

NYC 311 Complaints data analysis

Big data analysis C.Sc. 84030 : Final project

VISHAL BHARTI

311 Service

- 311 is New York City's main source of government information and non-emergency services.
- NYC receives 311 calls for non-emergency services from its residents, businesses and visitors.
- Response time for these calls is longer than those for emergency (911) calls.
- The data is open and updated frequently.
- Data used in the project 1/1/2010-12/7/2016 (~8.8 gb).
- Comprise of 14.2 million rows and 53 columns.

Hypothesis

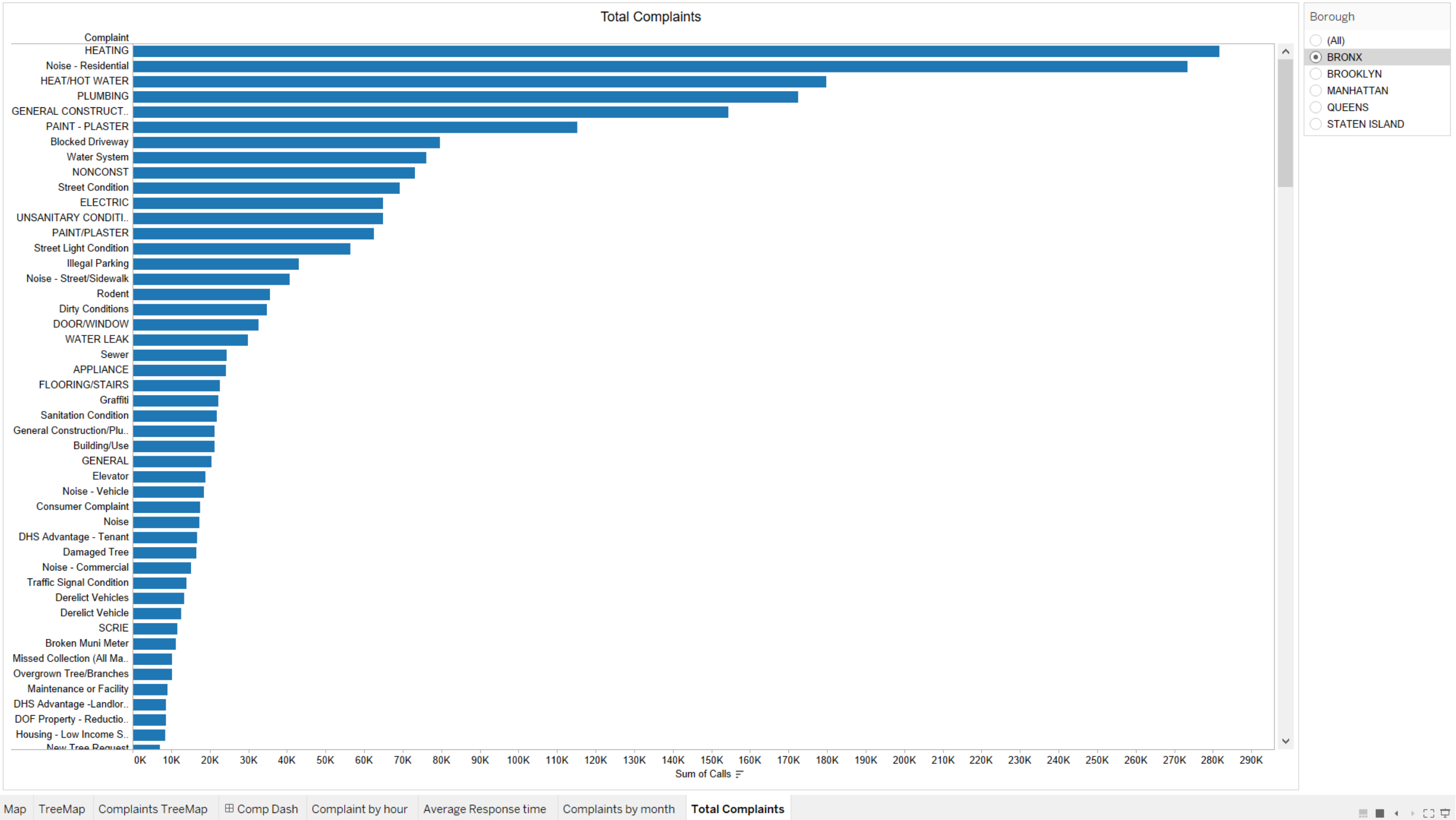
- Complaints will be vary for each location.
- Complaints will vary based on time of day, the day of week and time of the year.
- Noise complaints will be higher in residential areas.
- Parking violations complaints will increase during weekends.
- Street light conditions complaints will be higher during night time.
- Correlation between the number of complaints and response time.

