
COMS4060A/7056A: Assignment 2

Puseletso Nakana : 2341162
Ntando Ngobese : 2112256
Nthabiseng Thema : 2012016

Department of Computer Science
University of Witwatersrand

1 Question 1:

The analysis focused on three key metrics: trip duration, distance traveled, and average speed. For each metric, the interquartile range (IQR) method was used to identify outliers. This method considers values below $Q1 - 1.5IQR$ or above $Q3 + 1.5IQR$ as outliers, where $Q1$ and $Q3$ are the first and third quartiles respectively. The code calculated these thresholds for each metric and identified data points falling outside these ranges. Finally, all unique outliers from the three metrics were combined and removed from the dataset, resulting in a cleaned version of the data.

Number of outliers in trip duration: 17512

Number of outliers in distance: 33017

Number of outliers in speed: 14002

Number of unique outliers: 51322

Trip duration outliers (17,512): Extremely short or long trip durations could indicate:

Justification: Including these outliers could skew average trip times and impact analysis of typical travel patterns.

Distance outliers (33,017): Unusually short or long distances might represent: GPS errors

Incomplete trips

Justification: These outliers could distort average trip lengths and affect route analysis or fuel consumption estimates.

Justification for Removing Outliers Data Quality: Outliers often result from data entry errors or malfunctions in the taxi meters. Removing them helps in ensuring the quality and accuracy of the data.

Analysis Accuracy: Outliers can skew the results of statistical analyses and models. By removing them, you can obtain more accurate and reliable insights.

Realistic Insights: Outliers can distort the interpretation of patterns and trends. Removing them ensures that the insights you gain are more representative of typical trips.

33 2 Question 2

34 In this feature generation step, three new columns were added to the DataFrame:

35 Day of week: Extracted from the 'pickup_datetime' column using the dt.day_name() method. This
36 provides information about which day of the week each trip occurred.

37 Average speed: Calculated by dividing the 'distance_km' by the trip duration (converted from
38 seconds to hours). This gives the average speed of each trip in km/h.

39

	distance_km	day_of_week	speed_kmh
0	1.497580	Monday	11.848984
1	1.804374	Sunday	9.797504
2	6.381090	Tuesday	10.815406
3	1.484566	Wednesday	12.457894
4	1.187842	Saturday	9.830418

40

41 3 Question 3

42 3.1 3.1

43 Day of Week Extraction:

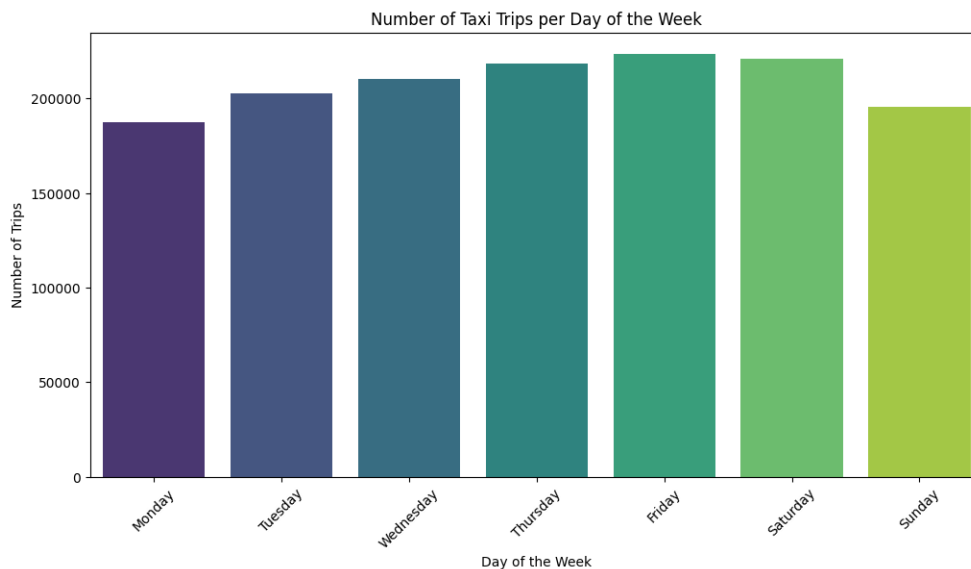
44 The code extracts the day of the week from the 'pickup_datetime' column using the dt.day_name()
45 method, creating a new 'day_of_week' column.

46 Trip Counting:

47 It then counts the number of trips for each day of the week using value_counts().

48 Reordering:

49 The days are reordered to follow the standard weekday sequence (Monday to Sunday) for better
50 interpretation.



51

52 3.2

53 Data Processing:

54 The NYC taxi dataset was loaded, and two new columns were created: 'hour_of_day' and
55 'day_of_week'. The data was then aggregated to count the number of pickups for each hour of
56 each day of the week.

