# Report: Optimising NYC Taxi Operations

Include your visualisations, analysis, results, insights, and outcomes. Explain your methodology and approach to the tasks. Add your conclusions to the sections.

## 1. Data Preparation

### 1.1. Loading the dataset

#### 1.1.1. Sample the data and combine the files
As suggested, we can begin the analysis using 250K to 300K records. However, applying a 5% sampling approach—by extracting 5% of data from each month—provides a more accurate and representative dataset. Using this method, the final sample consists of approximately 1.9 million records, which will support more precise prediction and analysis.

## 2. Data Cleaning

### 2.1. Fixing Columns

#### 2.1.1. Fix the index
The index was reset to align the Vendor ID column. Unwanted columns were removed here.

#### 2.1.2. Combine the two airport_fee columns
In the final sample dataset, two ambiguous columns were identified: **airport_fee** and **Airport_fee**. To resolve this, I applied a maximum-value approach—calculating the maximum value from both columns for each record. A new column, **airport_fee_combined**, was then created to store the consolidated values.

### 2.2. Handling Missing Values

#### 2.2.1. Find the proportion of missing values in each column

Highlighted columns in the below table contain blank values.

| Column Name | Missing % |
|---|---|
| VendorID | 0 |
| tpep_pickup_datetime | 0 |
| tpep_dropoff_datetime | 0 |
| **passenger_count** | **0.034255** |
| trip_distance | 0 |
| **RatecodeID** | **0.034255** |
| **store_and_fwd_flag** | **0.034255** |
| PULocationID | 0 |
| DOLocationID | 0 |
| payment_type | 0 |
| fare_amount | 0 |
| extra | 0 |
| mta_tax | 0 |
| tip_amount | 0 |
| tolls_amount | 0 |
| improvement_surcharge | 0 |
| total_amount | 0 |
| **congestion_surcharge** | **0.034255** |
| **airport_fee_combined** | **0.034255** |

### 2.2.2. Handling missing values in passenger_count

The **passenger_count** column contained a notable number of null values, which could directly influence key metrics such as average fare and other trip-based analyses. To address this, all null values were replaced with **1**, based on the reasonable assumption that any recorded trip would involve at least one passenger. This approach preserves the overall distribution of the data without introducing skewness.

### 2.2.3. Handle missing values in RatecodeID

To address the null or blank values in the **RatecodeID** column, all missing entries were replaced with **1**. The decision is based on the assumption that *Rate Code 1* represents the standard rate typically applied to most trips. By applying this default value, the analysis can proceed under the assumption that at least the standard rate was used, without altering the overall data profile or introducing bias.

**2.2.4.** **Impute NaN in congestion_surcharge**
Missing values in the **congestion_surcharge** column were addressed by replacing them with the **median** of the non-null values. Using the median helps avoid the influence of extreme outliers, ensuring the imputed values align with the central tendency of the data and maintaining the integrity of the column's distribution.

## 2.3. Handling Outliers and Standardising Values

**2.3.1.** **Check outliers in payment type, trip distance and tip amount columns**
**Payment Type:**
Outliers were identified where **payment_type** had a value of *0*, which is not a valid category. These invalid records were removed from the dataset.

**Trip Distance:**
Several types of outliers were detected in the **trip_distance** column:

- Trips with distances **< 0.1 miles** but fares **greater than $300** were considered inconsistent and removed.
- Trips with distances **> 250 miles** were treated as extreme outliers and excluded.
- Records showing **0 distance and 0 fare**, yet having different pickup and drop-off locations, were deemed invalid and removed from the dataset.

**Tip Amount:**
No filtering was applied for **tip_amount = 0**, as tipping is optional. However, extremely high tip values (outliers) were naturally minimized through **min–max standardization**, which scales all values between 0 and 1, reducing the influence of extreme tips on the analysis.
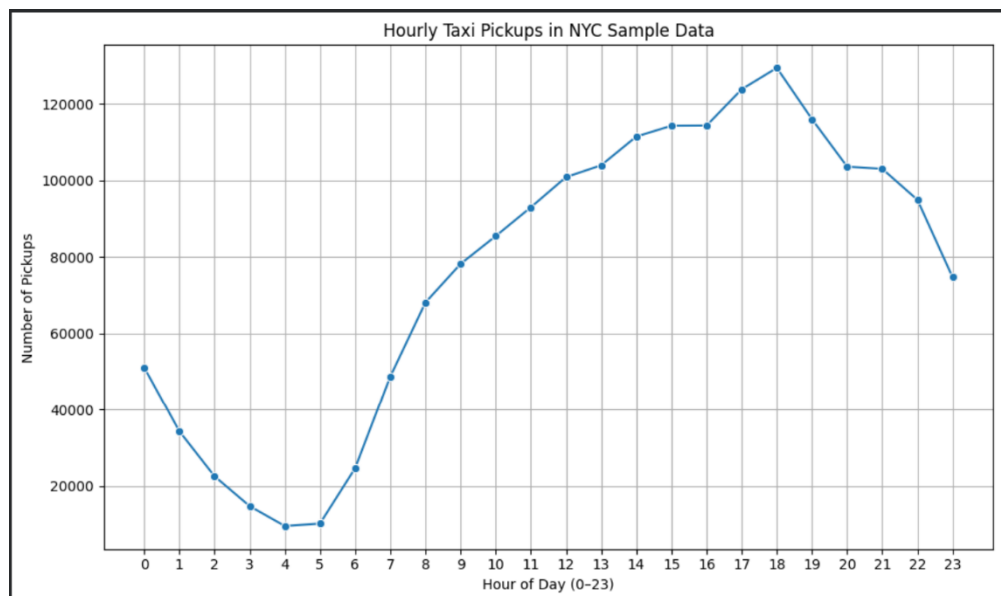
# 3. Exploratory Data Analysis

