# Urban Mobility Trends Unveiled: A Comprehensive Analysis of Lyft Bikes in San Francisco

Meta Team Members:

Divya Dodla (UA03330)

Naren Kandregula (PW71937)

Sachin Dudam (GG42052)

# Project Description

**GOAL :** To analyze lyft Bay Wheels System Data of San Francisco Bay Area from the years 2021-2023 by performing Big Data operations, feature engineering, generating visualizations (EDA), and building prediction model based on the principles of Big Data to predict the potential bike availability at specific stations.

**DATA :**

- Size of dataset: Approximately 1.24 GB.

- Number of Records: 6,625,912 individual records.

The large dataset size necessitates leveraging distributed computing capabilities, and the complexity of the analysis demands scalable solutions to extract actionable insights from the wealth of information within the Lyft Bay Wheels System Data

# Motivation

- The increasing popularity of bike-sharing systems contributes significantly to sustainable urban transportation.

- Analyzing the Bay Wheels System Data allows us to gain insights into user behaviors, station utilization, and broader trends, aiding urban planners and policymakers in making informed decisions to enhance the efficiency and accessibility of bike-sharing services.

# Implementation of Architecture

**Input Data Source**  **Platform**  **Storage**  **Processing**

Visualizations

Prediction Models

Spark MLLib

# Feature Engineering & Exploratory Data Analysis (EDA)

Missing values in categorical columns (Station Names and IDs) are filled with "Unknown" and rows with missing values in geographic coordinates (Latitude and Longitude) are removed.

**Potential Canceled Rides Analysis**

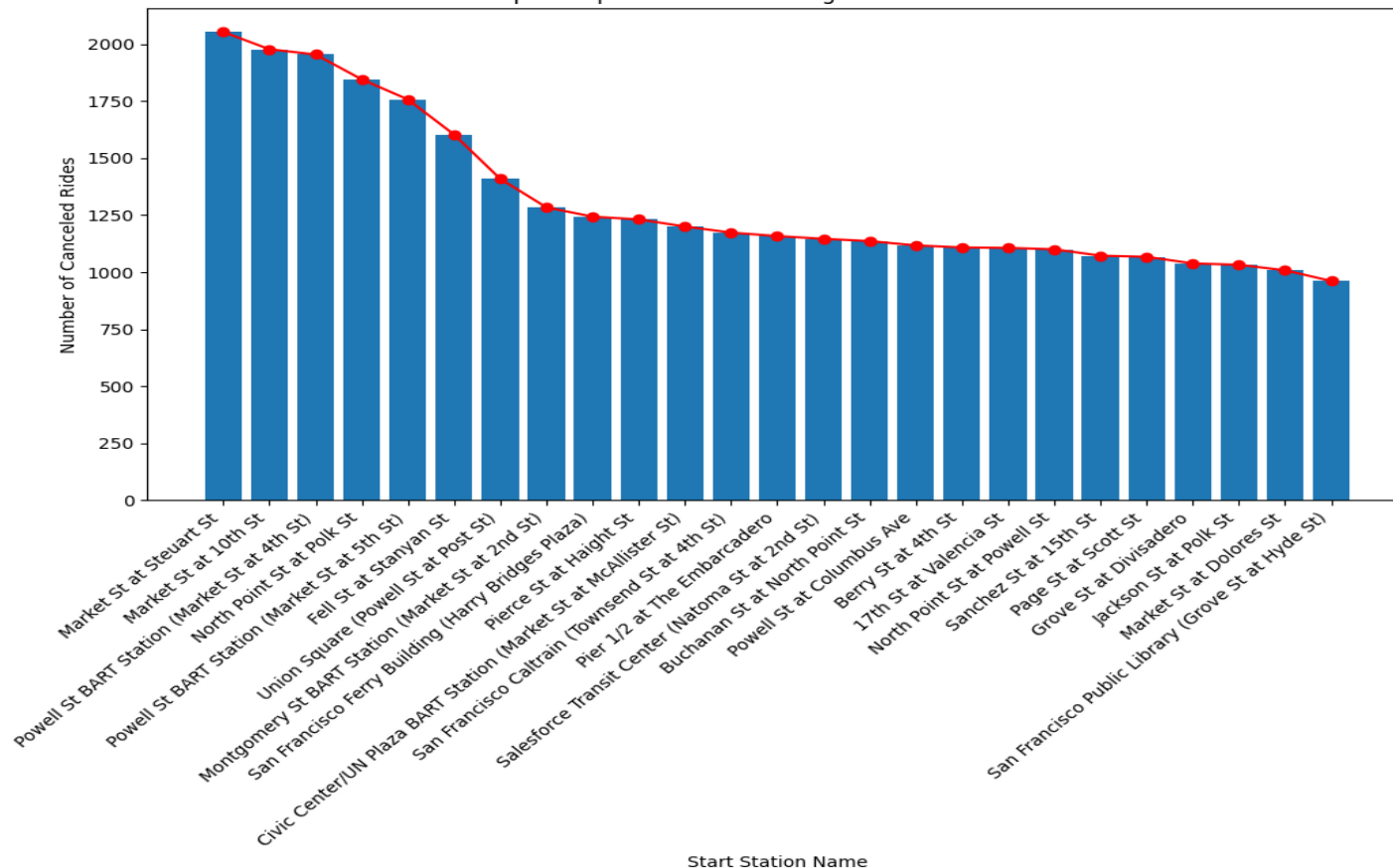To filter the DataFrame to retain only valid rides by excluding two cases:

- Keeping rides with a duration of 5 minutes or more.
- Excluding rides with less than 5 minutes duration where the start and end stations are the same.
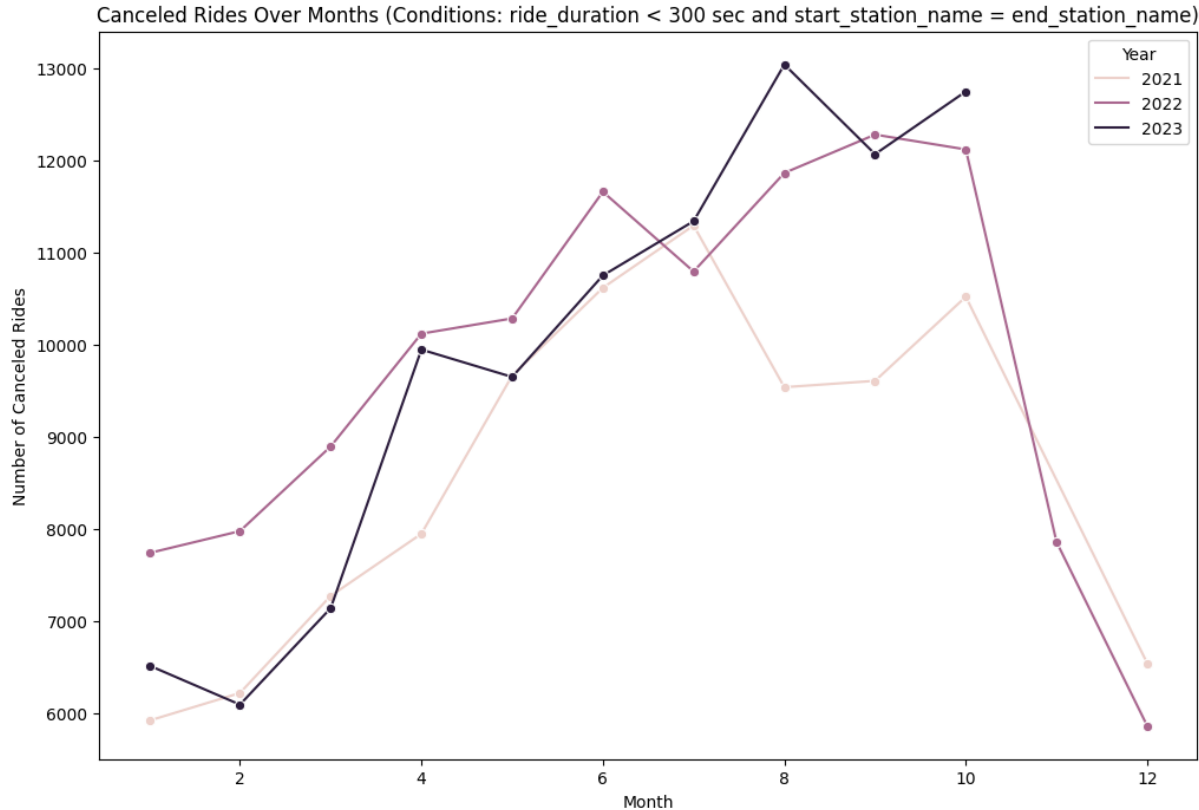
This ensures the Data Frame contains reliable ride data, filtering out potentially canceled or very short rides.
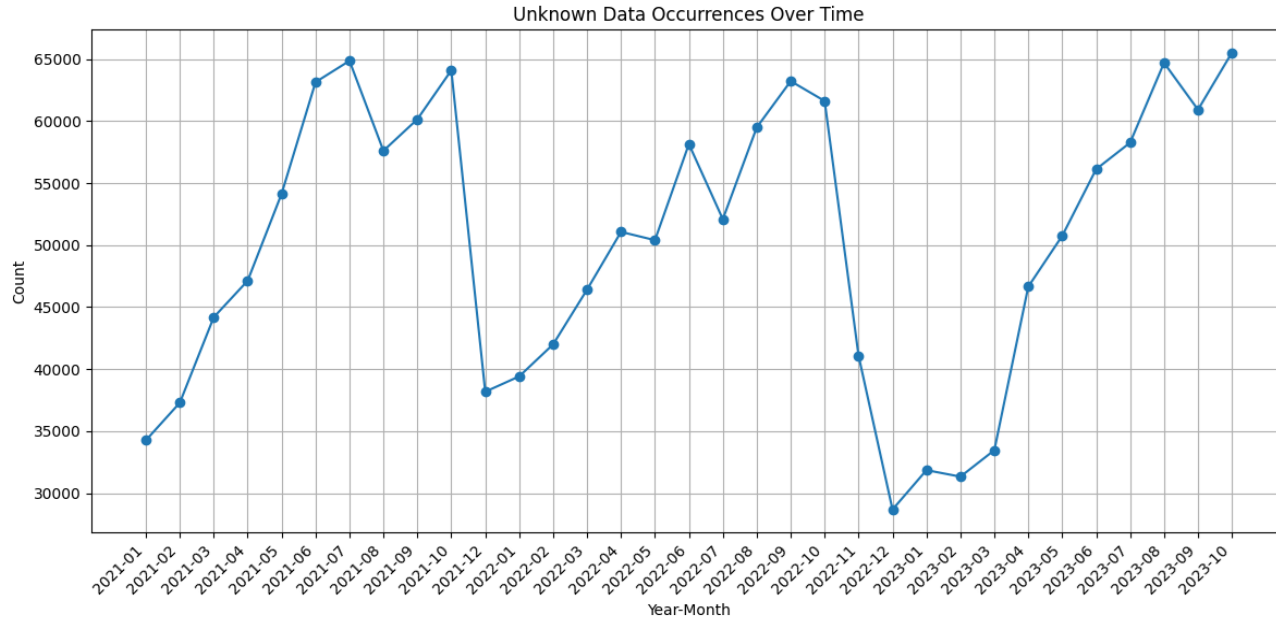
Number of potentially canceled rides: 311875

Total number of distinct stations: 603

Canceled Rides Over Months (Conditions: ride_duration < 300 sec and start_station_name = end_station_name)

The graph shows that there are more canceled rides in the summer months, and that the number of canceled rides has increased over time.