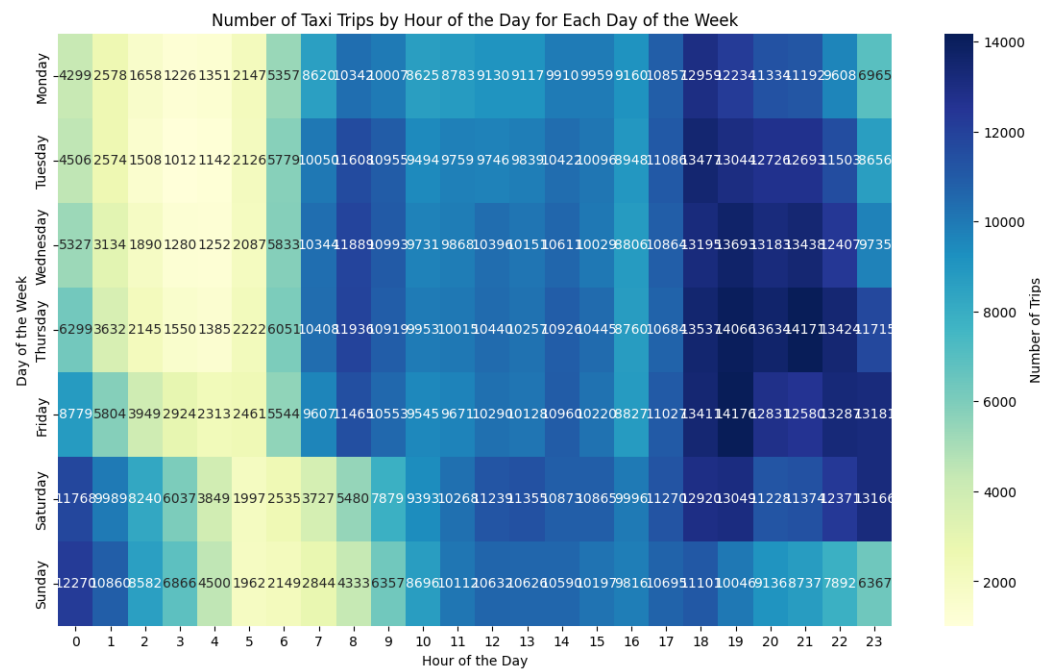
57 Data Analysis:

58 The code identified the most popular hour (with the highest number of pickups) for each day of the week.
59 A DataFrame 'most_popular_hours_df' was created to display these results, which is what we
60 see in the image.

day_of_week	Most Popular Hour	Count
Monday	18	12959
Tuesday	18	13477
Wednesday	19	13693
Thursday	21	14171
Friday	19	14176
Saturday	23	13166
Sunday	0	12270

61

62



63

64

65

66 3.3

67 Observations and Explanations You will see higher numbers of pickups during early morning 7-9
68 AM and evening hours 5-10 PM on weekdays, corresponding to typical work commute times.

69

70 You will see higher numbers of pickups during 6 PM Saturday evening to 1 AM on Sunday morning,
71 corresponding to late night outings on weekends.

72

73 Weekends vs. Weekdays: Weekends show different patterns, such as more late-night pickups
74 compared to weekdays. This could be due to social activities or nightlife.

75

76 3.4

77 Data Preparation:

78 Extracts hour of the day, day of the week, and date from the pickup_datetime. Defines specific
79 holiday dates for 2016.

80

81 Holiday Data Check:
82 Prints the holiday dates. Checks if these holidays are present in the dataset. Filters the dataset for
83 trips on holiday dates.

84 Holiday Analysis:
85 If holiday data is available:
86 Aggregates the number of pickups per hour for each holiday. Creates a heatmap to visualize the
87 distribution of trips by hour on holidays.
88

89 Regular Day Analysis:
90 Filters out holiday data to get regular days. Aggregates the number of pickups per hour for each day
91 of the week. Reorders days of the week and ensures all hours are represented. Creates a heatmap to
92 visualize the distribution of trips by hour on regular days.

