

# Detailed Report on NYC Taxi Trip Duration Analysis and Predictive Modeling

## Introduction:

New York City, one of the most bustling and dynamic cities in the world, is renowned for its iconic yellow taxis. These taxis play an essential role in the city's transportation system, carrying hundreds of thousands of passengers each day across numerous neighborhoods. The efficiency and reliability of taxi services are crucial aspects that significantly impact the lives of New Yorkers and visitors. Accurate prediction of taxi trip duration can provide valuable insights to taxi drivers, service providers, city planners, and passengers. It can help optimize routes, improve service reliability, manage the fleet effectively, and enhance overall customer satisfaction by providing accurate trip estimates.

Understanding and predicting taxi trip duration is a complex task due to the myriad factors that can influence it. These factors range from the time of day, day of the week, and season to weather conditions, traffic variations, road network structure, and pickup and dropoff locations. Dealing with such complexity requires advanced data analysis and machine learning techniques that can capture the nuanced relationships within the data.

In this context, the objective of this report is to analyze the New York City Taxi Trip Duration dataset, understand the factors influencing trip duration, and build a predictive model that can accurately estimate the duration of a taxi trip. The dataset includes details of individual rides made in New York City taxis in 2016. It provides an excellent opportunity to study and model the factors influencing taxi trip duration in an urban setting.

The report will present a detailed analysis of the dataset, elaborate on the exploratory data analysis performed, explain the data preprocessing and feature engineering steps undertaken, discuss the various predictive models implemented, and finally, conclude with insights and recommendations based on the findings. The ultimate goal is to derive meaningful insights that could help enhance the efficiency and reliability of taxi services in New York City.

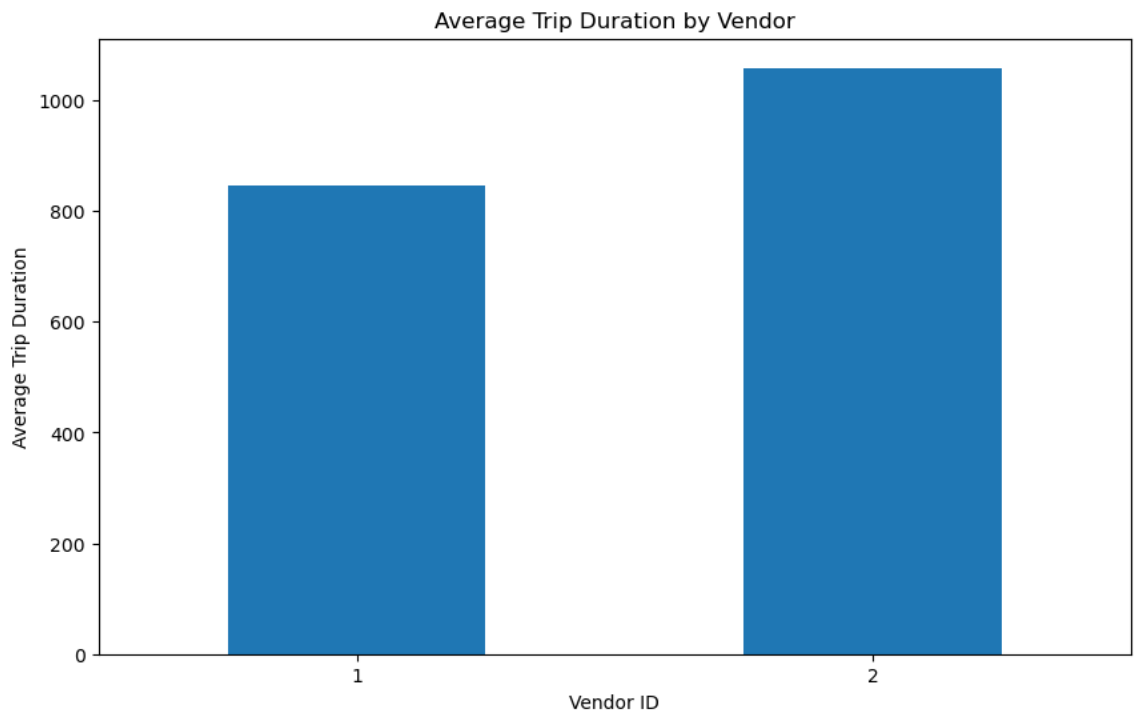
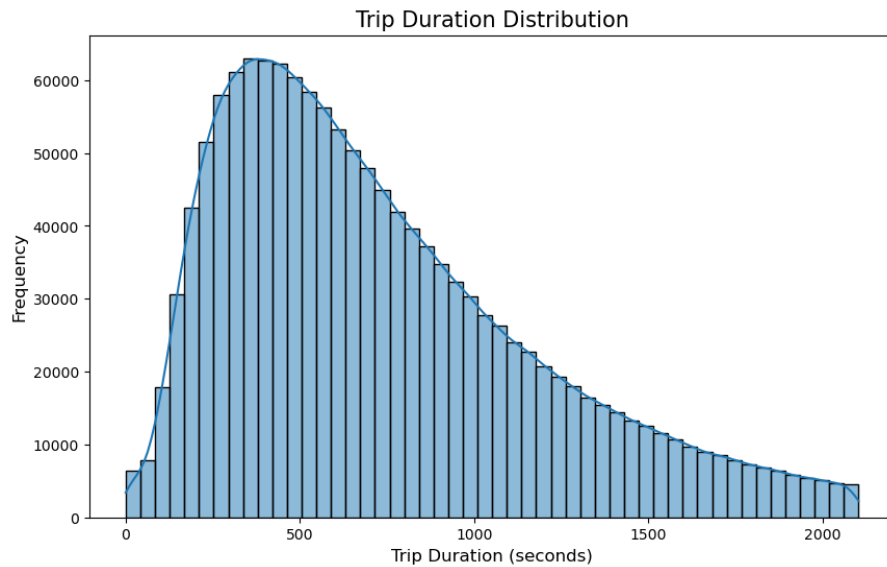
## Dataset Description:

The dataset provides information about individual taxi trips in New York City in 2016. It includes pickup times, drop-off times, pickup and drop-off geographic coordinates, the number of passengers, and the trip duration in seconds.

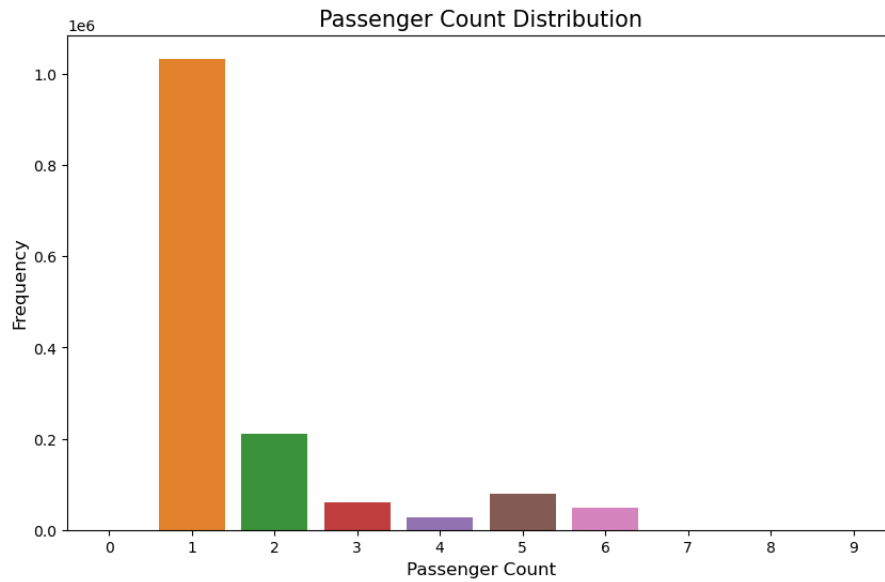
## Exploratory Data Analysis (EDA):

An in-depth EDA was carried out to understand the data's underlying structure, anomalies, and patterns. Key findings include:

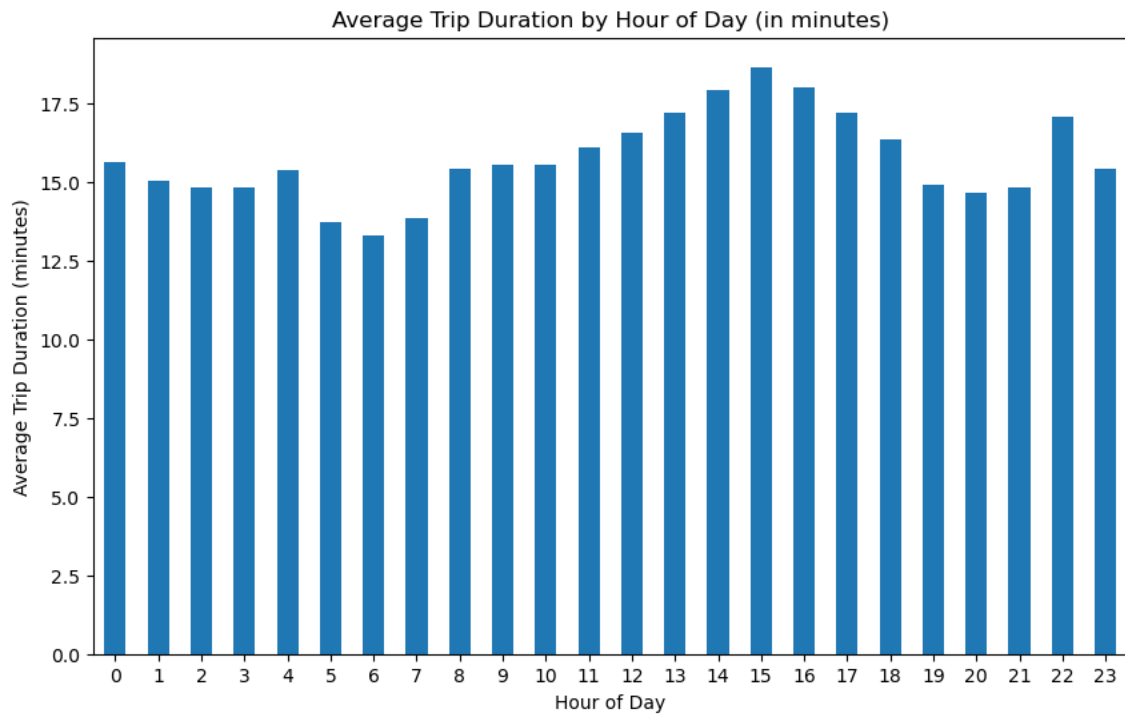
1. **Trip Duration:** The distribution of trip durations is right-skewed, with a majority of trips lasting between 5 to 20 minutes. There are a few extreme cases of very long trips, suggesting the presence of outliers.