Unknown Data Occurrences Over Time

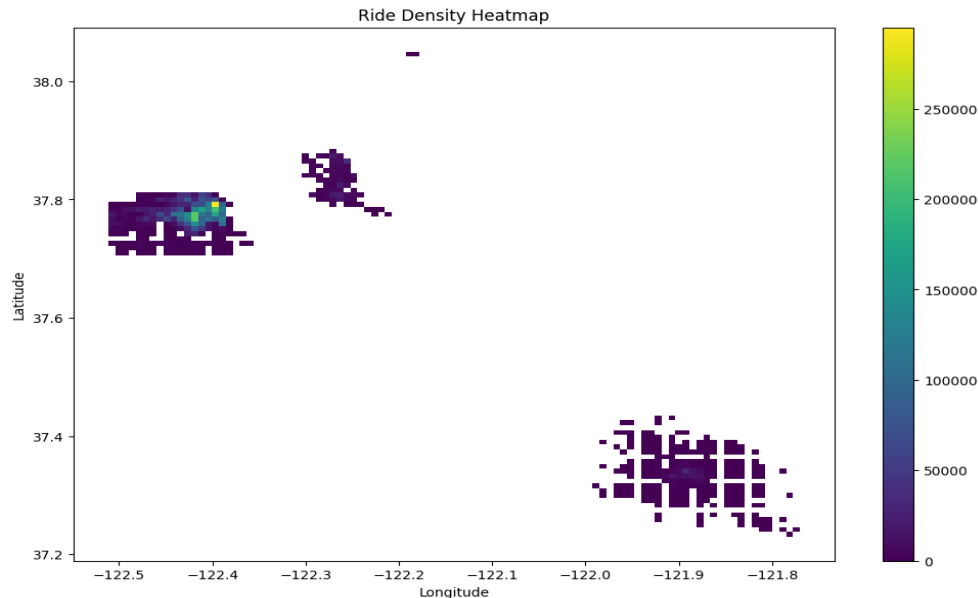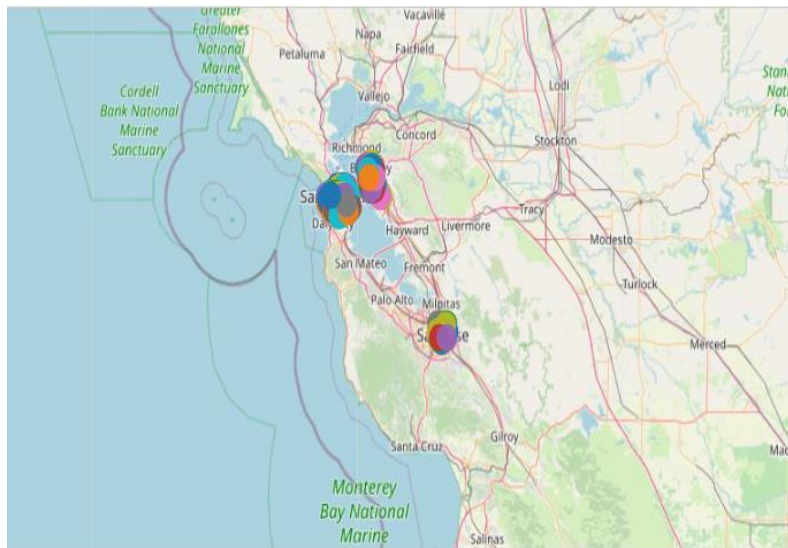The unknown data is likely due to a number of factors, such as:
- GPS errors
- User input errors
- Stations that are no longer in service
- Stations that are not properly identified in the dataset

Unknown data occurrences in Lyft Bay wheels rides dataset decreasing over time, but still significant.

