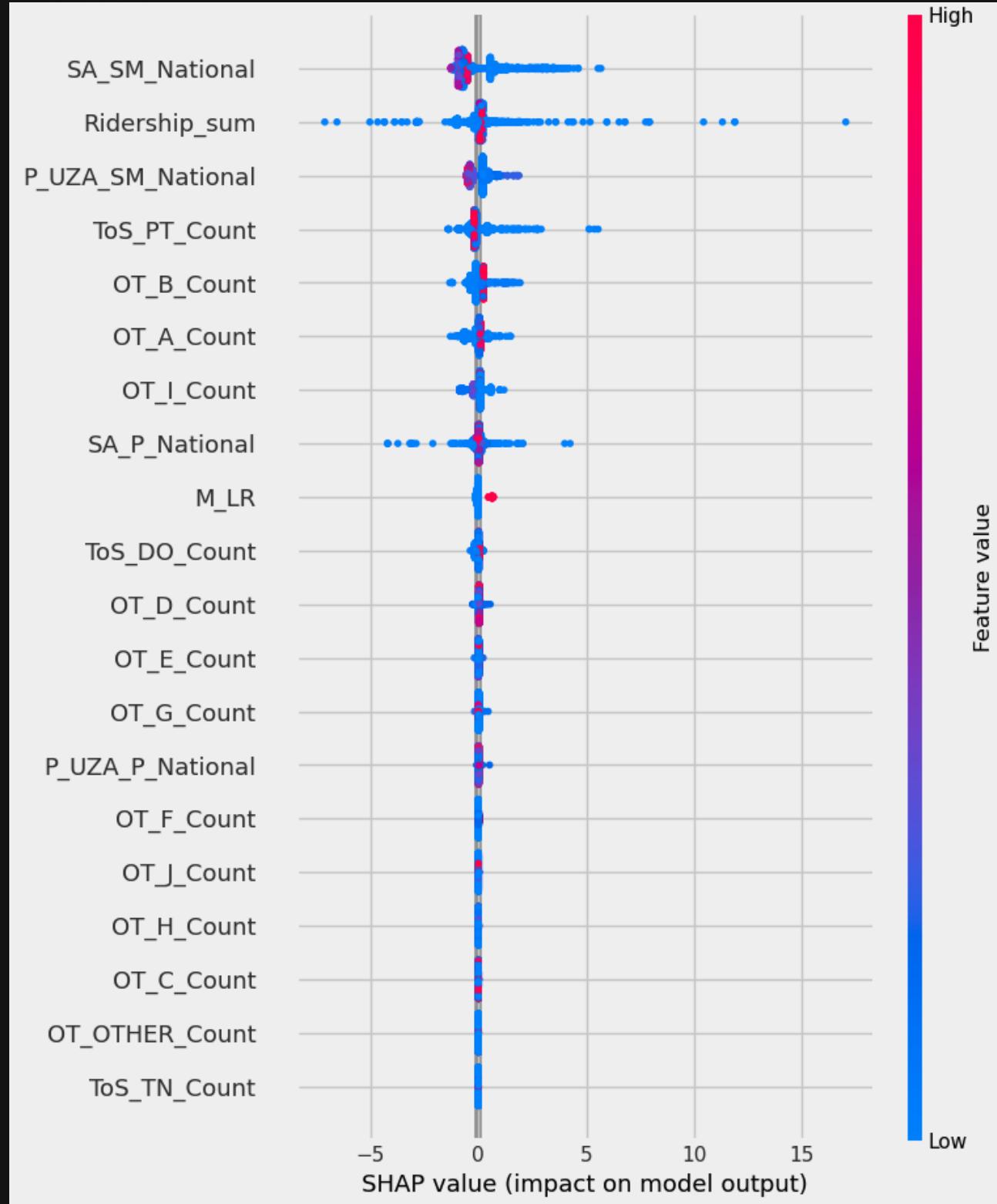
MODELING

Third Aggregation Date, mode

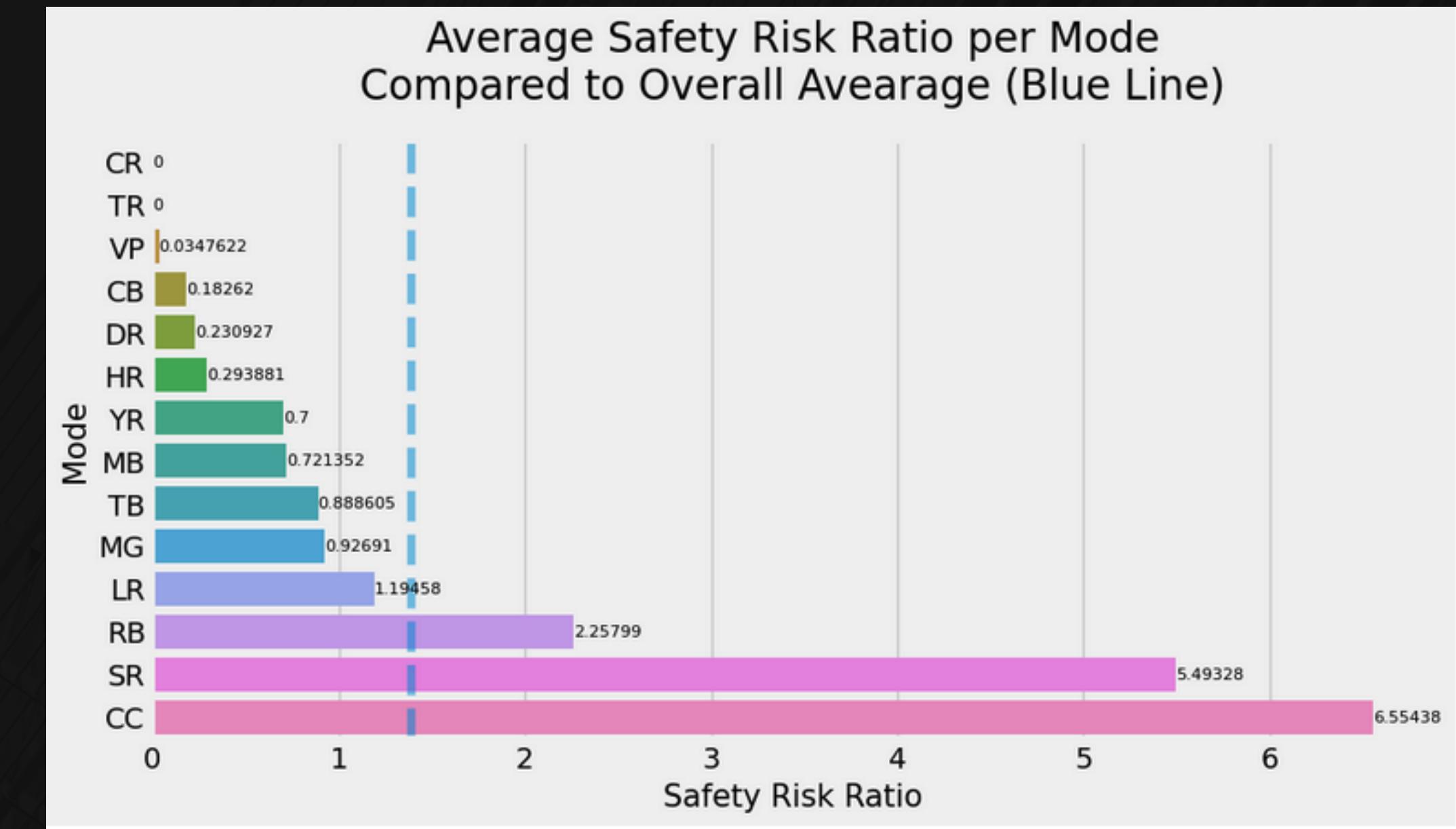
What features influence safety risk for the individual modes?

Which modes have the highest and lowest safety risk?