## 3.1. General EDA: Finding Patterns and Trends

**3.1.1.** **Classify variables into categorical and numerical**

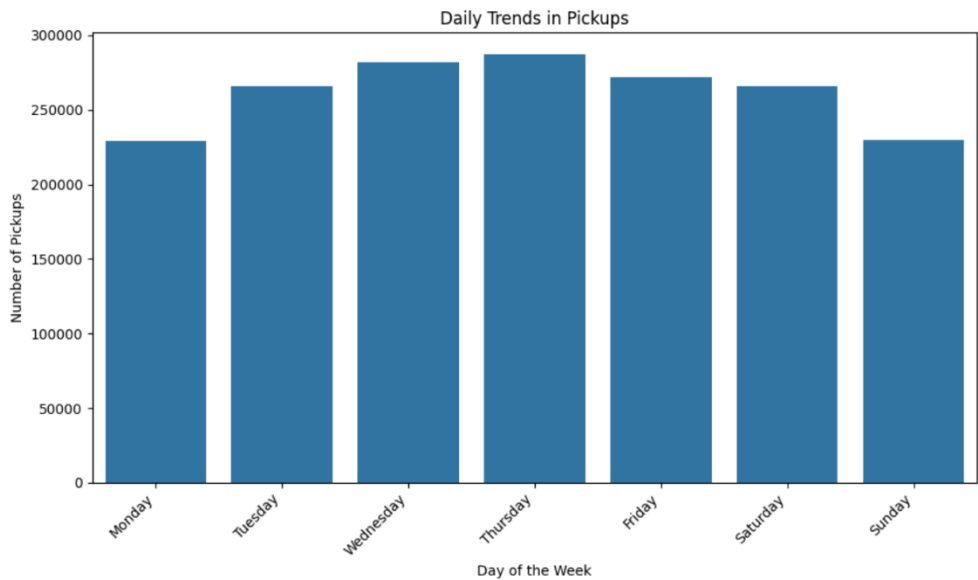| Variable Name | Type | Notes |
| --- | --- | --- |
| VendorID | Categorical | Coded identifier for vendor |
| tpep_pickup_datetime | Datetime | Timestamp of trip start |
| tpep_dropoff_datetime | Datetime | Timestamp of trip end |
| passenger_count | Numerical | Discrete count of passengers |
| trip_distance | Numerical | Continuous measurement (miles) |
| RatecodeID | Categorical | Coded rate type |
| PULocationID | Categorical | Pickup location code |
| DOLocationID | Categorical | Drop-off location code |
| payment_type | Categorical | Coded payment method |
| pickup_hour | Numerical | Derived from pickup timestamp |
| trip_duration | Numerical | Duration in minutes/seconds |
| fare_amount | Numerical | Monetary value |
| extra | Numerical | Monetary value |
| mta_tax | Numerical | Monetary value |
| tip_amount | Numerical | Monetary value |
| tolls_amount | Numerical | Monetary value |
| improvement_surcharge | Numerical | Monetary value |
| total_amount | Numerical | Monetary value |
| congestion_surcharge | Numerical | Monetary value |
| airport_fee | Numerical | Monetary value |

**3.1.2.  Analyse the distribution of taxi pickups by hours, days of the week, and months**
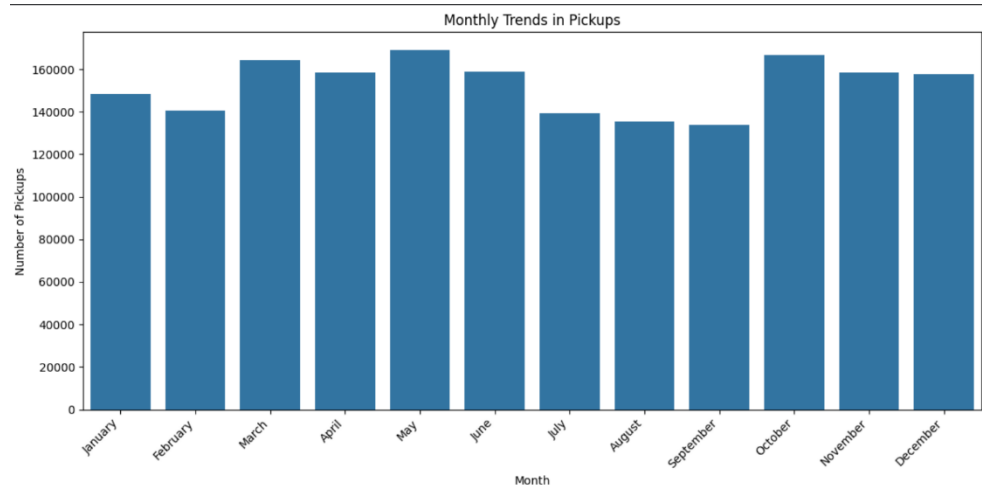
*Tax-Pickups by Hour*

The data shows a clear daily pattern, with minimal taxi demand before dawn, a sharp increase during the morning commute, and the highest activity occurring in the early evening. Demand then tapers off steadily into the night.

*Tax-Pickups by Day*



Taxi pickups steadily increase from Monday and peak mid-week (Wednesday–Thursday), reflecting higher weekday commuting activity. Demand slightly tapers on Friday and Saturday, followed by a noticeable drop on Sunday. Overall, weekdays show consistently higher pickup volumes compared to weekends.

*Tax-Pickups by Month*



Monthly Trends in Pickups

Taxi pickups fluctuate moderately throughout the year, with higher activity observed in spring and early summer (March–June). There is a dip during late summer and early fall, followed by another rise in October. Overall, seasonal variations are present but not extreme.
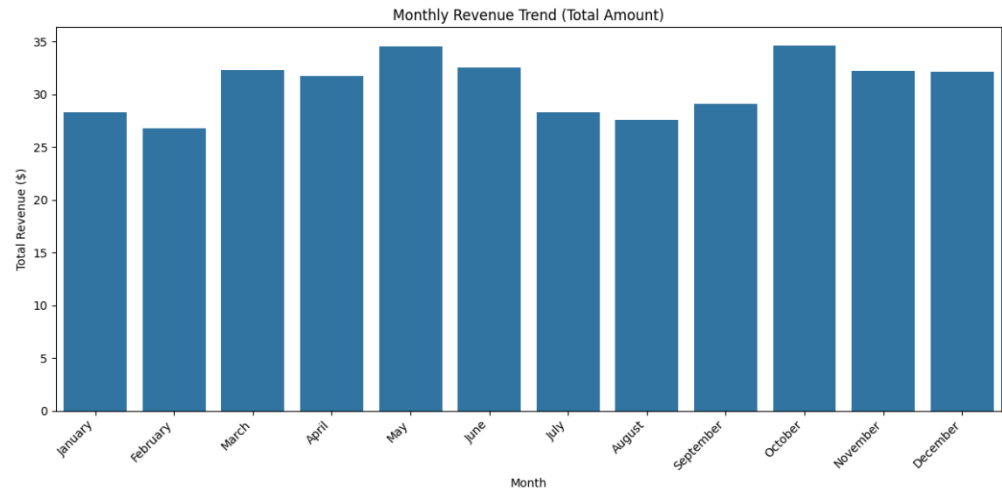
**3.1.3.    Filter out the zero/negative values in fares, distance and tips**
To ensure data quality, the dataset was filtered using the following criteria:

- Records with **fare_amount** or **total_amount** equal to zero were removed, as these typically indicate canceled or invalid trips.
- Trips with zero **trip_distance** despite having different pickup and drop-off locations were flagged as inconsistent and excluded.
- Zero **tip_amount** values were retained because tipping is optional, and many valid trips legitimately report no tip while still showing a correct total fare.