This insight suggests that Lyft is making progress in improving the accuracy and completeness of its dataset, but there is still room for improvement.
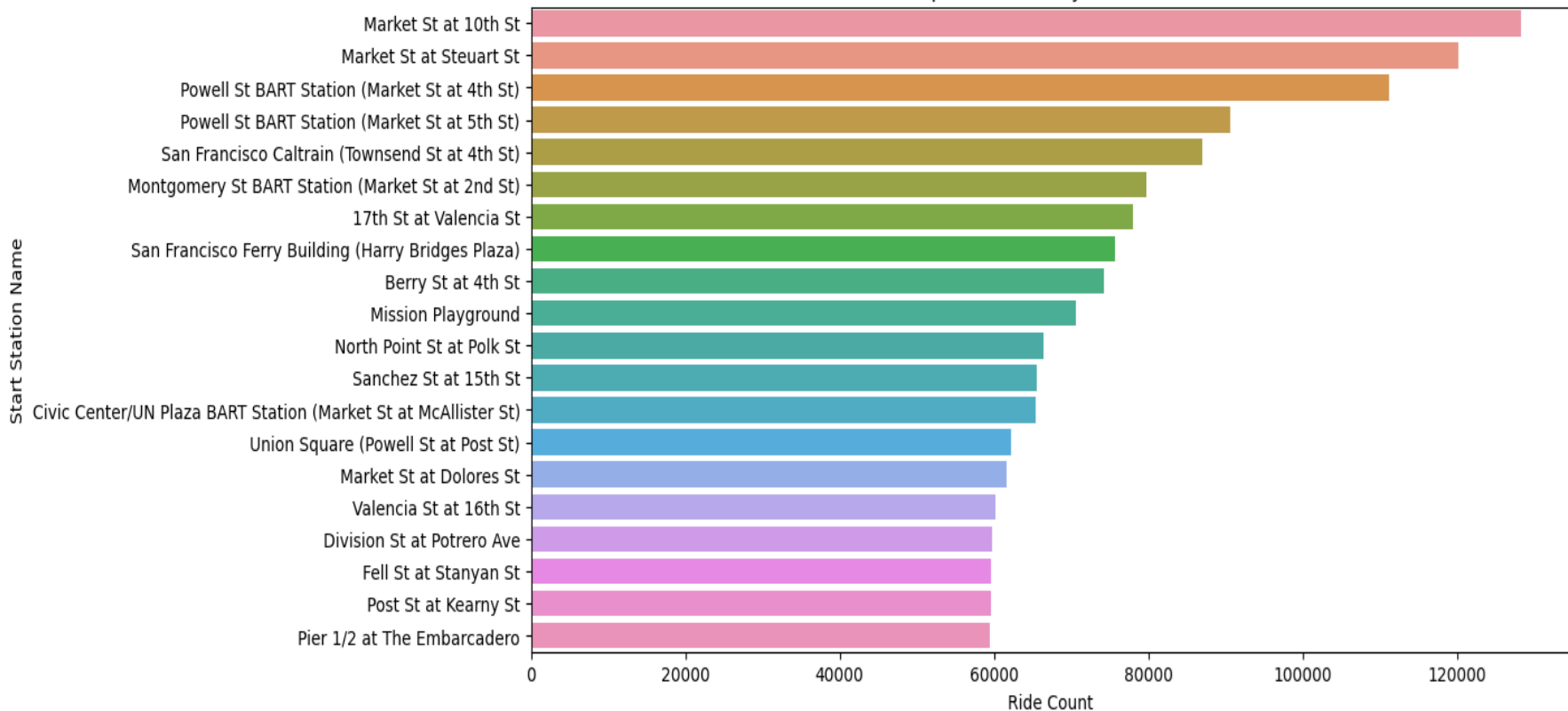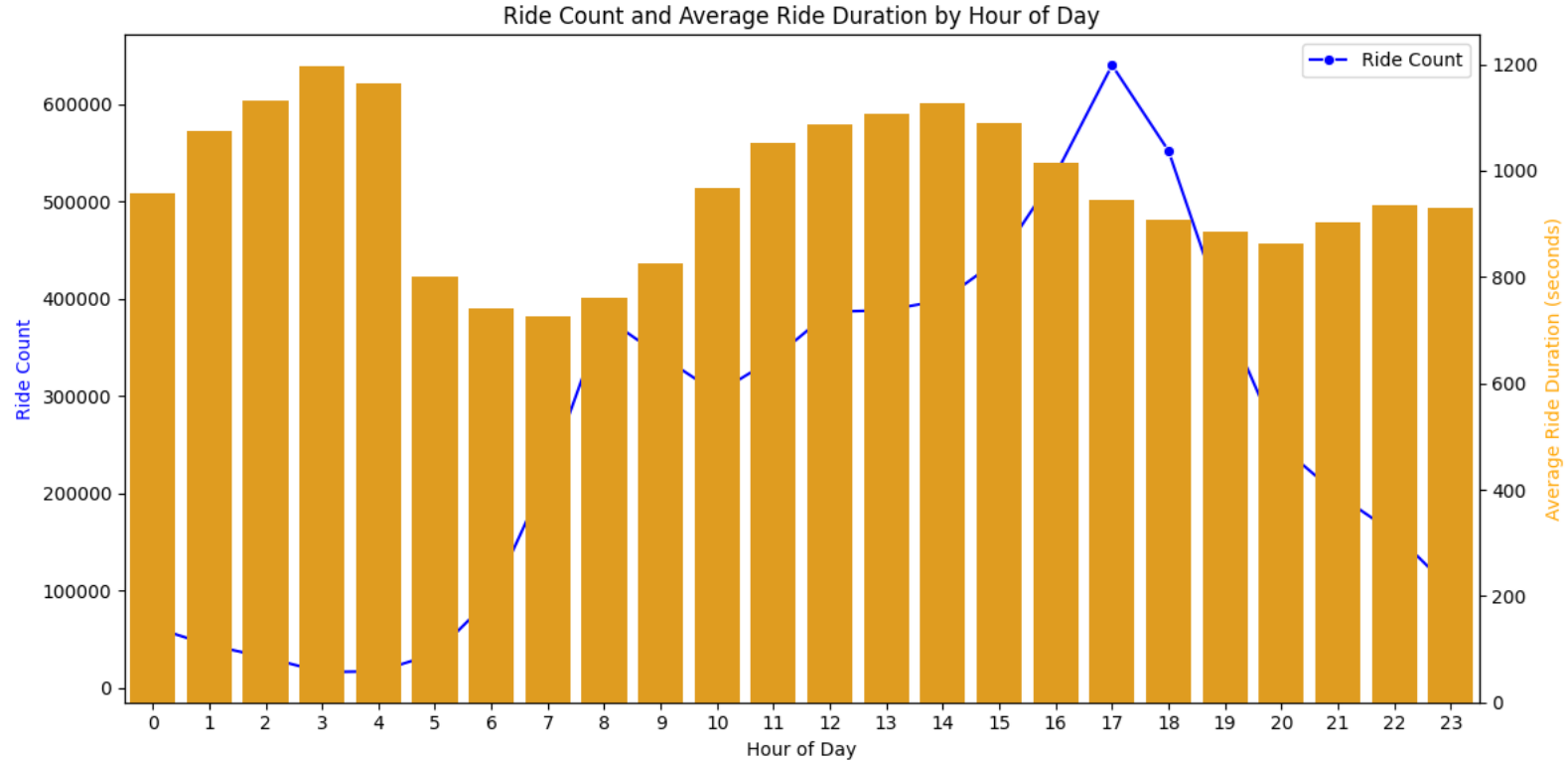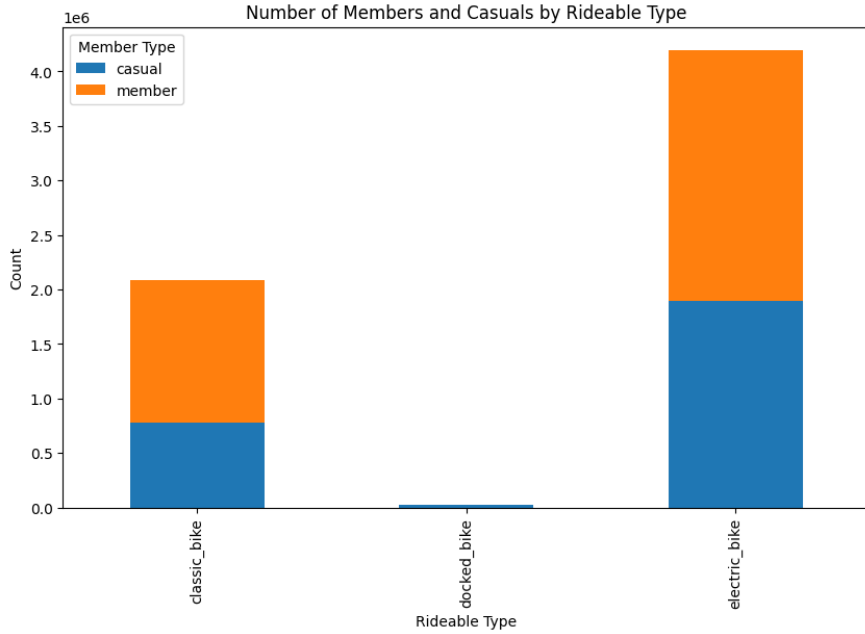
# Stations



