

Executive Summary

Project Overview:

Automatidata has undertaken a project with the New York City Taxi and Limousine Commission (NYC TLC) to analyze taxi trip data. The primary objective is to prepare the dataset for detailed analysis, including exploratory data analysis (EDA), visualization, and statistical testing. This summary outlines the findings from the initial data inspection and provides recommendations for further analysis.

Data Inspection and Summary:

1. Data Structure:

- The dataset consists of 18 columns, including integer, float, and object data types.
- Columns include variables such as trip_distance, fare_amount, total_amount, and datetime fields like tpep_pickup_datetime and tpep_dropoff_datetime.

2. Null Values and Data Types:

- No null values are present across all columns.
- Variables are of various types: integers (e.g., VendorID, passenger_count), floats (e.g., trip_distance, fare_amount), and objects (e.g., tpep_pickup_datetime, store_and_fwd_flag).

3. Data Distribution:

- **Trip Distance:** Ranges from 1.00 to 33.96 miles. Higher distances are plausible but may include some extreme values warranting further investigation.
- **Fare Amount:** Ranges from \$6.50 to \$1200.29. The highest fare appears unusually high and might be an outlier or error.
- **Total Amount:** Includes fare, tip, and additional charges. Extreme values should be validated against trip details.

4. Variable Analysis:

- **trip_distance** and **total_amount**: High values in these variables may indicate outliers or data entry errors. Normal ranges are within expected values, though the maximum values for total_amount stand out.
- **payment_type**: Analysis reveals different average tip amounts based on the payment method, with credit card payments generally showing different tipping behavior compared to cash payments.
- **VendorID**: Shows distribution of services between different vendors, which could affect fare amounts and service quality.

Recommendations for Further Analysis:

1. Variable Focus:

- **trip_distance**: Essential for predicting fare amounts and assessing the impact of trip length on overall revenue.
- **passenger_count**: Provides insights into the scale of trips and their potential influence on total fare, especially when combined with trip distance.

2. Outlier Handling:

- Further investigation is needed for extreme values in trip_distance and total_amount to determine if they are valid or erroneous.

3. Datetime Validation:

- Ensure consistency between tpep_pickup_datetime and tpep_dropoff_datetime to confirm data accuracy.

4. Categorical Analysis:

- Examine payment_type and VendorID in more detail to understand their effects on tipping behavior and service distribution.

Conclusion:

The initial data inspection provides a solid foundation for deeper analysis. Key variables such as trip_distance and passenger_count are critical for building predictive models and understanding the dataset's dynamics. Addressing outliers and validating data consistency will enhance the quality and reliability of future analyses.