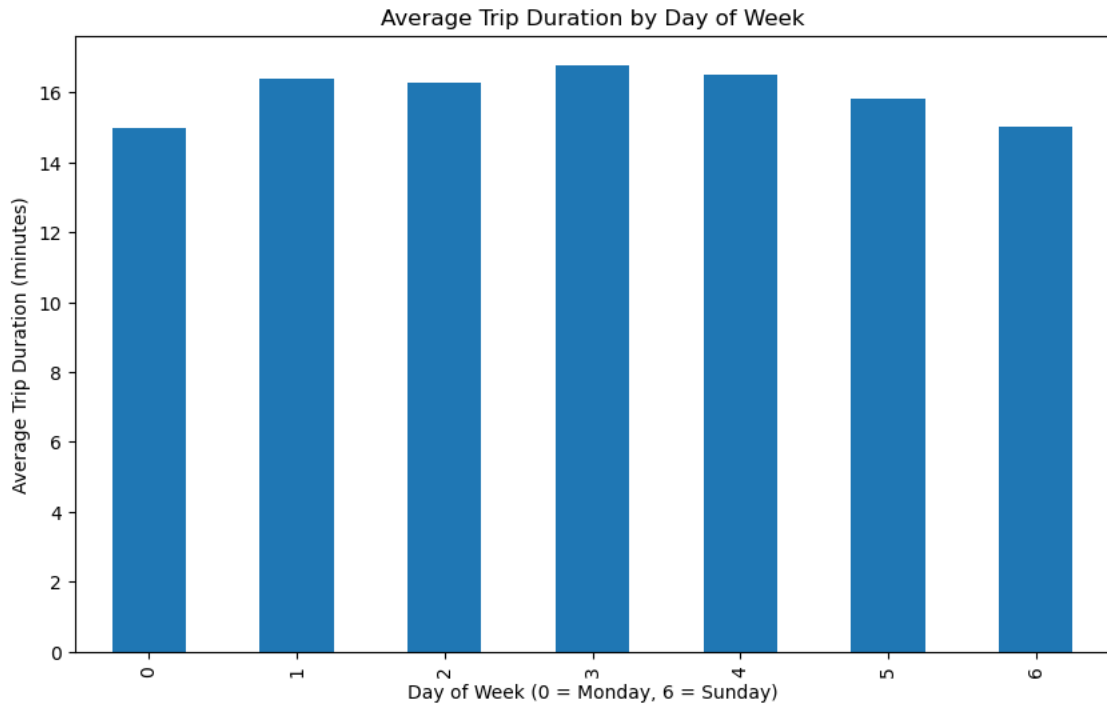


2. **Passenger Count:** Most of the taxi rides had one or two passengers. There was no substantial correlation between the number of passengers and the trip duration.

