# Analyzing Transportation Trends using NYC Taxi Dataset

## Problem Description:

### Introduction:

The transportation industry is an essential aspect of modern society, connecting people and goods to different destinations. However, the industry faces several challenges, including congestion, safety concerns, and the need to improve efficiency. To address these challenges, transportation companies and policymakers require accurate data and analysis to understand current trends and predict future transportation needs.

### Project Objectives:

The objective of this project is to analyze transportation trends using the New York City Taxi dataset and develop predictive models to forecast future trends. The project aims to achieve the following objectives:

1. To explore the NYC Taxi dataset and identify relevant variables that impact transportation trends.
2. To clean and preprocess the dataset to ensure data quality and prepare it for analysis.
3. To perform exploratory data analysis to gain insights into transportation trends and identify patterns in the data.
4. To develop predictive models using machine learning algorithms to forecast future transportation trends.
5. To evaluate the predictive models and choose the best model based on its performance.

### Dataset Link:

The New York City Taxi dataset used in this project is available on Kaggle at https://www.kaggle.com/c/nyc-taxi-trip-duration. The dataset contains information about taxi trips in New York City, including pick-up and drop-off locations, trip durations, passenger counts, and other relevant variables.

**Dataset Description:**

The NYC Taxi dataset contains the following variables:

1. id - a unique identifier for each taxi trip
2. vendor_id - a code indicating the provider associated with the trip record
3. pickup_datetime - the date and time when the meter was engaged
4. dropoff_datetime - the date and time when the meter was disengaged
5. passenger_count - the number of passengers in the vehicle
6. pickup_longitude - the longitude where the meter was engaged
7. pickup_latitude - the latitude where the meter was engaged
8. dropoff_longitude - the longitude where the meter was disengaged
9. dropoff_latitude - the latitude where the meter was disengaged
10. store_and_fwd_flag - this flag indicates whether the trip record was held in vehicle memory before sending to the vendor because the vehicle did not have a connection to the server
11. trip_duration - duration of the trip in seconds

**Deliverables:**

The deliverables for this project include:

1. A report outlining the findings and insights from the exploratory data analysis.
2. A set of machine learning models developed to predict future transportation trends.
3. A detailed evaluation of the predictive models and the choice of the best-performing model.
4. A presentation summarizing the project's objectives, methodology, and key findings.

# Possible Framework :

**Step 1: Data Collection and Preprocessing**

1.1. Collect the New York City Taxi dataset from Kaggle at https://www.kaggle.com/c/nyc-taxi-trip-duration.

1.2. Import the dataset into the Python environment using Pandas library.

1.3. Clean and preprocess the data to remove missing values, outliers, and other anomalies.

1.4. Convert the date and time variables into the appropriate format.

**Step 2: Exploratory Data Analysis (EDA)**

2.1. Conduct basic statistical analysis of the dataset, including mean, median, standard deviation, and other summary statistics.

2.2. Visualize the distribution of different variables using histograms, box plots, scatter plots, and other appropriate charts.

2.3. Identify any patterns, trends, or anomalies in the data using data visualization and statistical analysis.

2.4. Conduct correlation analysis to determine the relationship between different variables.

**Step 3: Feature Engineering and Selection**

3.1. Identify relevant features that impact transportation trends, such as pick-up and drop-off locations, trip durations, passenger counts, and other relevant variables.

3.2. Create new features that can provide additional insights into transportation trends, such as distance traveled, time of day, day of the week, and other relevant variables.

3.3. Use feature selection techniques to identify the most important features for predicting transportation trends.

**Step 4: Model Selection and Development**

4.1. Choose appropriate machine learning algorithms for predicting transportation trends, such as linear regression, decision trees, random forests, or neural networks.

4.2. Split the dataset into training and testing sets to evaluate the performance of different models.

4.3. Train different machine learning models using the training dataset.

4.4. Evaluate the performance of each model using appropriate metrics, such as mean squared error, root mean squared error, or R-squared.

4.5. Choose the best-performing model based on its performance on the testing dataset.

**Step 5: Model Deployment and Evaluation**

5.1. Deploy the best-performing model on new data to predict transportation trends.

5.2. Evaluate the accuracy of the predictions using appropriate metrics, such as mean squared error, root mean squared error, or R-squared.

5.3. Refine the model as necessary to improve its performance.

**Step 6: Reporting and Presentation**

6.1. Prepare a report outlining the findings and insights from the exploratory data analysis and predictive modeling.

6.2. Develop a presentation summarizing the project's objectives, methodology, and key findings.

6.3. Present the findings and insights to relevant stakeholders, such as transportation companies, policymakers, or the general public.

# Code Explanation :

**Here is the simple explanation for the code you can find at code.py file.**

**Step 1: Data Collection and Preprocessing**

In this section, we first import the dataset into a Pandas DataFrame using the **pd.read_csv()** function. We then perform some preprocessing steps such as removing missing values using the **dropna()** function, and removing outliers using a filter condition. Finally, we convert the date and time variables into appropriate format using the **pd.to_datetime()** function.

**Step 2: Exploratory Data Analysis (EDA)**

This section involves visualizing and understanding the dataset. We first perform basic statistical analysis using the **describe()** function to get an idea of the distribution of the data. We then create a histogram of the **trip_duration** variable to see its distribution. We also create a box plot of **trip_duration** by **passenger_count**, a scatter plot of **trip_duration** by **distance**, and a correlation matrix using **sns.boxplot()**, **plt.scatter()**, and **sns.heatmap()** functions respectively.

**Step 3: Feature Engineering and Selection**

In this section, we create new features such as **distance**, **pickup_hour**, and **pickup_day** based on the existing variables. We then perform feature selection by selecting the variables that are most strongly correlated with **trip_duration**.