Ride Density Heatmap

Lyft bike ride start stations are located throughout the San Francisco Bay Area, with a concentration in the downtown area.
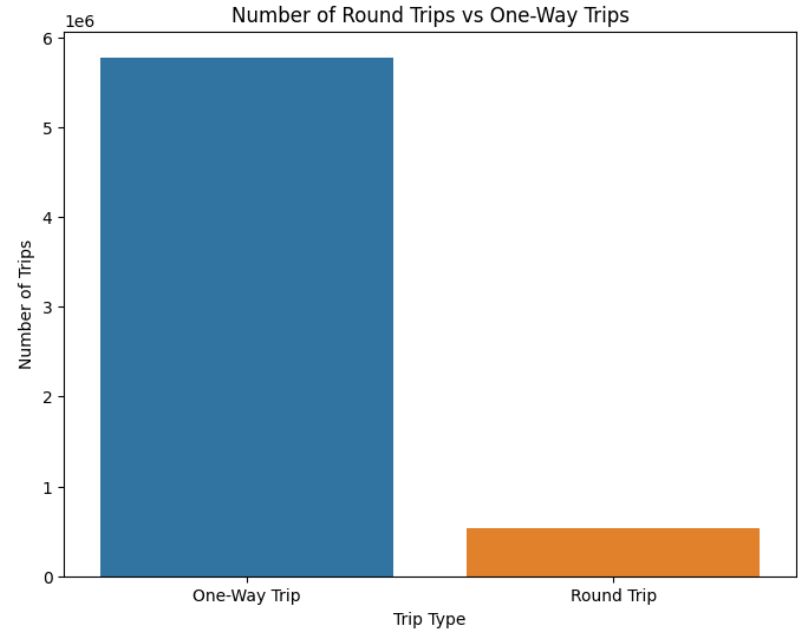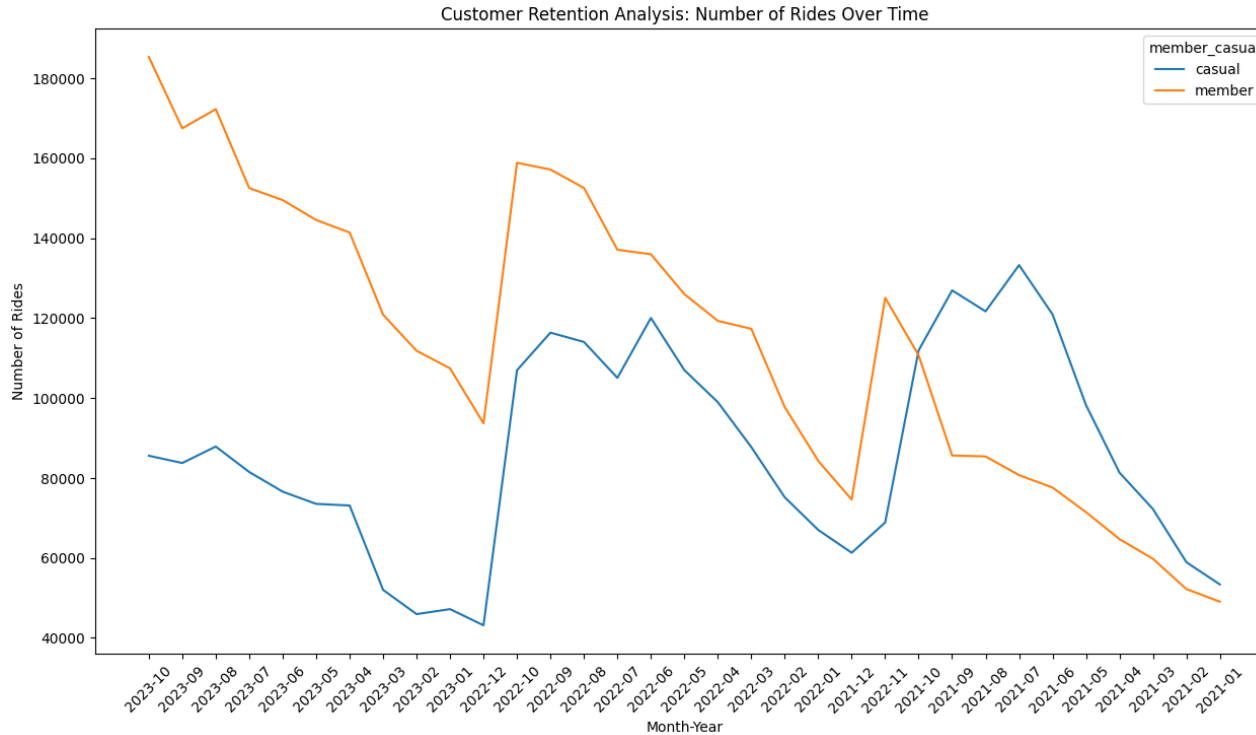
Top 20 Stations by Ride Count

Ride Count and Average Ride Duration by Hour of Day

This insight suggests that people are more likely to take shorter rides during the day, and longer rides at night.