93

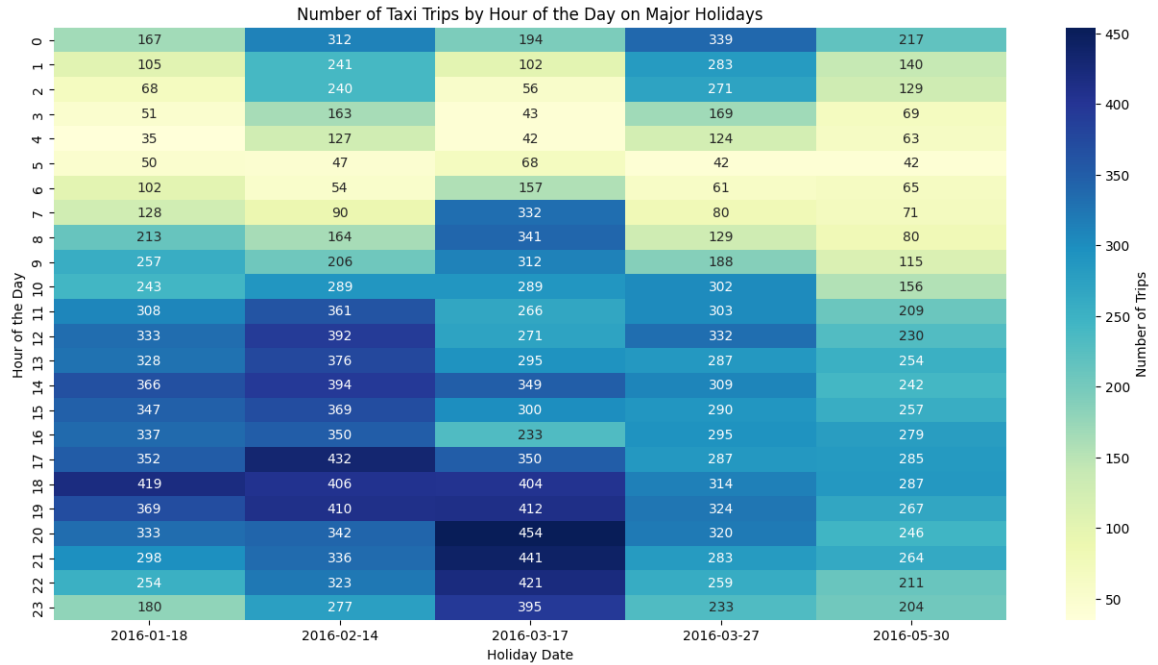
```
Holiday dates: [datetime.date(2016, 1, 18) datetime.date(2016, 2, 14)
datetime.date(2016, 3, 17) datetime.date(2016, 3, 27)
datetime.date(2016, 5, 30)]
Are holidays present in the dataset?
date
False    1421832
True      36812
Name: count, dtype: int64
Filtered holiday data:
   id  vendor_id  pickup_datetime  dropoff_datetime  \
25  id0129640    2 2016-02-14 13:27:56 2016-02-14 13:49:10
121 id2648478    1 2016-01-18 11:13:59 2016-01-18 11:18:56
162 id0762989    2 2016-03-17 08:24:27 2016-03-17 08:26:11
277 id2437858    2 2016-03-27 11:55:02 2016-03-27 12:05:06
346 id0861216    2 2016-01-18 13:00:37 2016-01-18 13:10:57

   passenger_count  pickup_longitude  pickup_latitude  dropoff_longitude  \
25                1      -73.956581      40.771358      -73.974968
121               1      -73.951576      40.766468      -73.960213
162               1      -73.977615      40.763573      -73.972572
277               5      -73.986832      40.761829      -73.977837
346               1      -73.990250      40.757286      -73.963982

   dropoff_latitude  stone_and_fwd_flag  trip_duration  hour_of_day  \
25      40.732792                N          1283          13
121      40.760540                N           297          11
162      40.765957                N           104           8
277      40.792122                N           604          11
346      40.756920                N           620          13

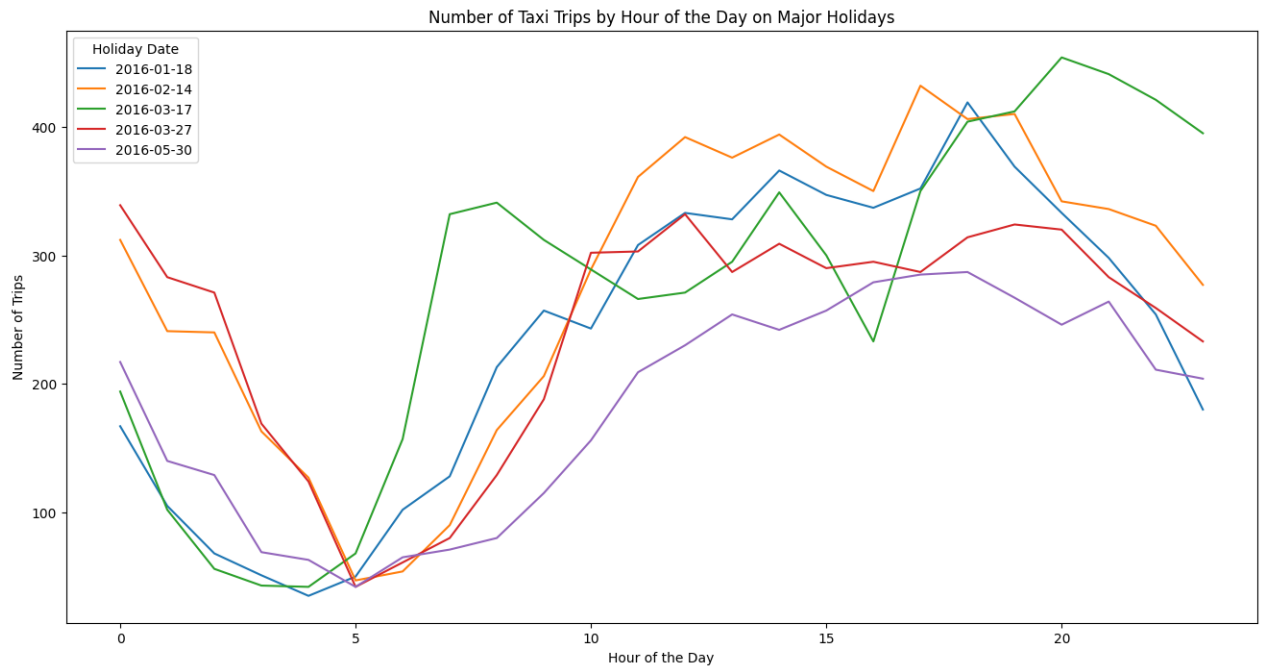
   day_of_week  date
25      Sunday 2016-02-14
121     Monday 2016-01-18
162     Thursday 2016-03-17
277      Sunday 2016-03-27
346     Monday 2016-01-18
```