For this project, we can use a regression model to predict the **trip_duration** of taxi trips based on the selected features. One suitable algorithm for this task is the Random Forest Regressor, which is a versatile algorithm that can handle non-linear relationships and interactions between variables.

To run this code, you will need to install the necessary libraries such as Pandas, NumPy, Matplotlib, and Seaborn. You can then run the code in a Python environment such as Jupyter Notebook or PyCharm. The dataset can be downloaded from the Kaggle website and saved in the working directory as **train.csv**.

# Future Work :

**Step 1: Additional Data Collection and Preprocessing**

In this step, we can collect additional data such as weather data, traffic data, and events data (e.g., concerts, sports games, etc.) to incorporate into our analysis. We can preprocess the data using similar techniques as in the original dataset, such as removing missing values and outliers.

**Step 2: Feature Engineering and Selection**

We can further engineer new features based on the additional data collected, such as weather conditions, traffic congestion levels, and event schedules. We can then perform feature selection to select the most important features for our prediction model.

**Step 3: Model Selection and Tuning**

In this step, we can try different regression models such as Linear Regression, Gradient Boosting, and Neural Networks to compare their performance with the Random Forest Regressor used in the original project. We can also tune the hyperparameters of the models using techniques such as Grid Search or Random Search to optimize their performance.

**Step 4: Model Evaluation and Interpretation**

Once we have trained and tuned our models, we can evaluate their performance using metrics such as Mean Squared Error (MSE) and Root Mean Squared Error (RMSE). We can also interpret the models using techniques such as feature importance analysis and partial dependence plots to understand how the selected features contribute to the predicted **trip_duration**.

**Step 5: Deployment and Monitoring**

Finally, we can deploy the model as a web service using a framework such as Flask or Django. We can also monitor the performance of the model in real-time using tools such as Prometheus and Grafana to ensure that it continues to perform well over time.

To implement these steps, we can follow similar steps as in the original project and update the code accordingly. We can also use version control tools such as Git to manage changes and collaborate with others.

# Exercise Questions :

**1. Can you explain how you handled missing values in the dataset?**

To handle missing values, I used the **dropna()** function in Pandas to remove any rows with missing values. This is a common approach when the percentage of missing values is relatively small, as removing them is unlikely to significantly affect the overall analysis. Another approach is to impute the missing values using techniques such as mean imputation or regression imputation.

**2. How did you select the features for your prediction model?**

I selected the features using a combination of domain knowledge and statistical analysis. I first created new features based on the existing variables, such as the **distance** feature which I calculated using the Haversine formula. I then used correlation analysis to select the features that were most strongly correlated with the target variable **trip_duration**. This helped me to identify features such as **pickup_hour** and **pickup_day** which were found to have a relatively strong correlation with **trip_duration**.

**3. Can you explain how the Random Forest Regressor works?**

The Random Forest Regressor is an ensemble learning algorithm that works by building multiple decision trees and aggregating their predictions. Each decision tree is trained on a random subset of the training data and a random subset of the features, which helps to reduce overfitting and improve generalization. The final prediction is then made by averaging the predictions of all the individual decision trees.

**4. How did you evaluate the performance of your prediction model?**

I evaluated the performance of my prediction model using the Mean Squared Error (MSE) and Root Mean Squared Error (RMSE) metrics. The MSE measures the average squared difference between the predicted values and the actual values, while the RMSE is the square root of the MSE and provides a measure of the average deviation of the predicted values from the actual values. I also used a scatter plot to visually inspect the predictions and see how well they aligned with the actual values.

**5. How would you improve the performance of the prediction model?**

To improve the performance of the prediction model, I would consider several approaches. One approach is to collect additional data, such as weather data or traffic

data, which could help to capture additional factors that affect trip duration. Another approach is to try different regression algorithms or ensemble methods, such as Gradient Boosting or Stacked Regression, which may be better suited for capturing non-linear relationships between variables. Finally, I would also explore tuning the hyperparameters of the model using techniques such as Grid Search or Random Search to optimize its performance.

# Concept Explanation :

The algorithm we used is called the Random Forest Regressor. Now, you might be thinking, "What does a forest have to do with predicting taxi trip duration?" Well, let me explain!

Imagine you're lost in a forest and trying to find your way out. You come across a clearing with five different paths leading in different directions. You're not sure which path to take, so you decide to ask five different forest rangers for their opinion.

Each ranger has their own way of navigating the forest, with their own set of experiences and knowledge. One ranger tells you to take the path with the most sunlight, while another suggests following the river. A third ranger recommends following the bird calls, while a fourth suggests following the moss on the trees. The fifth ranger tells you to just trust your instincts and go with your gut.

You're not sure which ranger to believe, so you decide to take the average of all their suggestions and follow the path that seems to make the most sense based on their combined opinions.

This is essentially how the Random Forest Regressor works. Instead of forest rangers, we have decision trees - models that make predictions based on a series of if-then statements. Each decision tree is like a different ranger, with its own set of experiences and knowledge. The Random Forest combines the predictions of many different decision trees to make a final prediction.

For example, let's say we want to predict the price of a house based on its square footage, number of bedrooms, and neighborhood. We build multiple decision trees, each one trained on a different subset of the data and using a different subset of the features. Each tree makes its own prediction, based on its own set of if-then statements. The final prediction is then made by taking the average of all the individual tree predictions.

Just like how the average of the forest ranger opinions can help guide you out of the forest, the average of the decision tree predictions can help guide us in making accurate predictions about house prices (or in our case, taxi trip durations).

I hope that helps explain the Random Forest Regressor in a fun and easy-to-understand way!