FEATURE IMPORTANCE



MODAL SAFETY RISK



CONCLUSIONS

Best Models

- XGBoost Regressor
- Random Forest Regressor

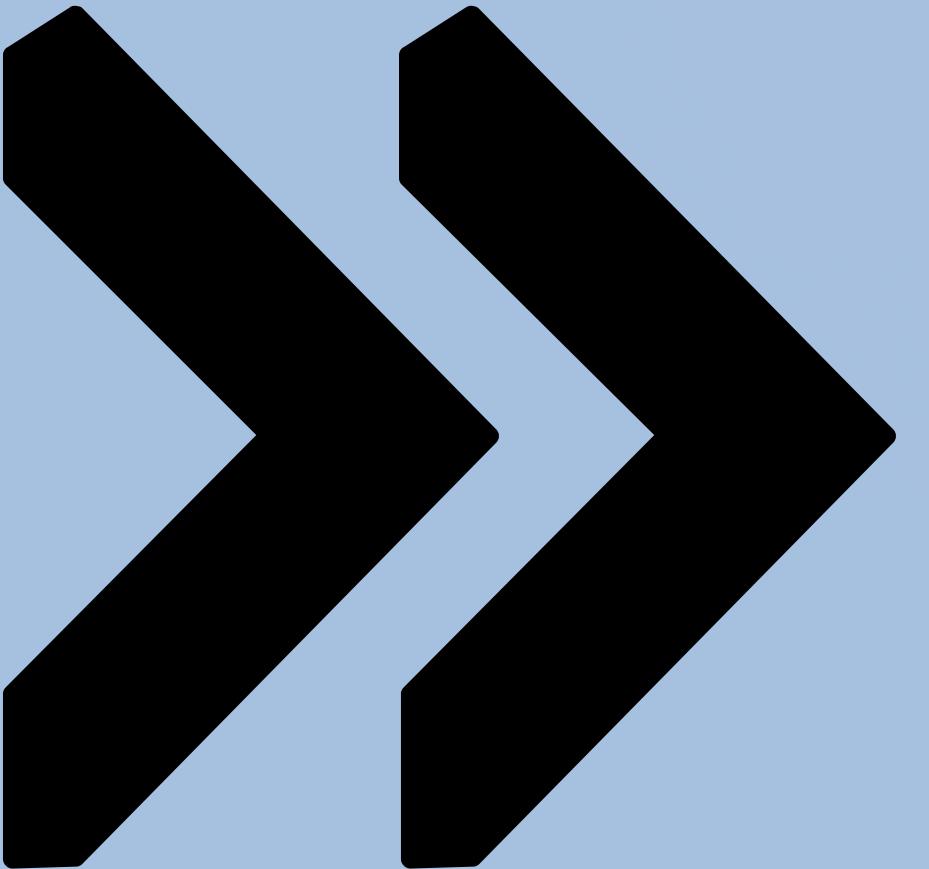
Aggregations

- 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
- 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