94



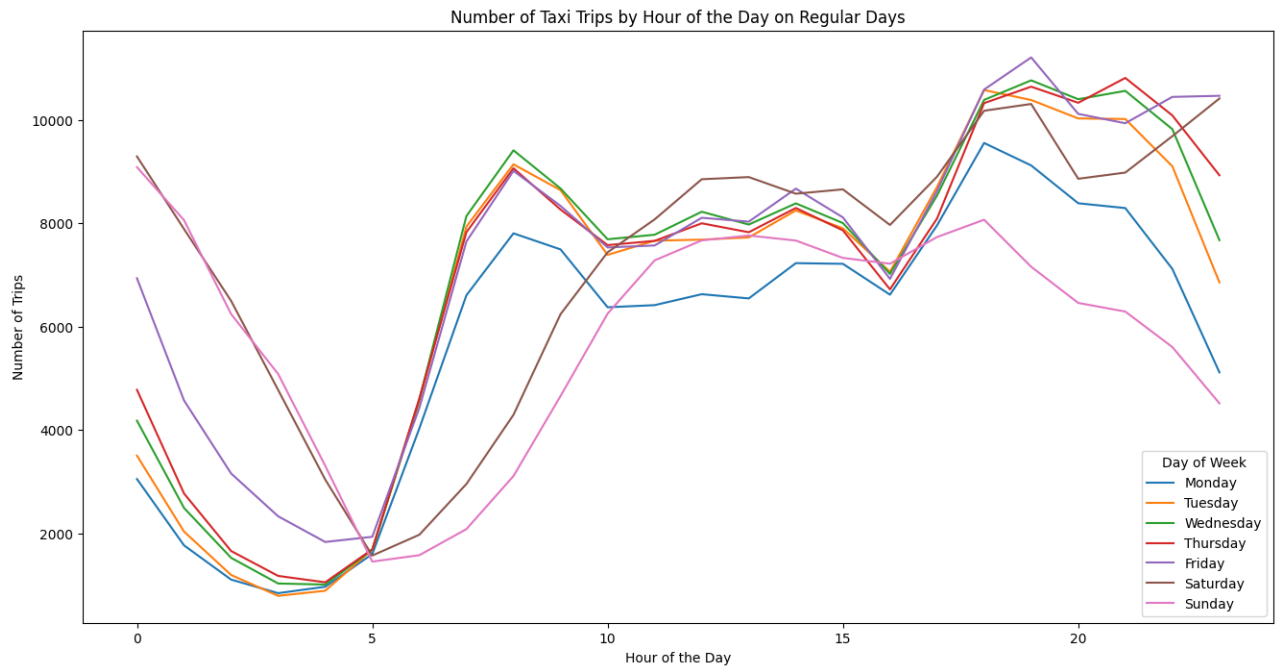
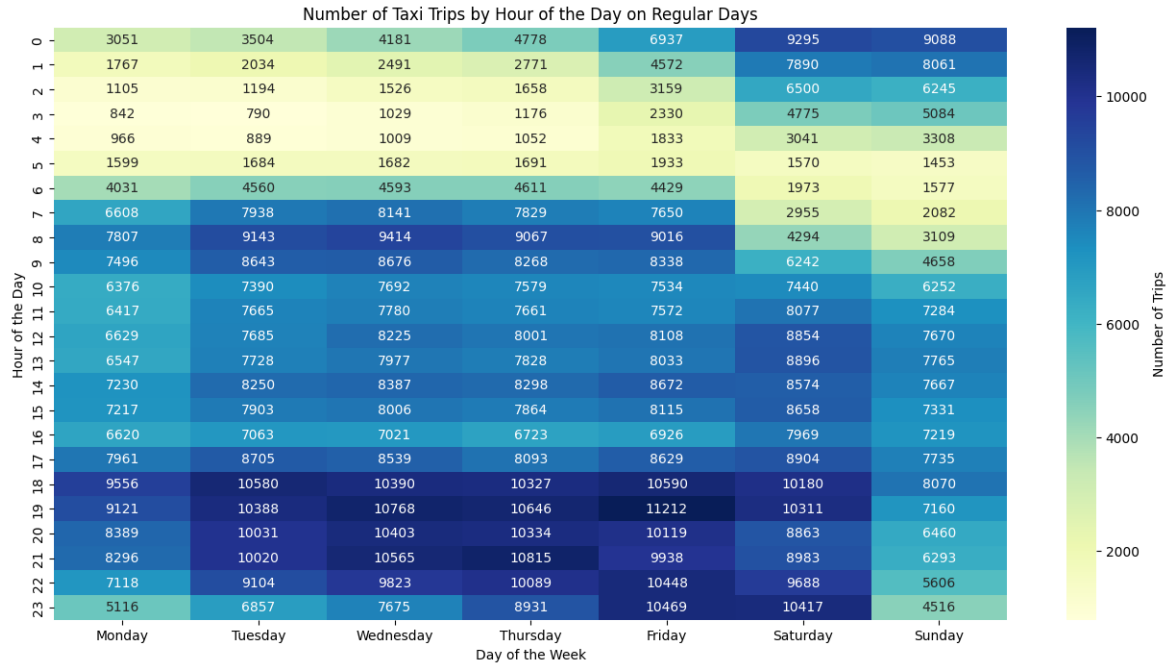
95