This approach cleans the dataset effectively while preserving real-world behaviors—such as passengers choosing not to tip.
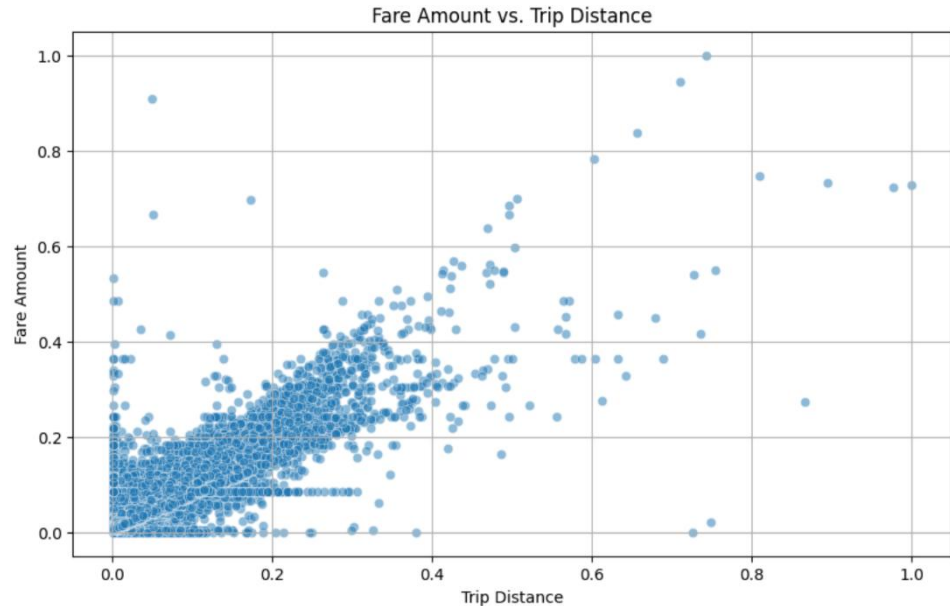
### 3.1.4. Analyse the monthly revenue trends



Revenue shows a steady increase from February through May, peaking during late spring and early summer. A dip occurs in July and August, followed by a strong rebound in October, the highest revenue month. Overall, revenue patterns align closely with seasonal travel and activity levels.

### 3.1.5. Find the proportion of each quarter's revenue in the yearly revenue

Quarterly revenue shows noticeable fluctuations, with 2023 Q2 and Q4 recording the highest average total amounts. Revenue dips in 2023 Q3, indicating a seasonal slowdown before rising again at year-end. Overall, 2023 maintains consistently stronger performance compared to late 2022.

| Quarter | Average Total Amount ($) |
|---------|--------------------------|
| 2022 Q4 | 0 |
| 2023 Q1 | 23.69 |
| 2023 Q2 | 26.75 |
| 2023 Q3 | 22.76 |
| 2023 Q4 | 26.81 |

### 3.1.6. Analyse and visualise the relationship between distance and fare amount
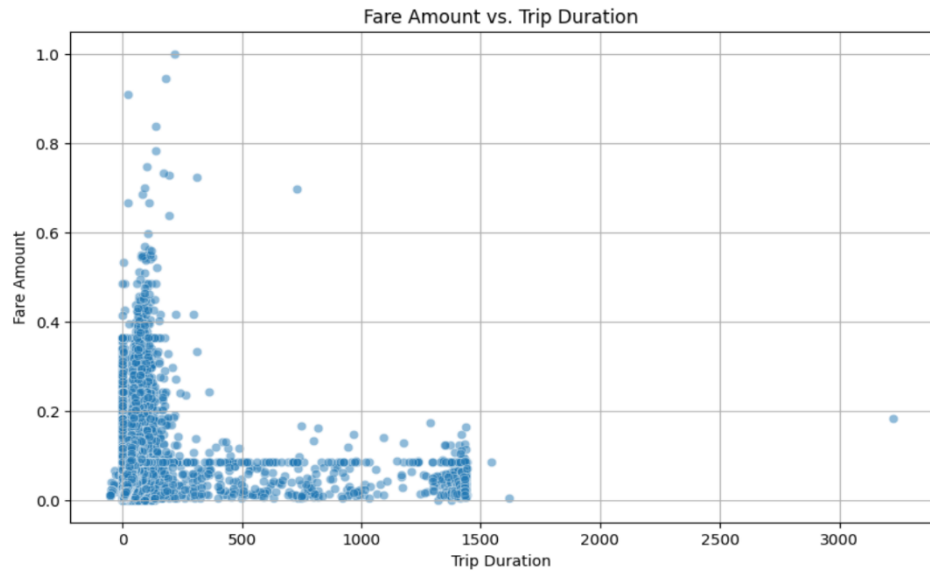
Fare Amount vs. Trip Distance

The scatter plot shows a clear upward trend, indicating that higher trip distances generally correspond to higher fare amounts. Most data points are concentrated at shorter distances, but even there, the positive relationship holds.

This visual pattern supports the **strong correlation (0.94)** between distance and fare.

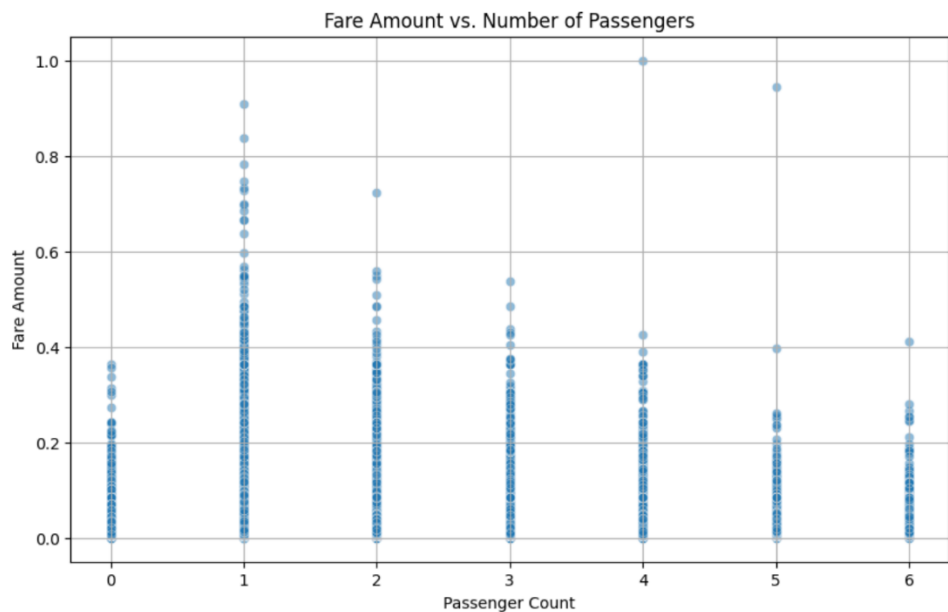### 3.1.7.    Analyse the relationship between fare/tips and trips/passengers

**Fare Amount Vs. Trip Duration**

The scatter plot shows a weak relationship between trip duration and fare amount, with most points scattered widely and no strong upward trend. The correlation value of **0.27** confirms only a mild positive association. This suggests that fare pricing is driven more by **distance** than by **duration**.

Fare Amount vs. Trip Duration

**Fare Amount vs. Number of Passengers**

The scatter plot shows no meaningful relationship between passenger count and fare amount, with fare values spread similarly across all passenger counts. The very low **correlation of 0.05** confirms that the number of passengers has almost no impact on fare calculations. This indicates fares are primarily distance-based rather than influenced by how many people share the ride.