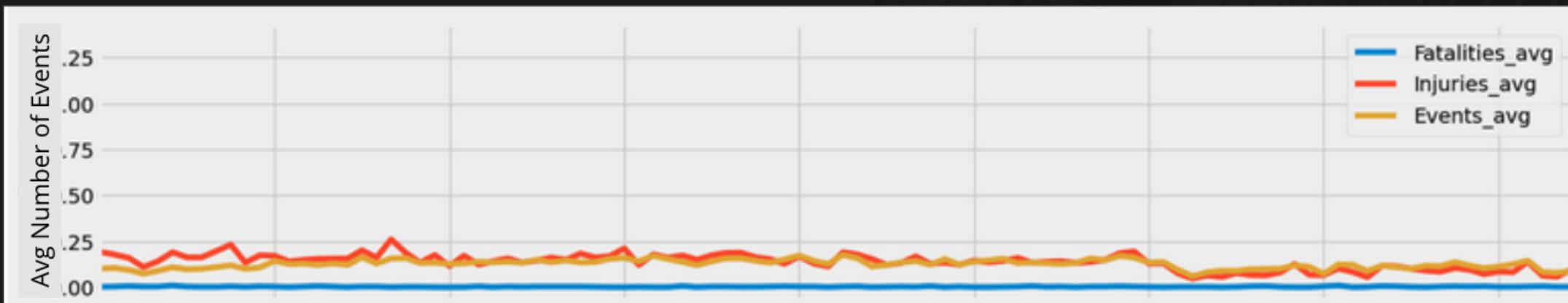
*Correlation Does Not Mean Causation



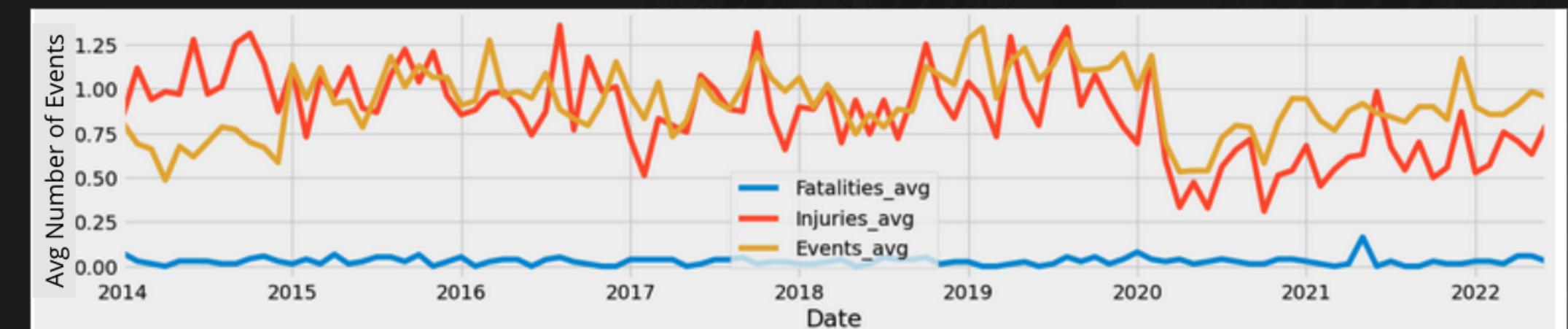
DON'T ADD LIGHT RAIL

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



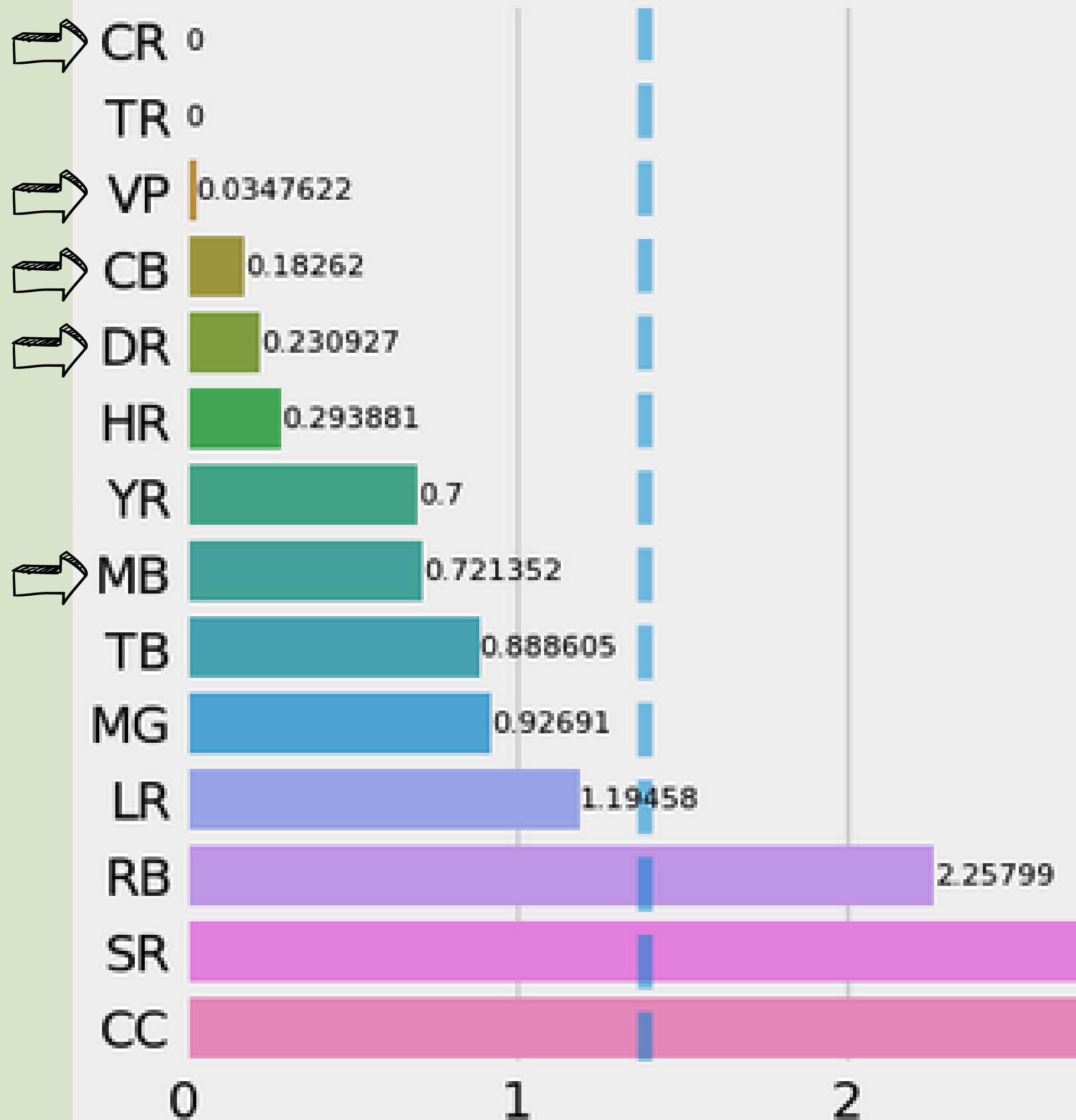