96



97

98



3.5

Data Preparation:

Creates a copy of the original dataframe to avoid warnings. Calculates the distance of each trip using the Haversine formula. Ensures datetime columns are in the correct format.

Feature Engineering:

Computes trip duration in seconds. Calculates speed in km/h for each trip. Extracts the hour of the

110 day from the pickup datetime.

111

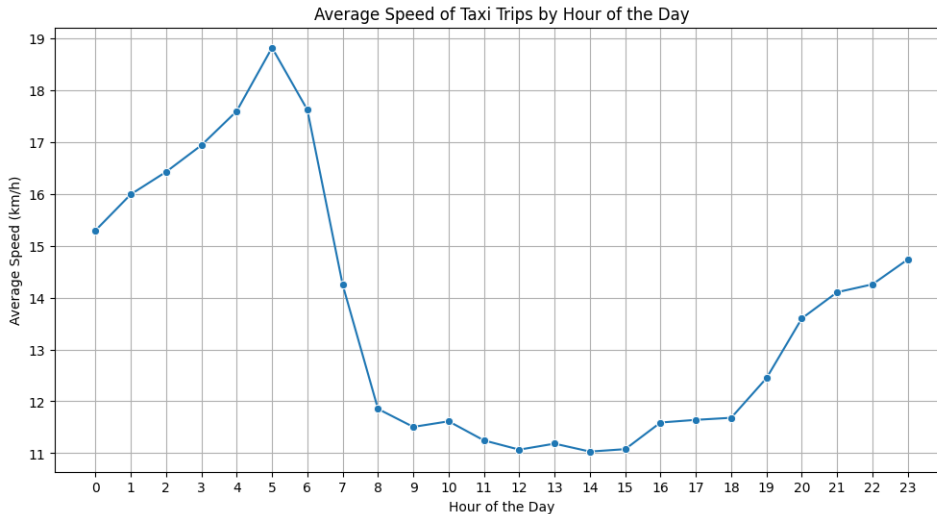
112 Data Cleaning:

113 Removes outliers identified earlier in the analysis.

114 Analysis:

115 Groups the data by hour of the day and calculates the mean speed for each hour.

116 Visualization:



117

118

119 The fastest time of day is 5:00 with an average speed of 18.82 km/h.

120 4 Question 4 : Location clusters

121 4.1 Heatmaps

122 To analyze the distribution of trip pickups, we generated a heatmap based on different time periods.

- 123 • **Weekday Rush Hours:** Clear peaks are visible on weekdays, with the highest number of
124 trips occurring during evening rush hour (6-8 PM) and a secondary peak during morning
125 rush hour (8-9 AM).
- 126 • **Weekend Late Night Activity:** Saturdays and Sundays show distinctly different patterns,
127 with the busiest times occurring late night to early morning (11 PM to 3 AM), especially on
128 Saturday night/Sunday morning.
- 129 • **Daily Cycle:** Across all days, trip numbers are lowest between 4-5 AM, then gradually
130 increase throughout the day, peaking in the evening on weekdays and late night on weekends.
- 131 • **Friday Surge:** Friday consistently shows the highest overall activity, particularly in the
132 evening and night hours, likely due to a combination of commuter traffic and the start of
133 weekend social activities.
- 134 • **Weekday vs Weekend Contrast:** Weekdays display structured patterns aligned with typical
135 work schedules, while weekends show a more evenly distributed trip pattern throughout the
136 day and night, reflecting leisure and social behaviors.