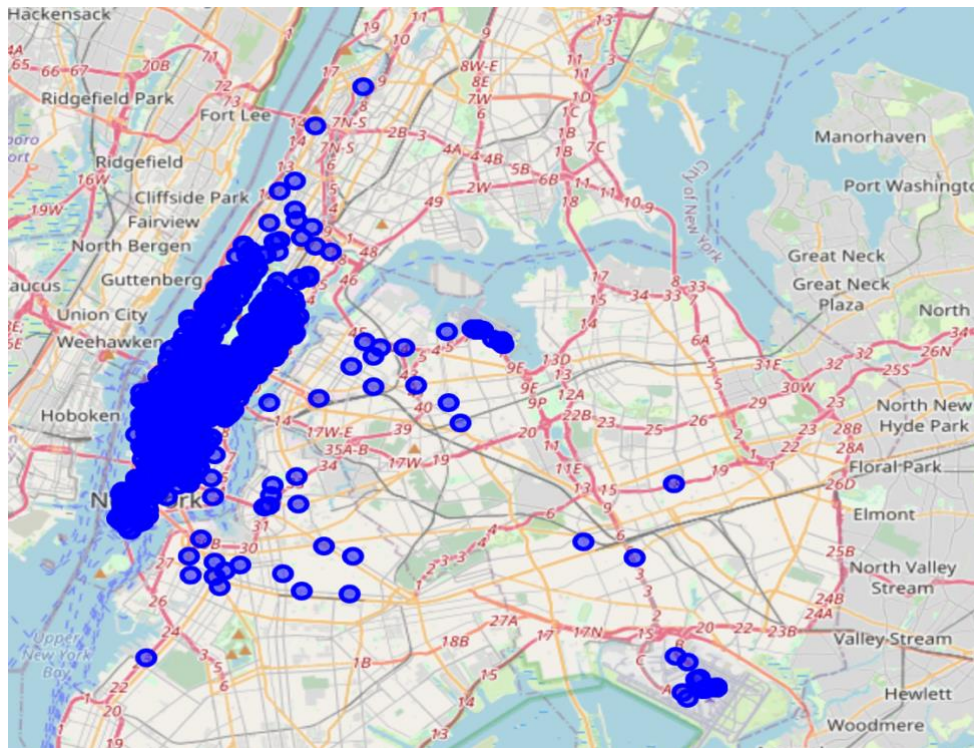


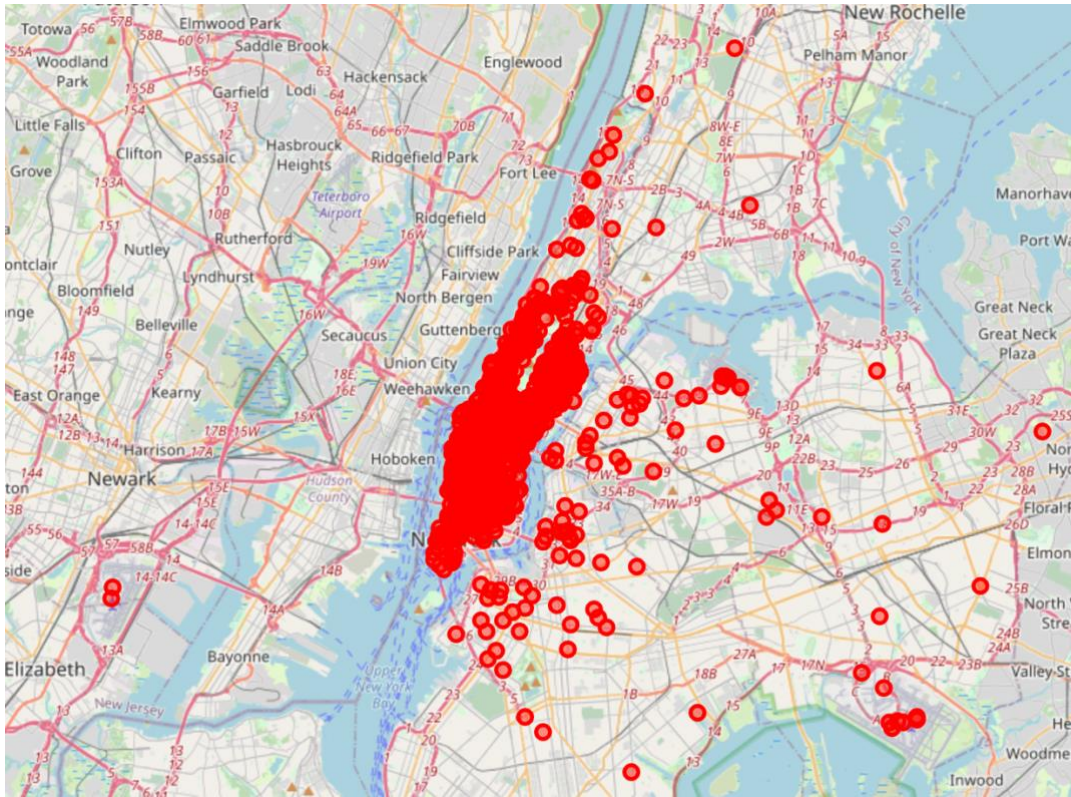
3. **Time Analysis:** We broke down the pickup time to hours, days, and months. We found that the average trip duration varies significantly depending on the hour of the day, with longer durations observed during peak hours. This variation is likely attributed to changing traffic conditions throughout the day.