Number of Members and Casuals by Rideable Type

```
+--------------+-------+
|rideable_type|  count|
+--------------+-------+
|electric_bike|4197496|
| classic_bike|2085406|
|  docked_bike|  25252|
+--------------+-------+
```
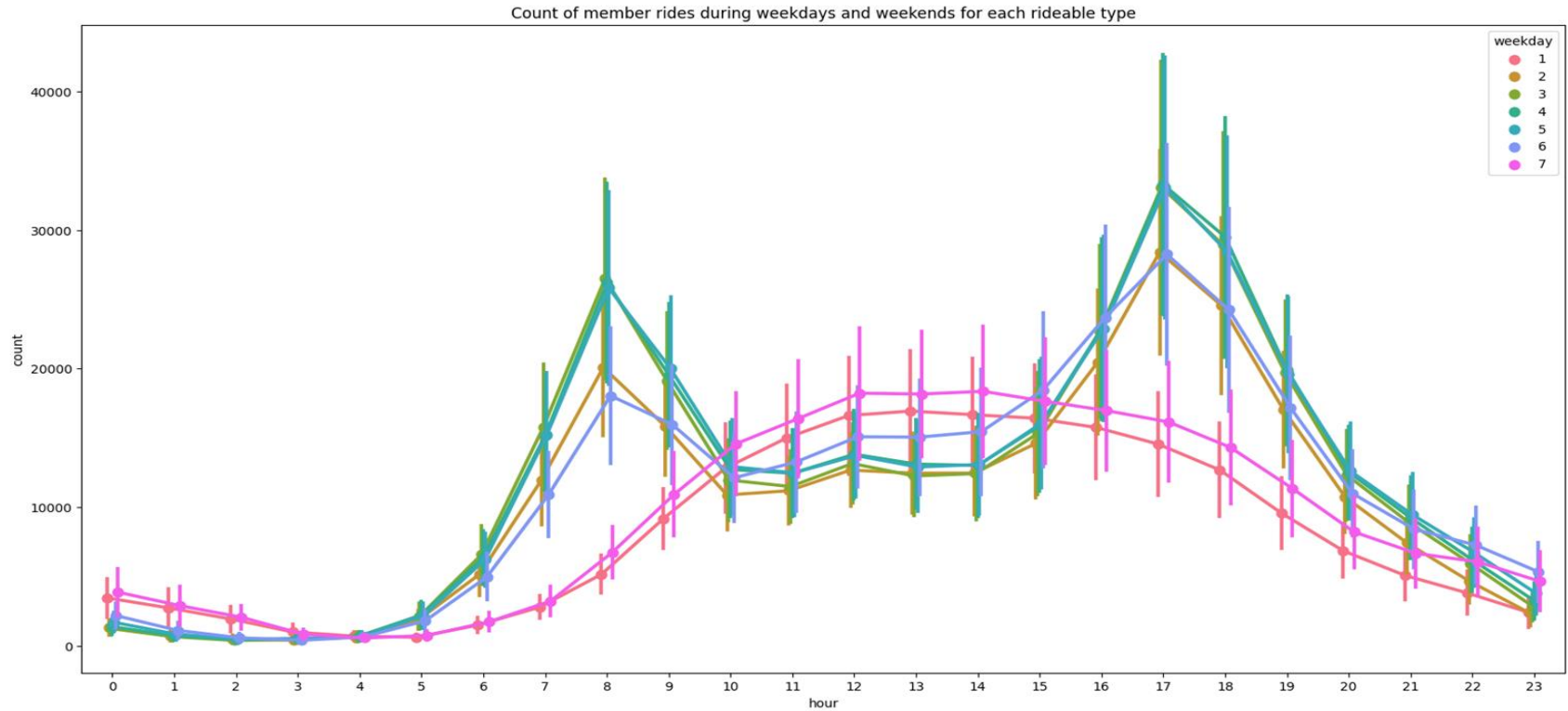
Number of Round Trips vs One-Way Trips

```
+--------------+-------+
|member_casual|  count|
+--------------+-------+
|       casual|2702435|
|       member|3605719|
+--------------+-------+
```
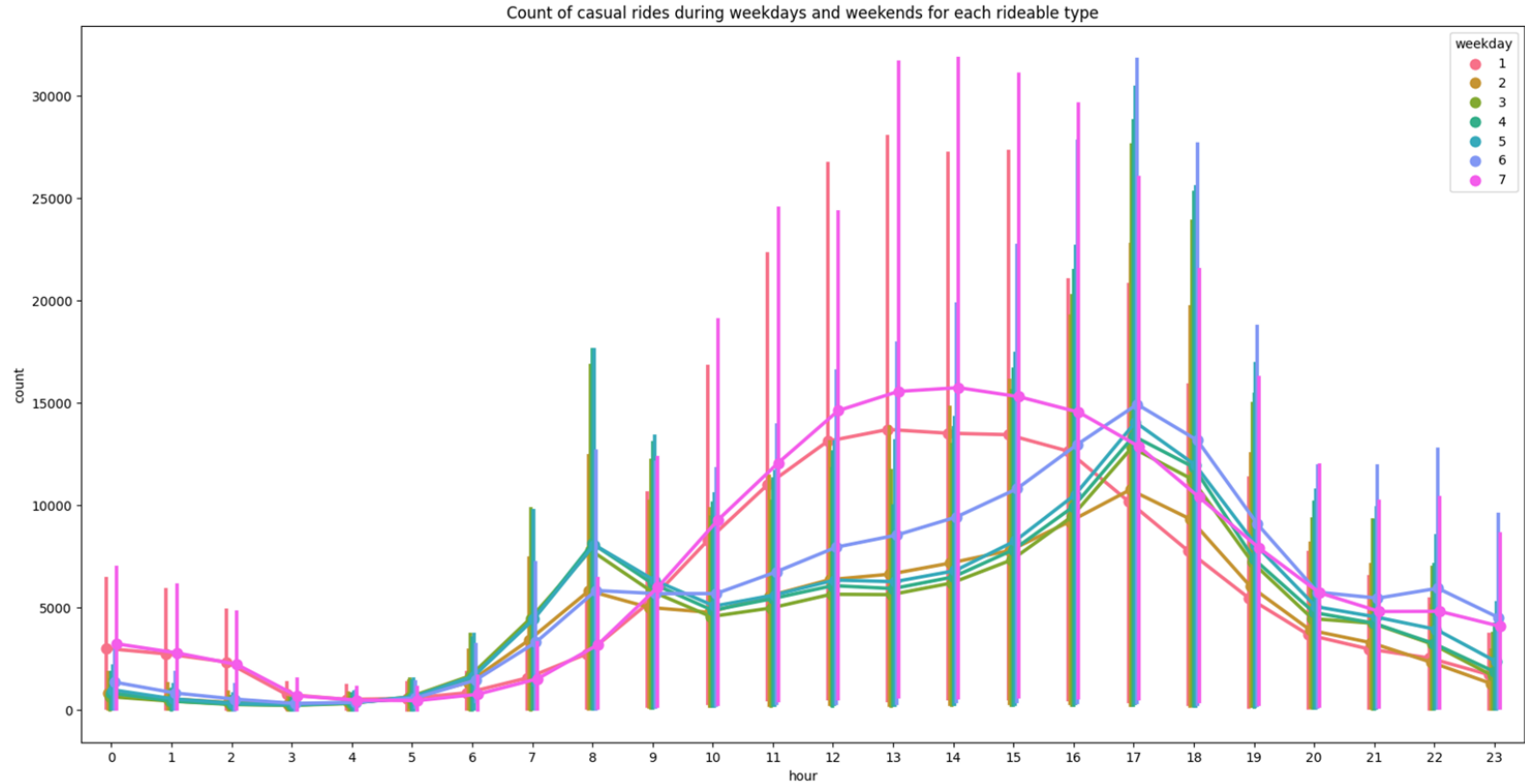
Customer Retention Analysis: Number of Rides Over Time