137 4.2 Hotspots

138 To identify popular taxi hotspots, we used DBSCAN clustering on two different time frames: from
139 23:00 on a Friday evening to 02:00 on a Saturday morning, and from 17:00 to 20:00 on a Thursday

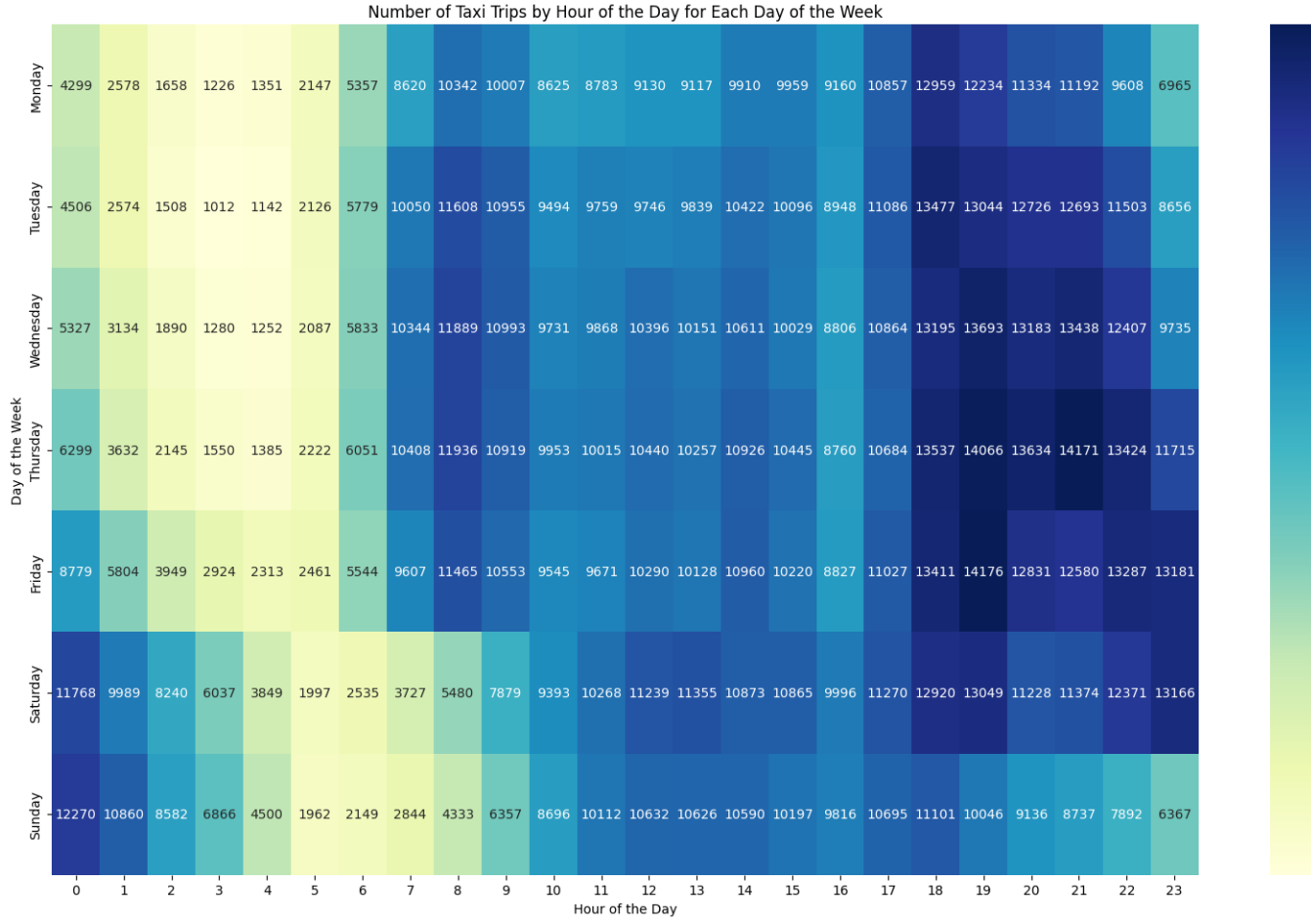


Figure 1: Heatmap of Trip Pickups by Time of Day and Hour

evening. We defined a hotspot as a cluster of at least 15 pickups within a maximum distance of 75 meters. The DBSCAN algorithm determined the number of clusters based on these parameters.