Fare Amount vs. Number of Passengers

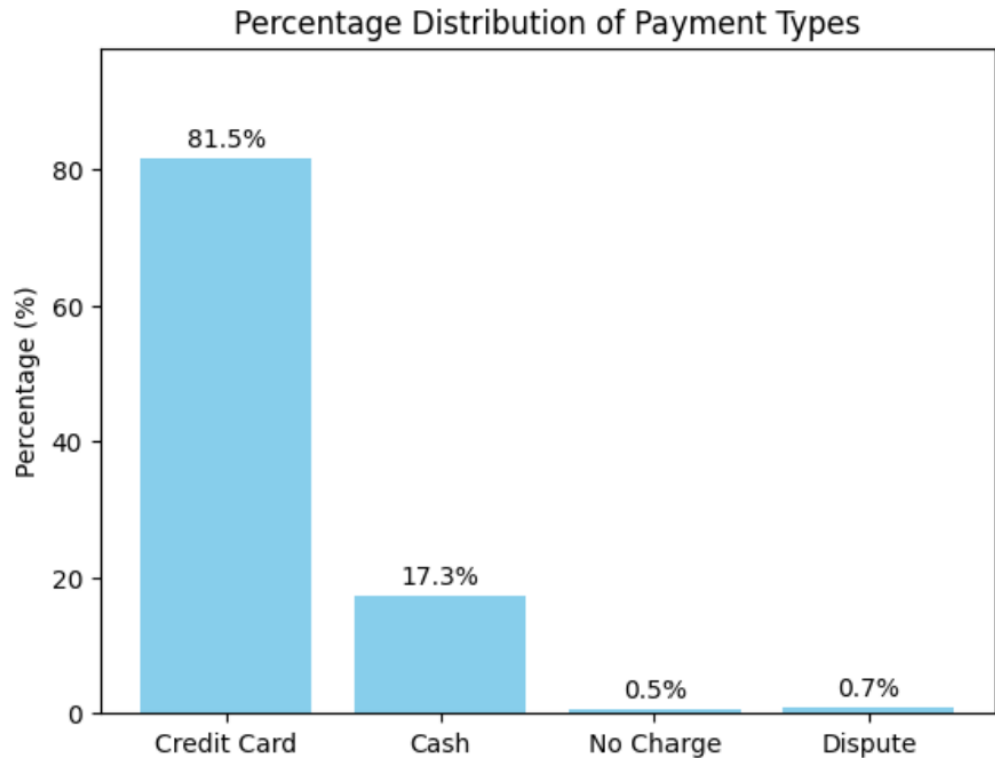**Tip Amount vs. Trip Distance**

The scatter plot shows a moderate upward trend, indicating that tip amounts tend to increase with trip distance, although the relationship is not very strong. The correlation of **0.57** confirms a moderate positive association. Overall, longer trips generally lead to slightly higher tips, but tip behavior still varies widely among riders.
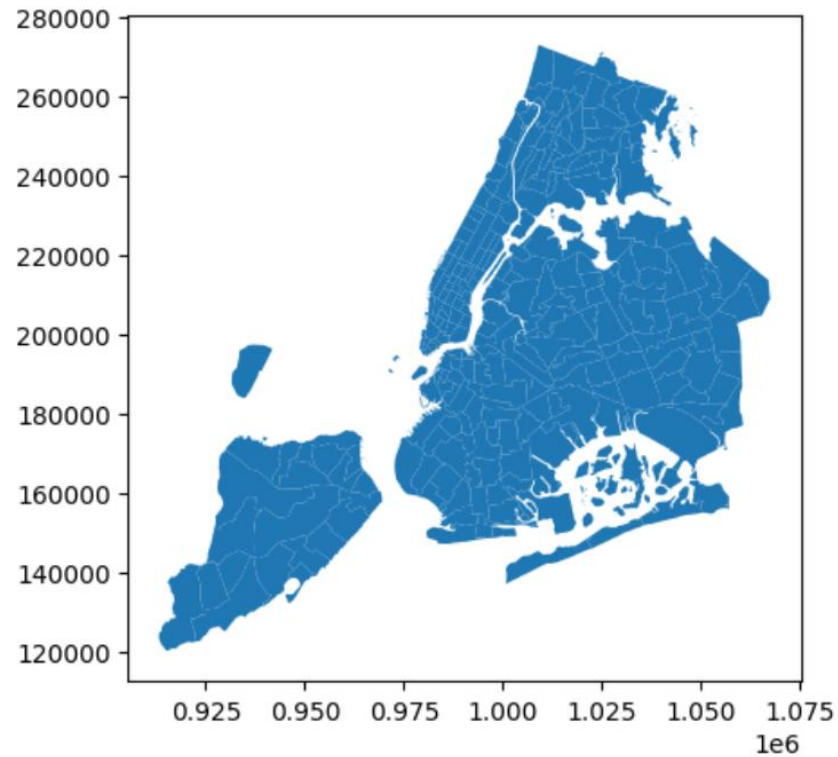


Tip Amount vs. Trip Distance

### 3.1.8. Analyse the distribution of different payment types



Percentage Distribution of Payment Types

Credit cards account for **81.5%** of all payments, showing a strong preference for digital transactions. Cash makes up **17.3%**, while **No Charge** and **Dispute** categories together contribute less than **1.5%**. Overall, the data highlights a predominantly cashless taxi ecosystem.

### 3.1.9. Load the taxi zones shapefile and display it

| | OBJECTID | Shape_Leng | Shape_Area | zone | LocationID | borough | geometry |
|---|---|---|---|---|---|---|---|
| 0 | 1 | 0.116357 | 0.000782 | Newark Airport | 1 | EWR | POLYGON ((933100.918 192536.086, 933091.011 19... |
| 1 | 2 | 0.433470 | 0.004866 | Jamaica Bay | 2 | Queens | MULTIPOLYGON (((1033269.244 172126.008, 103343... |
| 2 | 3 | 0.084341 | 0.000314 | Allerton/Pelham Gardens | 3 | Bronx | POLYGON ((1026308.77 256767.698, 1026495.593 2... |
| 3 | 4 | 0.043567 | 0.000112 | Alphabet City | 4 | Manhattan | POLYGON ((992073.467 203714.076, 992068.667 20... |
| 4 | 5 | 0.092146 | 0.000498 | Arden Heights | 5 | Staten Island | POLYGON ((935843.31 144283.336, 936046.565 144... |