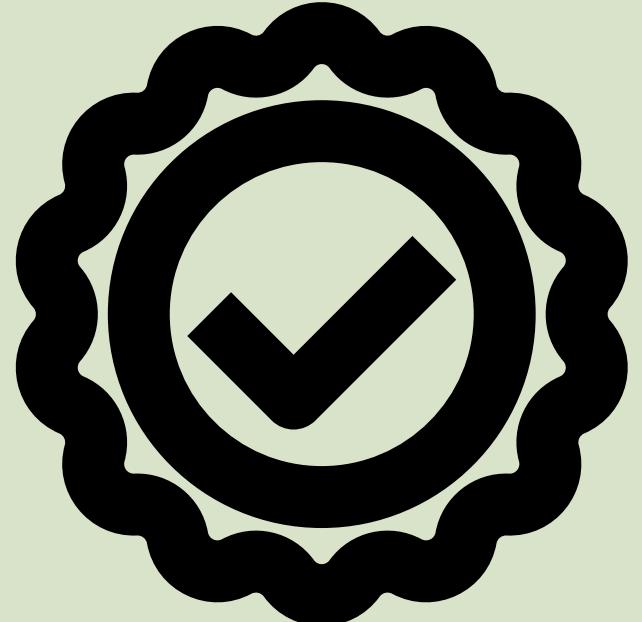
ADD SERVICE TO EXISTING MODES

- ✓ 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

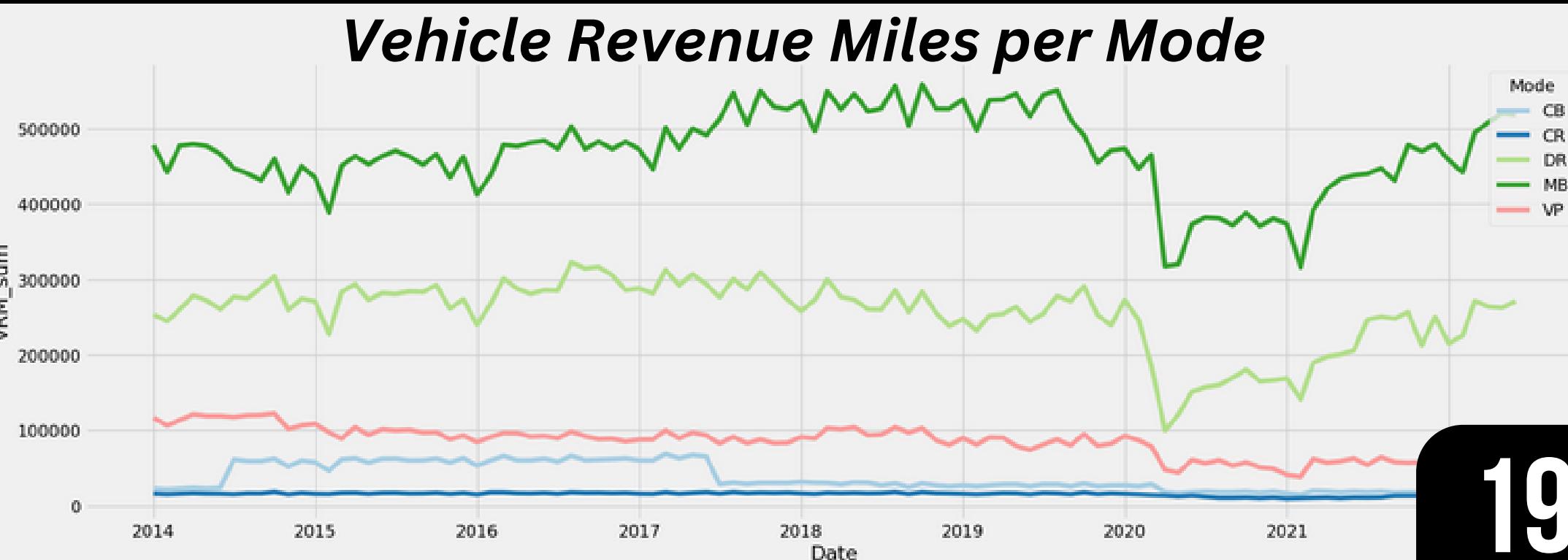
CB: Commuter Bus

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