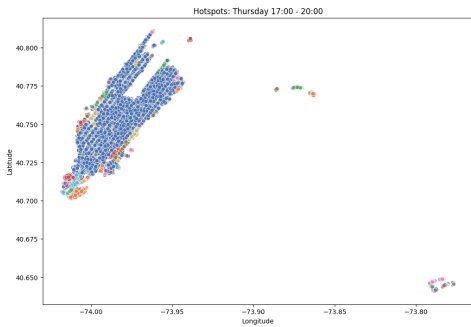


Figure 2: Hotspot Locations: Thursday 17:00 - 20:00

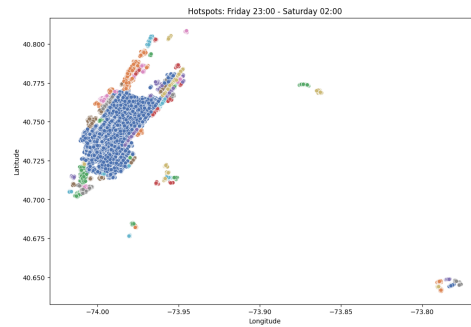


Figure 3: Hotspot Locations: Friday 23:00 - Saturday 02:00

Figure 4: Hotspot Locations for Different Time Periods

fig. 1 illustrates:

- **Cluster Comparison:** More hotspots were identified for Friday night/Saturday morning (87) than Thursday evening (71), suggesting more dispersed taxi activity during weekend nights.
- **Hotspot Variation:** Different cluster sizes and colors suggest varying levels of taxi demand across locations, with some areas showing more intense activity than others.

5 Question 5 : Airports

We examined the average travel time from the Empire State Building to JFK and Newark Airports. The radii for each location were set using the 75th percentile of the distance distribution to effectively cover typical travel distances. For the Empire State Building, a radius of 3.67 km was chosen. For JFK Airport, a radius of 21.49 km was used, and for Newark Airport, 19.58 km. This approach ensures the radii reflect the common travel distances while minimizing the influence of outliers.

Table 1: Coordinates and radii for the selected locations.

Location	Coordinates	Radius (km)
Empire State Building	(40.756724, -73.983806)	3.67
JFK Airport	(40.647929, -73.777813)	21.49
Newark Airport	(40.689442, -74.173242)	19.58

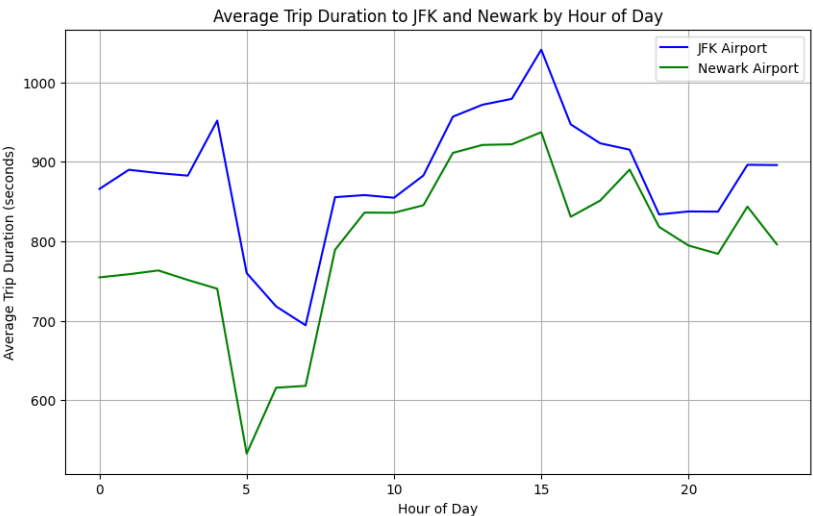


Figure 5: Average Trip Duration from the Empire State Building to JFK and Newark Airports

154 The plot in fig. 5 illustrates the average trip duration from the Empire State Building to JFK and
155 Newark Airports throughout the day. The results indicate that trips to JFK Airport generally have
156 longer average durations compared to trips to Newark Airport. This trend is observed consistently
157 throughout different times of the day. The results reveal several key patterns:

158 • **Daily Variations:**

- 159 – Both JFK and Newark Airport trip durations show significant fluctuations by the hour.
160 Notably, average trip durations increase during the early morning hours, particularly
161 between 6 and 7 AM. This increase is likely due to peak traffic times, as people
162 commute to work or catch early flights.
- 163 – After 3 PM, there is a general decrease in average trip durations to both airports. This
164 decline may be attributed to reduced traffic congestion as the day progresses and rush
165 hour subsides.

166 • **Geographical and Traffic Factors:**

- 167 – The consistently higher average travel times to JFK compared to Newark are likely
168 due to the greater geographical distance and potentially more congested or less direct
169 routes to JFK Airport.

170 6 Question 6: Boroughs

171 6.1 Neighbourhoods

172 Using the provided shapefile, we identified the neighborhoods for the trip start and end locations using
173 GeoPandas. The unique neighborhoods for the pickup locations and dropoff locations respectively
174 are:

```
175 ['Lincoln Square', 'Murray Hill-Kips Bay', 'Midtown-Midtown South',  
176 'SoHo-TriBeCa-Civic Center-Little Italy', 'Upper West Side', 'Gramercy', ...]  
  
177 ['Upper East Side-Carnegie Hill', 'West Village',  
178 'Battery Park City-Lower Manhattan',...,]
```

179 6.2 Choropleth

180 The map in fig. 6 highlighted neighborhoods with varying counts of pickup and dropoff locations.
181 Areas with more pickups or dropoff were shown in darker colors, while those with fewer pickups
182 appeared lighter.