### 3.1.10.  Merge the zone data with trips data

A merge operation was performed to combine the zones dataset with the trip dataset by matching the **LocationID** field from the zones data with the **PULocationID** field in the trip data.
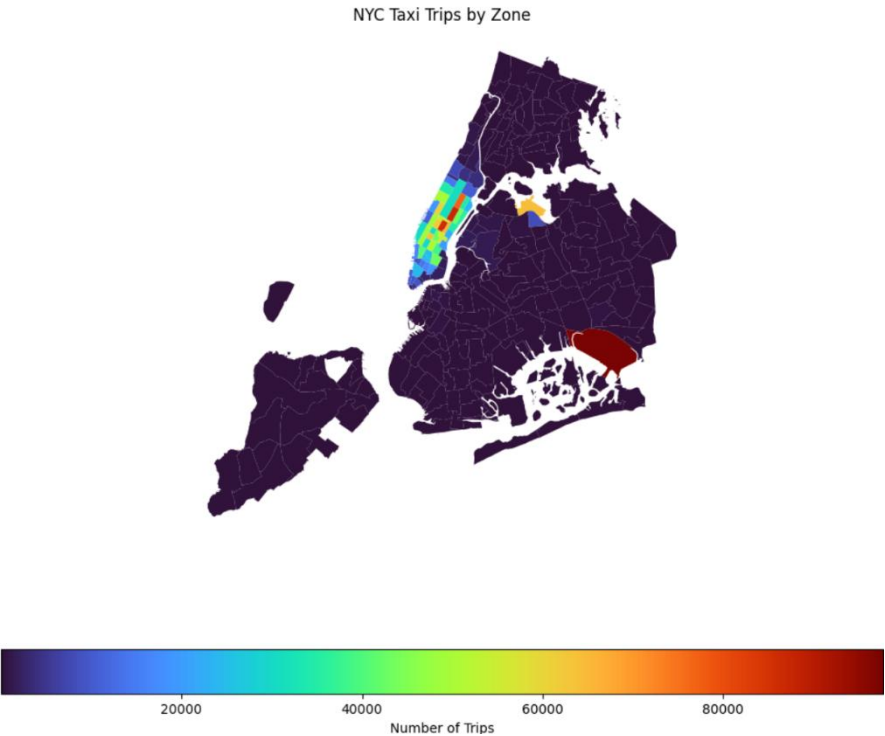
### 3.1.11.  Find the number of trips for each zone/location ID

|   | PULocationID | TripCount |
|---|---|---|
| **0** | 1 | 234 |
| **1** | 2 | 2 |
| **2** | 3 | 39 |
| **3** | 4 | 1836 |
| **4** | 5 | 25 |

### 3.1.12.   Add the number of trips for each zone to the zones dataframe

| | OBJECTID | Shape_Leng | Shape_Area | zone | LocationID | borough | geometry | PULocationID | TripCount |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0.116357 | 0.000782 | Newark Airport | 1 | EWR | POLYGON ((933100.918 192536.086, 933091.011 19... | 1.0 | 234.0 |
| 1 | 2 | 0.433470 | 0.004866 | Jamaica Bay | 2 | Queens | MULTIPOLYGON (((1033269.244 172126.008, 103343... | 2.0 | 2.0 |
| 2 | 3 | 0.084341 | 0.000314 | Allerton/Pelham Gardens | 3 | Bronx | POLYGON ((1026308.77 256767.698, 1026495.593 2... | 3.0 | 39.0 |
| 3 | 4 | 0.043567 | 0.000112 | Alphabet City | 4 | Manhattan | POLYGON ((992073.467 203714.076, 992068.667 20... | 4.0 | 1836.0 |
| 4 | 5 | 0.092146 | 0.000498 | Arden Heights | 5 | Staten Island | POLYGON ((935843.31 144283.336, 936046.565 144... | 5.0 | 25.0 |

### 3.1.13.   Plot a map of the zones showing number of trips



JFK Airport emerges as the most prominent hotspot, highlighted in deep red, while central Manhattan shows moderate activity in green and yellow tones. The majority of outer-borough zones appear in darker shades, reflecting comparatively low taxi demand.

### 3.1.14.   Conclude with results

- Rigorous data cleaning and standardization ensured a reliable foundation for analysis by removing inconsistencies, correcting numeric fields, and resolving missing values.
- Travel patterns reveal that **airport zones and Midtown Manhattan consistently attract the highest trip volumes**, highlighting their role as major transportation hubs.

- **Weekday commuting peaks** dominate the morning and evening hours, while **weekends shift activity toward late-night travel**, reflecting different rider behavior.
- The data shows that **most rides involve just one or two passengers**, and **credit card transactions overwhelmingly lead payment choices**, underscoring the system's shift toward cashless mobility.
- A clear **seasonal pattern** emerges, with **Q3 registering the highest trip activity**, likely influenced by tourism and summer travel.
- A strong positive relationship between **distance and fare amounts** confirms that trip length remains the primary driver of pricing.

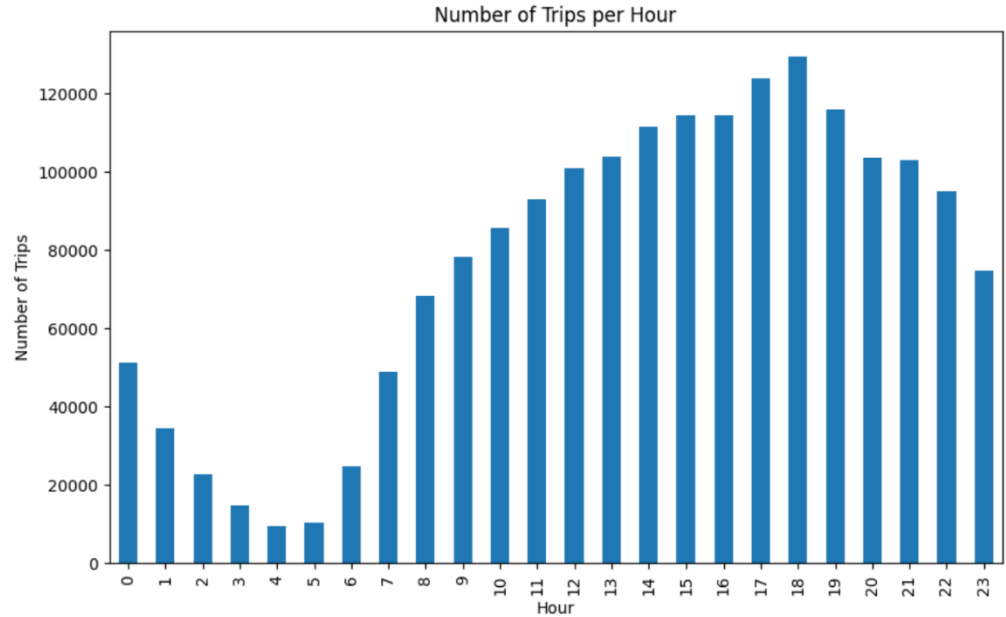## 3.2. Detailed EDA: Insights and Strategies

### 3.2.1. Identify slow routes by comparing average speeds on different routes

These records show trips with extremely low speeds, often caused by very short distances or heavy congestion. Many occur within the same pickup and drop-off zones, indicating minimal vehicle movement during the trip.

| PULocationID | DOLocationID | Hour | Trip Speed |
|---|---|---|---|
| 223 | 223 | 4 | 0.000006 |
| 194 | 194 | 15 | 0.000016 |
| 193 | 264 | 10 | 0.000119 |
| 215 | 215 | 16 | 0.000123 |
| 80 | 68 | 17 | 0.000218 |
| 237 | 264 | 2 | 0.000297 |
| 216 | 181 | 9 | 0.000358 |
| 235 | 181 | 9 | 0.000381 |
| 135 | 47 | 13 | 0.000408 |
| 51 | 48 | 9 | 0.000413 |

### 3.2.2. Calculate the hourly number of trips and identify the busy hours

The busiest hour for taxi trips is **6 PM,** with **approximately 129,425 rides**, reflecting peak evening commute demand. Trip volumes steadily build throughout the day and reach their highest point during this rush-hour window.