Lyft member rides increasing over time, while casual rider rides remain stable.
This insight suggests that Lyft is becoming more popular among regular riders, which could be due to a number of factors, such as the convenience and affordability of the service.

Count of member rides during weekdays and weekends for each rideable type

Average number of member rides higher on weekdays than on weekends.

Count of casual rides during weekdays and weekends for each rideable type

Average number of Casual rides higher on weekends than weekdays.

Bay Wheels rides peak in the morning and evening, lowest in the middle of the day and at night.

Total Rides During Holidays and Non-Holidays by Rideable Type and Member Type

The holiday column includes weekends (Saturday and Sunday) and other specified holidays, adopted from Office Holidays - California, USA.

Distance Traveled by Member and Casual Riders by Rideable Type



Average speed Traveled by Member and Casual Riders by Rideable Type

Casual riders travel further distances than member riders, regardless of rideable type.

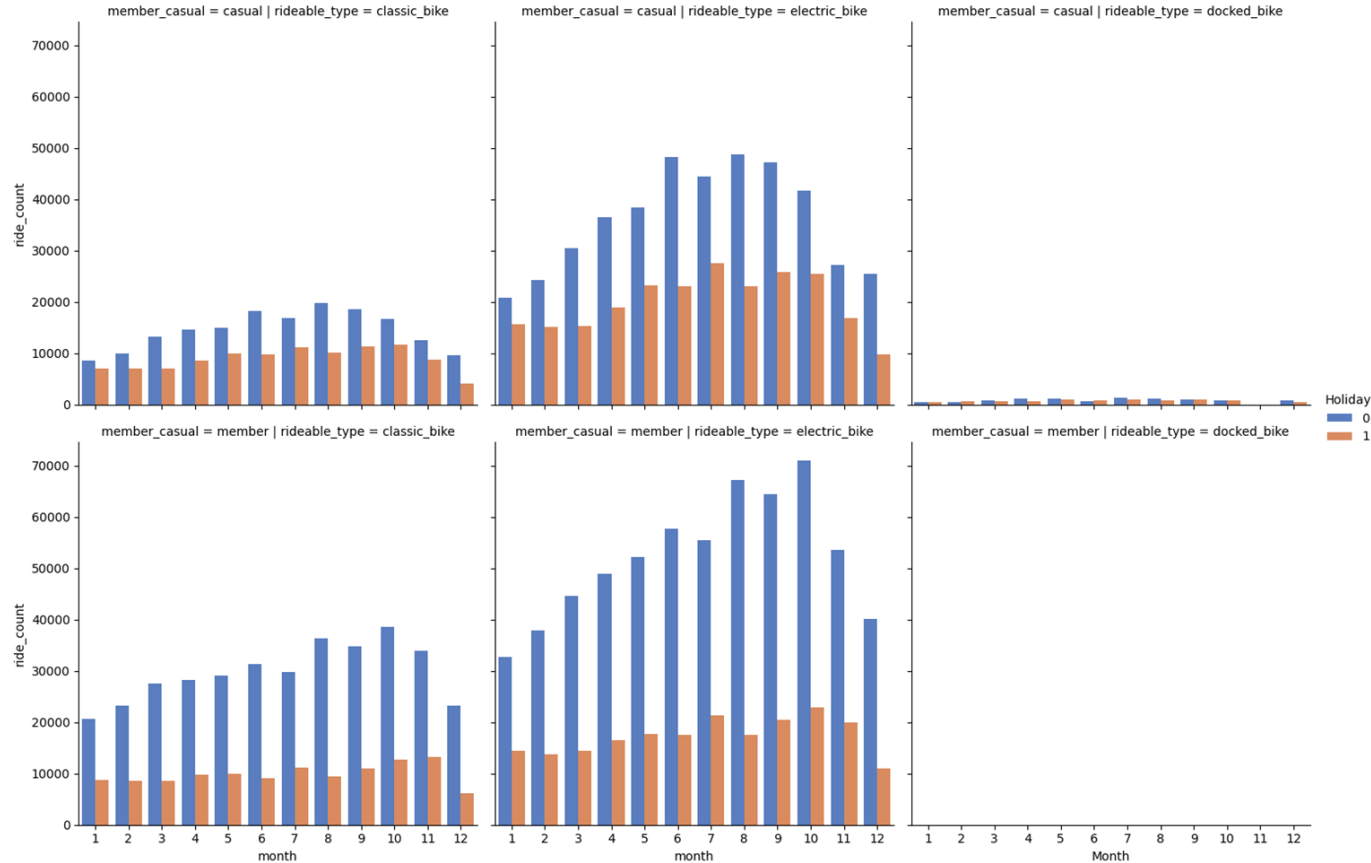Member riders travel faster than casual riders.