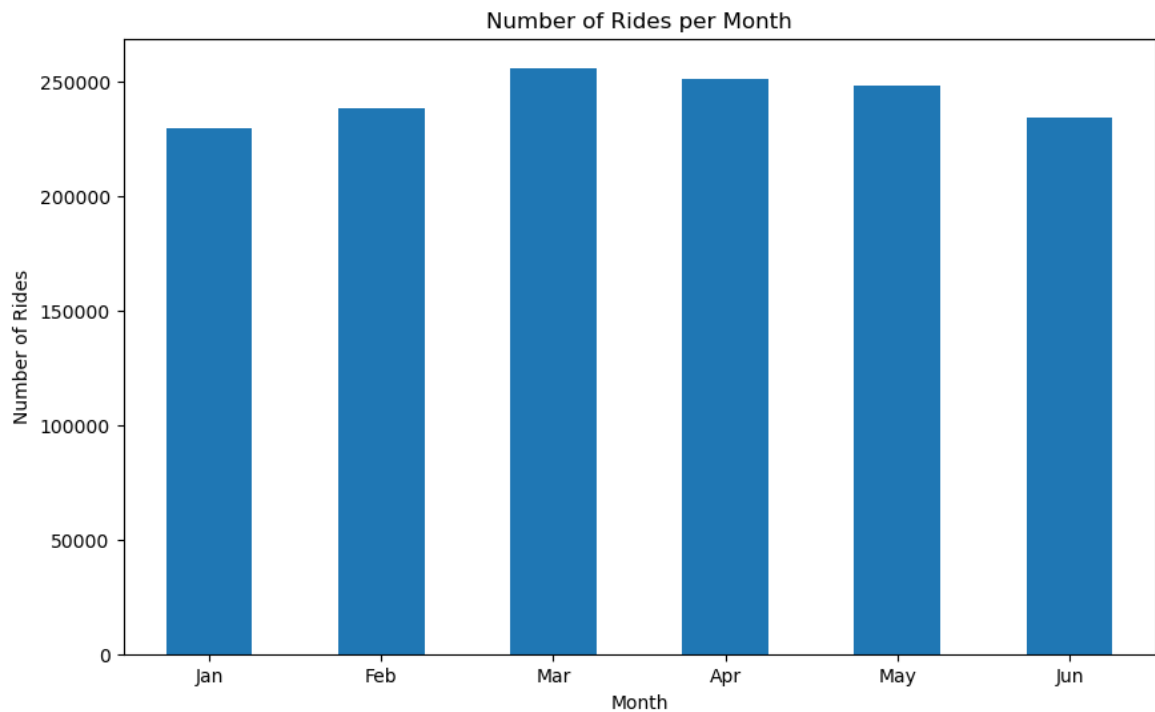


4. **Location Analysis:** Visualization of pickup and drop-off locations (blue for pick-up and red for drop-off) revealed a concentration of trips within Manhattan and the two airports (JFK and LaGuardia).



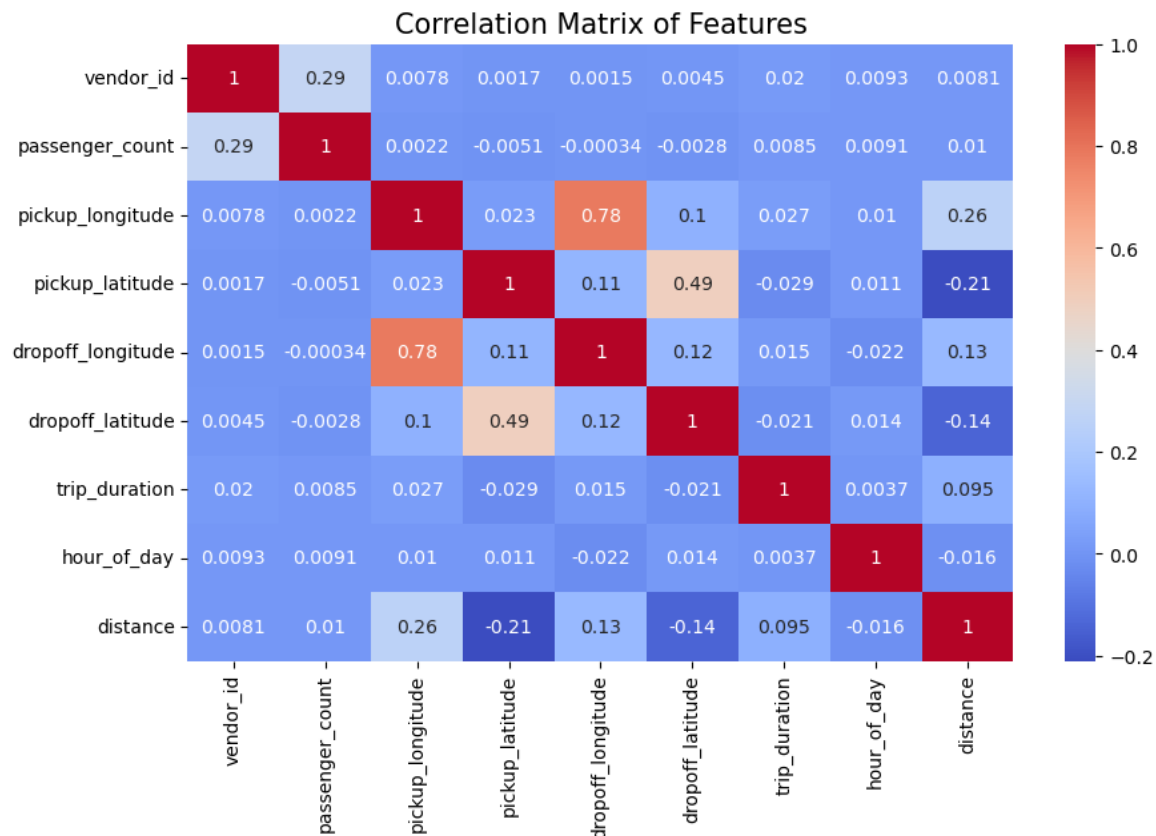


5. **Monthly Analysis:** The number of taxi rides varied across different months, suggesting the influence of seasonal trends or specific events on taxi demand.