Number of Trips per Hour

### 3.2.3. Scale up the number of trips from above to find the actual number of trips

**Top 5 Busiest Hour**

| Hour | Trip Count |
|------|------------|
| 18 | 2,588,500 |
| 17 | 2,476,800 |
| 19 | 2,320,340 |
| 16 | 2,287,060 |
| 15 | 2,286,020 |

### 3.2.4. Compare hourly traffic on weekdays and weekends

Weekday traffic shows clear morning and evening peaks, driven by commuter activity, with the highest demand occurring around 5–6 PM. In contrast, weekend traffic remains much steadier throughout the day, with no sharp spikes and overall lower trip volumes. This reflects a shift from work-driven travel on weekdays to more flexible, leisure-oriented movement on weekends.

Hourly Traffic Patterns: Weekdays vs. Weekends

### 3.2.5. Identify the top 10 zones with high hourly pickups and drops

```
Top 10 Pickup Zones
   LocationID  Total_Trips                          zone
0         132        97726                   JFK Airport
1         237        87136         Upper East Side South
2         161        85710                 Midtown Center
3         236        76747         Upper East Side North
4         162        66033                  Midtown East
5         138        64277             LaGuardia Airport
6         186        64098    Penn Station/Madison Sq West
7         230        61498       Times Sq/Theatre District
8         142        60079            Lincoln Square East
9         170        54697                   Murray Hill
Top 10 Dropoff Zones
   LocationID  Total_Trips                          zone
0         236        81277         Upper East Side North
1         237        77782         Upper East Side South
2         161        71446                 Midtown Center
3         230        56044       Times Sq/Theatre District
4         170        54134                   Murray Hill
5         162        51983                  Midtown East
6         142        51462            Lincoln Square East
7         239        51427         Upper West Side South
8         141        48264              Lenox Hill West
9          68        46350                  East Chelsea
```

### 3.2.6. Find the ratio of pickups and dropoffs in each zone

| | LocationID | Trip_Counts | Pickup_Zone |
|---|---|---|---|
| 0 | 79 | 15441 | East Village |
| 1 | 132 | 14776 | JFK Airport |
| 2 | 249 | 12650 | West Village |
| 3 | 48 | 10431 | Clinton East |
| 4 | 148 | 9648 | Lower East Side |
| 5 | 114 | 8942 | Greenwich Village South |
| 6 | 230 | 8177 | Times Sq/Theatre District |
| 7 | 186 | 7071 | Penn Station/Madison Sq West |
| 8 | 164 | 6190 | Midtown South |
| 9 | 138 | 6096 | LaGuardia Airport |

| zone | ratio |
|---|---|
| ...urst | 8.051823 |
| ...port | 4.552383 |
| ...port | 2.879407 |
| ...Bay | 2.000000 |
| ...Vest | 1.588983 |
| ...outh | 1.394740 |
| ...Park | 1.372204 |
| ...lage | 1.336353 |
| ...East | 1.270281 |
| ...strict | 1.210277 |

**Pickup Zones**

**Top 10 Drop off Zones**

### 3.2.7. Identify the top zones with high traffic during night hours

**Top 10 Pickup Zones at Night**

| | LocationID | Pickup_Trips | Dropoff_Trips | zone | ratio |
|---|---|---|---|---|---|
| 253 | 58 | 2 | 51.0 | Country Club | 0.039216 |
| 251 | 221 | 2 | 49.0 | Stapleton | 0.040816 |
| 92 | 1 | 234 | 5186.0 | Newark Airport | 0.045121 |
| 196 | 257 | 39 | 748.0 | Windsor Terrace | 0.052139 |
| 95 | 112 | 217 | 4026.0 | Greenpoint | 0.053900 |
| 226 | 16 | 17 | 314.0 | Bayside | 0.054140 |
| 176 | 198 | 55 | 1009.0 | Ridgewood | 0.054509 |
| 228 | 64 | 13 | 230.0 | Douglaston | 0.056522 |
| 254 | 176 | 1 | 17.0 | Oakwood | 0.058824 |
| 250 | 27 | 2 | 33.0 | Breezy Point/Fort Tilden/Riis Beach | 0.060606 |

**Top 10 Dropoff Zones at Night**

| | LocationID | Trip_Counts | Dropoff_Zone |
|---|---|---|---|
| 0 | 79 | 8430 | East Village |
| 1 | 48 | 6970 | Clinton East |
| 2 | 170 | 6303 | Murray Hill |
| 3 | 107 | 5796 | Gramercy |
| 4 | 68 | 5693 | East Chelsea |
| 5 | 141 | 5193 | Lenox Hill West |
| 6 | 263 | 5084 | Yorkville West |
| 7 | 249 | 5072 | West Village |
| 8 | 230 | 4630 | Times Sq/Theatre District |
| 9 | 90 | 4426 | Flatiron |

**3.2.8.** **Find the revenue share for nighttime and daytime hours**

| | |
|---|---|
| Night-time revenue Share | 12.12% |
| Daytime revenue Share | 87.88% |

**3.2.9.** **For the different passenger counts, find the average fare per mile per passenger**

| | passenger_count | trip_distance | total_amount | avg_fare_per_mile | avg_fare_per_mile_per_pax |
|---|---|---|---|---|---|
| 0 | 1.0 | 27062.921320 | 45854.534037 | 0.590191 | 0.590191 |
| 1 | 2.0 | 6698.462588 | 10471.182964 | 0.639704 | 0.319852 |
| 2 | 3.0 | 1594.012817 | 2552.829972 | 0.624410 | 0.208137 |
| 3 | 4.0 | 962.307660 | 1504.093358 | 0.639793 | 0.159948 |
| 4 | 5.0 | 485.788287 | 811.195864 | 0.598854 | 0.119771 |
| 5 | 6.0 | 313.773334 | 526.314686 | 0.596171 | 0.099362 |