The geolocation distance calculation using the Haversine formula, as implemented through a User-Defined Function (UDF) and inspired by the article 'Calculating Distance Between Two Geolocations in Python', adds a new 'distance(miles)' column to the DataFrame using the latitude and longitudinal information and duration of rides.

Cardinal Direction Popularity

The largest segment, representing rides heading North, indicates a higher preference or need for travel in this direction. It might also reflect commuting patterns, which could be vital for urban planning and expansion of services. Understanding why certain directions are preferred could help Lyft optimize its network to better serve the community's transportation needs.

Inspired by the method detailed in the mapscaping reference, a User-Defined Function (UDF) is implemented to calculate the bearing between geographical coordinates, offering valuable insights into the direction of movement for analyzing ride patterns and trends in the dataset.

Correlation Matrix

The correlation matrix shows how different features of Lyft Bay wheels rides are related to each other.

# Predicting Potential Bike Availability at Specific Stations

**Bike Availability Prediction Workflow:**

- Calculate hourly ride counts, create DataFrames, and join.

- Handle missing values, aggregate counts, and compute availability.

- Merge availability info with the original data.

- Select features, transform for ML, and train Machine Learning Algorithms.

Overall, the ML encompasses preprocessing, feature engineering, and machine learning to predict potential bike availability.

**Availability column Assumptions:**

- Positive differences imply increased availability, assuming immediate reflection.

- Negative differences set to 0 uniformly, potentially overlooking fluctuations.

- Binary classification (0 or 1) simplifies availability, missing nuanced variations.

- Null values filled with 0 assume no rides mean no availability, potentially inaccurately handling missing data.

# Machine Learning Model Performance Evaluation

The ROC (Receiver Operating Characteristic) curve is a metric for evaluating the performance of classification models. Accuracy is the proportion of correctly classified instances out of the total instances. Here are the ROC curve scores and Accuracy scores for different machine learning models:

- Logistic Regression: 0.644 (ROC), 0.843(Accuracy)
- Random Forest: 0.688(ROC), 0.782(Accuracy)
- Multilayer Perceptron Classifier: 0.627(ROC), 0.843(Accuracy)

A higher ROC score and Higher Accuracy indicates better model performance in distinguishing between classes. The Logistic Regression Classifier demonstrates the highest ROC and Accuracy scores among the evaluated models, suggesting superior discriminatory power.

# Insights & Urban Mobility Impact

- Lyft is making it easy for people to find a bike to ride, no matter where they are in the city. This can help to encourage people to choose biking as a transportation option, which can have a number of benefits, such as reducing traffic congestion and improving air quality.

- Electric bikes are popular among casual riders, covering longer distances at higher speeds compared to classic bikes.Peak ride hours and daily patterns reveal user behavior, aiding operational planning and resource allocation.

- Ride counts fluctuate during holidays; understanding these patterns helps optimize services and promotions.

- Cardinal direction analysis uncovers preferred routes, guiding station placements for enhanced accessibility.

- Multilayer Perceptron model outperforms, providing reliable predictions of ride availability.

- Data-driven insights reveal general user trends, guiding business strategies and user experience improvements.

- The Urban mobility trends from this analysis will contribute valuable insights to urban planners, policymakers, and transportation enthusiasts interested in understanding and optimizing the dynamics of bike-sharing systems.

# Future Scope

- Apply the developed framework and insights to other cities with similar bike-sharing systems, adapting the model to local variations. This could contribute to a scalable solution applicable in diverse urban environments.
- Can develop a model to suggest optimal routes for rides based on historical traffic patterns, road closures, and other relevant factors.
- Can create a model to optimize the allocation of vehicles to different areas based on predicted demand and user behaviour.

By focusing on these areas, Lyft can potentially improve service reliability, customer satisfaction, and operational efficiency, ultimately enhancing the user experience and possibly reducing the rate of canceled rides.

# References

*System data: Bay Wheels*. Lyft. (n.d.). https://www.lyft.com/bikes/bay-wheels/system-data.

Shaw, T. (2020, January 1). *Federal holidays in California in 2023*. Office Holidays.
https://www.officeholidays.com/countries/usa/california.

Bhardwaj, A. (2020, June 14). *Calculating distance between two geolocations in python*. Medium.
https://towardsdatascience.com/calculating-distance-between-two-geolocations-in-python-26ad3afe287b.

*How to calculate bearing between two coordinates - December 4, 2023*. mapscaping.com. (2023, November 7).
https://mapscaping.com/how-to-calculate-bearing-between-two-coordinates/#:~:text=In%20Microsoft%20Excel%2C%20you%20can,x%2Daxis%20and%20the%20point.

Thank You…