## Correlation Analysis:

The correlation analysis was conducted to measure how each variable interacts with every other variable in the dataset. The correlation coefficients range from -1 to +1, where -1 indicates a perfect negative correlation, +1 a perfect positive correlation, and 0 no correlation.



Key insights from the correlation matrix are:

- Vendor ID:** The vendor id showed a very low correlation with all other features, indicating that the vendor's identity has a minor role in affecting the characteristics of a trip.
- Passenger Count:** Similarly, passenger count showed a minimal correlation with the other variables. It suggests that the number of passengers in a ride doesn't majorly influence the trip's other aspects.
- Pickup and Dropoff Locations:** Pickup longitude and dropoff longitude had a high correlation (0.78), suggesting that trips tend to start and end within the same longitudinal zones. The latitude pair also showed a moderate correlation (0.49), suggesting a similar pattern.



4. **Distance:** The calculated trip distance showed a moderate positive correlation with pickup longitude (0.26) and dropoff longitude (0.13). However, it displayed a moderate negative correlation with pickup latitude (-0.21) and dropoff latitude (-0.14). These correlations could be a result of the city's geographical layout and commonly used routes.
5. **Trip Duration:** The duration of the trip had a small positive correlation with distance (0.09). This indicates that longer distances tend to lead to longer trip durations, which is an expected finding. However, the correlation is not very strong, indicating the existence of other contributing factors such as traffic conditions and route selection.
6. **Hour of the Day:** The hour of the day showed a negligible correlation with trip duration (0.0037), suggesting that the time of the day might not be a major determinant for the duration of the trip.

The correlation analysis provides valuable insights into the relationships between different aspects of taxi trips. However, it's worth noting that correlation doesn't imply causation and only captures linear relationships between variables. Therefore, while these insights are useful in feature selection and initial understanding, more sophisticated techniques such as regression analysis and machine learning models should be used for prediction tasks.

## Data Preprocessing and Feature Engineering:

To prepare the data for modeling, we undertook several preprocessing and feature engineering steps. The outliers in trip duration were removed. New features were engineered from the date-time and geographic data. Specifically, we extracted the hour of the day, day of the week, and month from the pickup time. We also calculated the Haversine distance (the great-circle distance between two points on a sphere given their longitudes and latitudes) to get the travel distance for each trip.

## Predictive Modeling:

We employed several regression models to predict the trip duration:

```
Linear Regression – Train RMSE: 9.265788444331472, Test RMSE: 9.411068548558989
Decision Tree – Train RMSE: 0.11768473355061854, Test RMSE: 7.366492250968637
Random Forest – Train RMSE: 1.9759691799867796, Test RMSE: 5.278336779331367
Extra Trees Regressor – Train RMSE: 0.11768473355061856, Test RMSE: 5.390310653896816
XGBoost – Train RMSE: 5.255748178616791, Test RMSE: 5.3327928979550965
```

1. **Linear Regression:** The Linear Regression model didn't perform very well. The model demonstrated underfitting, evidenced by a relatively high Root Mean Squared Error (RMSE) for both the training and testing data. This result indicates that the model may be too simple to capture the complexities of the dataset.
2. **Decision Tree and Extra Trees Regressor:** These models appeared to overfit the training data, as indicated by a very low RMSE on the training data and a significantly higher RMSE on the testing data.
3. **Random Forest and XGBoost:** The Random Forest and XGBoost models performed the best among all models. They were able to fit the training data well while also demonstrating good generalizability to unseen data, as the difference between training and testing RMSE was small.

## Conclusion and Recommendations:

The analysis of the NYC Taxi Trip Duration dataset provided several significant insights. Key temporal and geographical factors such as the hour of the day, day of the week, month, and travel distance significantly impact the trip duration.

The Random Forest and XGBoost models were found to be most successful in capturing these relationships, delivering reliable predictions on the testing data.

Going forward, further improvements can be made by exploring additional feature engineering opportunities, refining model hyperparameters, and using ensemble methods to leverage the strengths of different models.

Moreover, the integration of external data sources such as real-time traffic data or weather conditions, which are known to affect travel time, could offer even more sophisticated and accurate predictive models.