**3.2.10.** **Find the average fare per mile by hours of the day and by days of the week**

| | hour | fare_per_mile |
|---|---|---|
| 0 | 0 | 0.556726 |
| 1 | 1 | 0.549024 |
| 2 | 2 | 0.546853 |
| 3 | 3 | 0.574301 |
| 4 | 4 | 0.631854 |
| 5 | 5 | 0.649081 |
| 6 | 6 | 0.614117 |
| 7 | 7 | 0.540669 |
| 8 | 8 | 0.492234 |
| 9 | 9 | 0.482615 |
| 10 | 10 | 0.478502 |
| 11 | 11 | 0.465453 |
| 12 | 12 | 0.463898 |
| 13 | 13 | 0.470036 |
| 14 | 14 | 0.474274 |
| 15 | 15 | 0.469299 |
| 16 | 16 | 0.443713 |
| 17 | 17 | 0.432011 |
| 18 | 18 | 0.429186 |
| 19 | 19 | 0.450692 |
| 20 | 20 | 0.495361 |
| 21 | 21 | 0.510383 |
| 22 | 22 | 0.526080 |
| 23 | 23 | 0.552887 |

| | pickup_day | fare_per_mile |
|---|---|---|
| 0 | Friday | 0.475572 |
| 1 | Monday | 0.502020 |
| 2 | Saturday | 0.494441 |
| 3 | Sunday | 0.536380 |
| 4 | Thursday | 0.463126 |
| 5 | Tuesday | 0.468968 |
| 6 | Wednesday | 0.462874 |

### 3.2.11. Analyse the average fare per mile for the different vendors

Average fare per mile stays fairly consistent across vendors, though Vendor 1 generally charges slightly higher rates during most hours. Both vendors show early-morning peaks around 4–6 AM, followed by lower, more stable fares throughout the day. Late evenings see a gradual rise again, likely due to reduced traffic and longer-distance trips.

Average Fare per Mile by Vendor

### 3.2.12. Compare the fare rates of different vendors in a distance-tiered fashion



Average Fare per Mile by Vendor and Distance Tier

### 3.2.13. Analyse the tip percentages

**Time period with lowest tip %**

|   | Time_period | tip_percentage |
|---|---|---|
| 1 | Early Morning | 30.319907 |
| 2 | Evening | 30.487970 |
| 0 | Afternoon | 30.874937 |
| 3 | Morning | 31.048132 |

**Distance Tier with lowest tip %**

| | Distance_tier | tip_percentage |
|---|---|---|
| 0 | 0 to 2 Miles | 30.705468 |

**Pax segment with lowest tip %**

| | passenger_count | tip_percentage |
|---|---|---|
| 0 | 1.0 | 30.315061 |
| 1 | 2.0 | 30.668905 |
| 2 | 3.0 | 30.725878 |
| 4 | 5.0 | 31.038703 |
| 5 | 6.0 | 31.043807 |

### 3.2.14. Analyse the trends in passenger count

**Hourly Trend**

Passenger volumes are lowest during the early morning hours but rise sharply after 6 AM, peaking in the late afternoon and early evening. The highest passenger activity occurs around 6–7 PM, aligning with the evening commute. After this peak, counts gradually decline through the night.
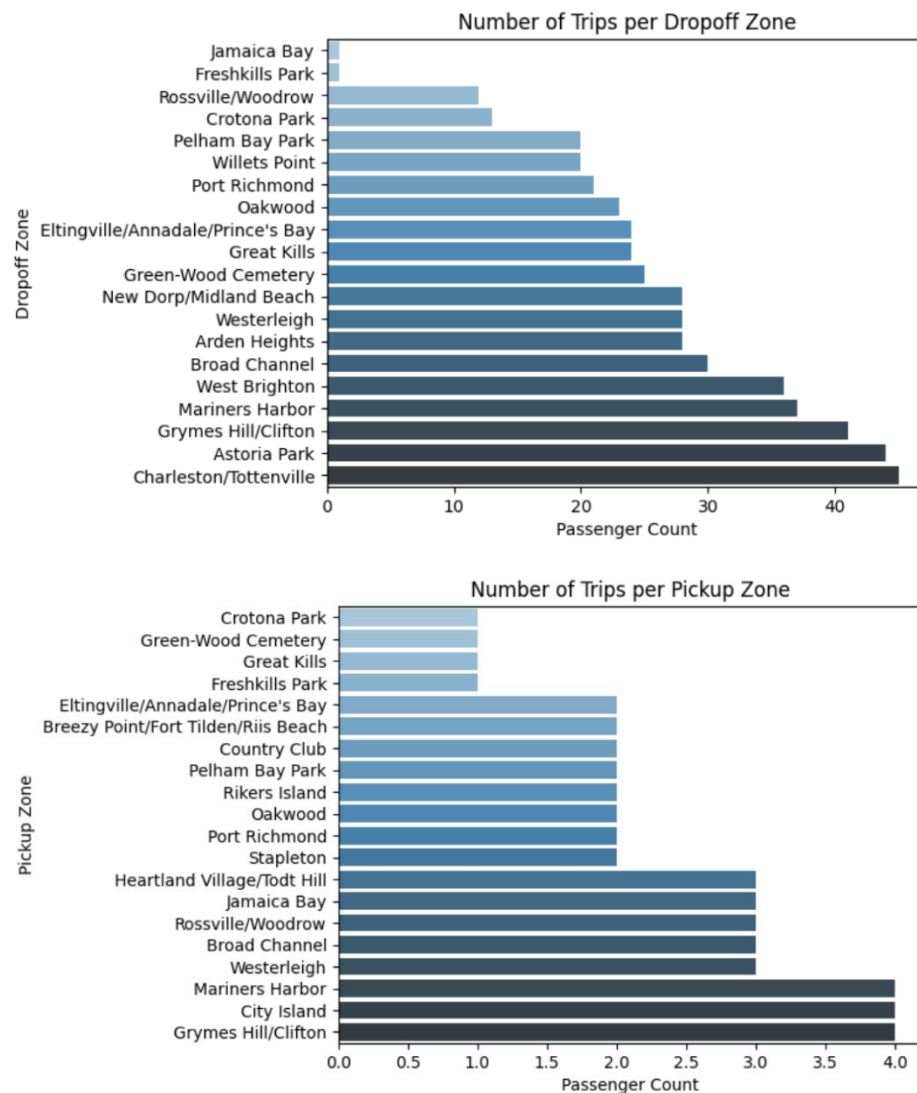
**Day wise Trend**

Average passengers per trip remain fairly stable across days, but weekends—especially Saturday and Sunday—consistently show higher passenger counts. Early mornings across all days have the lowest averages, while afternoon and evening hours see the highest. Overall, weekends tend to carry slightly more shared or group rides compared to weekdays