183
184 The choropleth maps for pickups and dropoffs revealed that the shading intensity was the same for
185 both cases across neighborhoods. This observation indicates that the spatial distribution of pickups
186 and dropoffs is consistent. In other words, neighborhoods with high pickup counts also tend to have
187 high dropoff counts, and vice versa. Implying the following :

- 188 • **Similar Distribution:** The uniformity in shading suggests that the distribution of transporta-
189 tion activity is balanced across neighborhoods. Areas with high levels of both pickups and
190 dropoffs are likely major transportation hubs or popular destinations, while those with low
191 activity are consistently less active in both respects.
- 192 • **High Traffic Areas:** Neighborhoods showing intense colors in both maps are likely experi-
193 encing high volumes of both pickups and dropoffs, indicating significant transportation
194 activity.
- 195 • **Uniform Patterns:** The consistent shading indicates that there is no significant difference in
196 where trips start and end. This might be due to even distribution of trips or uniform service
197 patterns across neighborhoods.

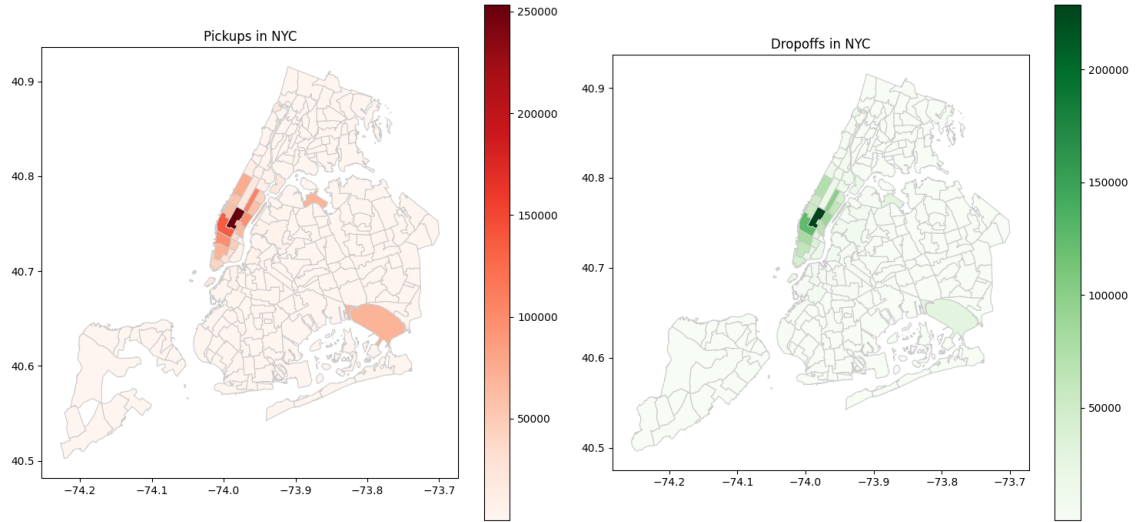


Figure 6: Choropleth for Pickup and Dropoff trips in NYC

6.3 Quietest Neighborhoods Between Midnight and 5 AM

To determine the quietest neighborhood(s) between midnight and 5 AM, we filtered trips that started and ended within this time window and counted the number of trips in each neighborhood. We define "quiet" as having the fewest trips. The table below shows the top 3 neighborhoods with the fewest trips, indicating they are the least busy during this time period. The quietest neighborhood, with only a single trip, is Oakwood-Oakwood Beach.

Table 2: Top 3 quiet neighborhoods with the fewest trips between midnight and 5 AM.

Neighborhood	Trip Count
Oakwood-Oakwood Beach	1
Annadale-Huguenot-Prince's Bay-Eltingville	2
Old Town-Dongan Hills-South Beach	2

204 6.4 Busiest Neighborhoods Between Midnight and 5 AM

205 To determine the busiest neighborhood(s) between midnight and 5 AM, we filtered trips that started
 206 and ended within this time window and counted the number of trips in each neighborhood. We define
 207 "busiest" as having the most trips. The table below shows the top 3 neighborhoods with the most trips,
 208 indicating they are the busiest during this time period. The busiest neighborhood, with the highest
 209 number of trips, is Midtown-Midtown South.

Table 3: Top 3 busiest neighborhoods with the most trips between midnight and 5 AM.

Neighborhood	Trip Count
Midtown-Midtown South	14,507
Hudson Yards-Chelsea-Flatiron-Union Square	11,141
West Village	8,727