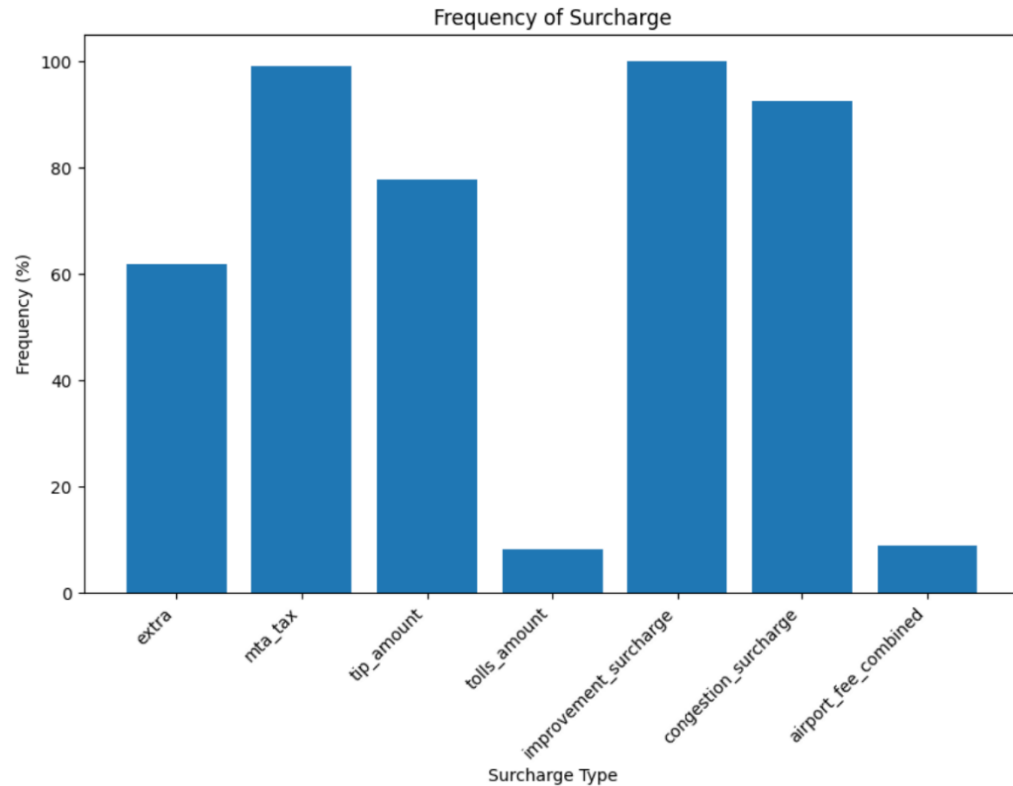Average Passenger Count by Hour and Day of the Week

### 3.2.15. Analyse the variation of passenger counts across zones

The pickup and dropoff charts both highlight low-activity zones, but notable differences appear in their rankings. Pickup activity is more evenly distributed across zones, with places like Grymes Hill/Clifton, City Island, and Mariners Harbor showing relatively higher counts. In contrast, dropoff patterns are more concentrated, with zones such as Charleston/Tottenville, Astoria Park, and Grymes Hill/Clifton receiving significantly more dropoffs compared to others.

## Number of Trips per Dropoff Zone



## Number of Trips per Pickup Zone



**3.2.16.** **Analyse the pickup/dropoff zones or times when extra charges are applied more frequently.**

The chart shows that **mta_tax** and **improvement_surcharge** appear in nearly all trips, indicating they are mandatory fees. **Congestion surcharge** is also highly frequent, reflecting heavy use of high-traffic zones. In contrast, **tolls** and **airport fees** occur far less often, suggesting they apply only to specific routes or locations.

Frequency of Surcharge

# 4. Conclusions

## 4.1. Final Insights and Recommendations

### 4.1.1. Recommendations to optimize routing and dispatching based on demand patterns and operational inefficiencies.

**Key Takeaways**

- **When people travel:** Demand shoots up during office rush hours, weekends, and a few busy months of the year. Areas with nightlife stay active late into the night.
- **How money flows:** Fares mostly depend on how far and how long the trip is. Shared rides often come with discounts. Tip amounts change based on the type of trip and customer experience.
- **Where people go**: Airports, major hubs, and tourist spots attract the most rides. Some areas get more pickups than drop-offs and vice-versa. Nightlife areas turn into hotspots after dark.
- **Vendors and extra charges**: Different vendors follow different pricing patterns. Some surcharges appear often, and long-distance trips may follow tiered prices.

**Optimization Recommendations**

- **Demand Strategy**
  - Prioritize coverage in high-demand areas and peak timings.
  - Strengthen late-night operations in nightlife corridors.
  - Develop tailored offerings for groups and shared-ride customers.
- **Supply Strategy**
  - Increase fleet presence in high-demand zones during peak hours.
  - Implement dynamic pricing aligned with demand and trip attributes.
  - Guide drivers to reposition strategically to balance supply.
  - Offer incentives for servicing low-demand periods or underserved locations.
- **Customer Experience**
  - Elevate service quality through training, monitoring, and standards.
  - Provide multiple payment options for convenience.
  - Promote ride-sharing as a cost-effective and eco-friendly choice.
- **Continuous Improvement**
  - Continuously track performance and refine strategies using analytics and customer feedback.
  - Work with city authorities to address operational and regulatory challenges.

**Conclusion**

By leveraging data-driven insights into demand, supply, and customer behavior, NYC taxi operators can enhance efficiency, meet rider expectations, and deliver a consistently positive travel experience. Through smarter planning and continuous improvement, the city's taxi ecosystem can thrive and adapt to changing transportation needs.

**4.1.2.** **Suggestions on strategically positioning cabs across different zones to make best use of insights uncovered by analysing trip trends across time, days and months.**

**Strategic Cab Positioning**

- **Time-Based Deployment**: Adjust cab availability according to daily patterns—morning and evening rush hours, late-night surges, quieter midday periods, and monthly demand shifts.

- **Day-Specific Approach**: Target business hubs on weekdays and shift focus to entertainment and residential areas on weekends. Stay flexible for festivals, events, and city gatherings.
- **Zone-Focused Positioning**: Prioritize zones with consistently high demand, correct pickup–dropoff imbalances, and maintain strong coverage in nightlife hotspots after dark.
- **Data-Led Decisions**: Utilize real-time data, predictive algorithms, and ride-hailing platforms to guide dynamic cab placement.
- **Collaborative Operations**: Maintain continuous driver communication and coordinate with city authorities to streamline traffic flow and enhance service reliability.
- **Tech-Powered Optimization**: Use GPS tracking, heatmaps, and analytics dashboards to gain actionable insights and refine positioning strategies.

By integrating these smart positioning methods, taxi operators can better match supply with demand, reduce passenger wait times, and improve overall service efficiency across NYC.

**4.1.3. Propose data-driven adjustments to the pricing strategy to maximize revenue while maintaining competitive rates with other vendors.**

**Data-Driven Pricing Adjustments**
- **Dynamic Fare Management**: Modify fares in real time based on demand, driver availability, and traffic patterns. Increase prices during peak periods and introduce incentives or discounts during slower hours.
- **Tiered & Zoned Pricing:** Keep short-trip rates competitive while using tiered pricing for longer journeys. Introduce zone-based variations to better reflect regional demand and trip complexity.
- **Shared Ride Incentives:** Promote group travel and shared-ride options to boost occupancy, reduce costs for passengers, and improve overall fleet efficiency.
- **Smarter Surcharge Strategies**: Review surcharge patterns and apply peak surcharges only when justified. Ensure passengers receive clear, upfront communication to maintain trust.
- **Competitor Benchmarking:** Stay aware of market pricing trends, adjust rates strategically, and emphasize unique service benefits to justify premium pricing when appropriate.
- **Continuous Optimization:** Use ongoing data collection, A/B price testing, and performance tracking to refine pricing models and balance revenue growth with customer satisfaction.