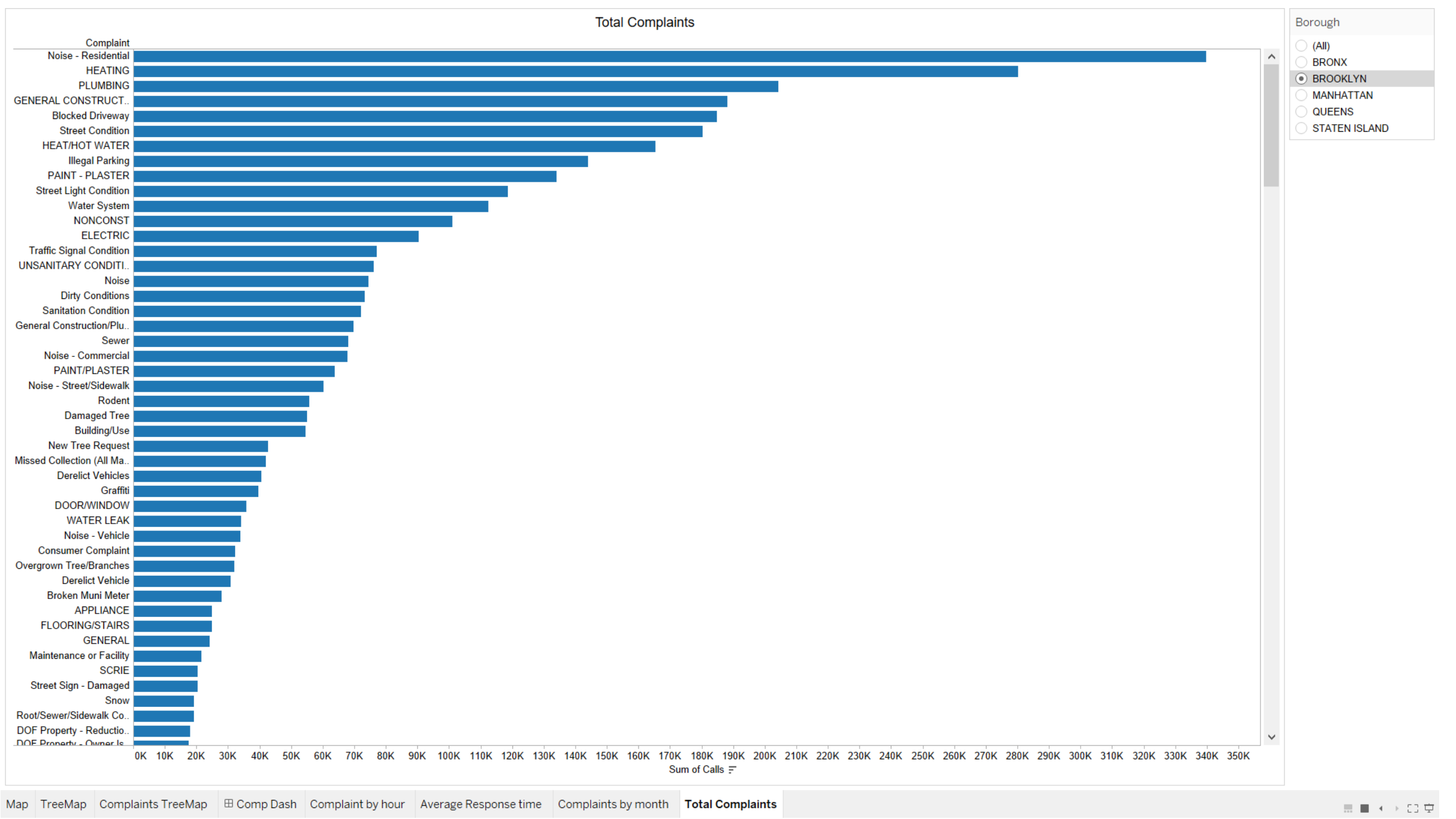
Methodology

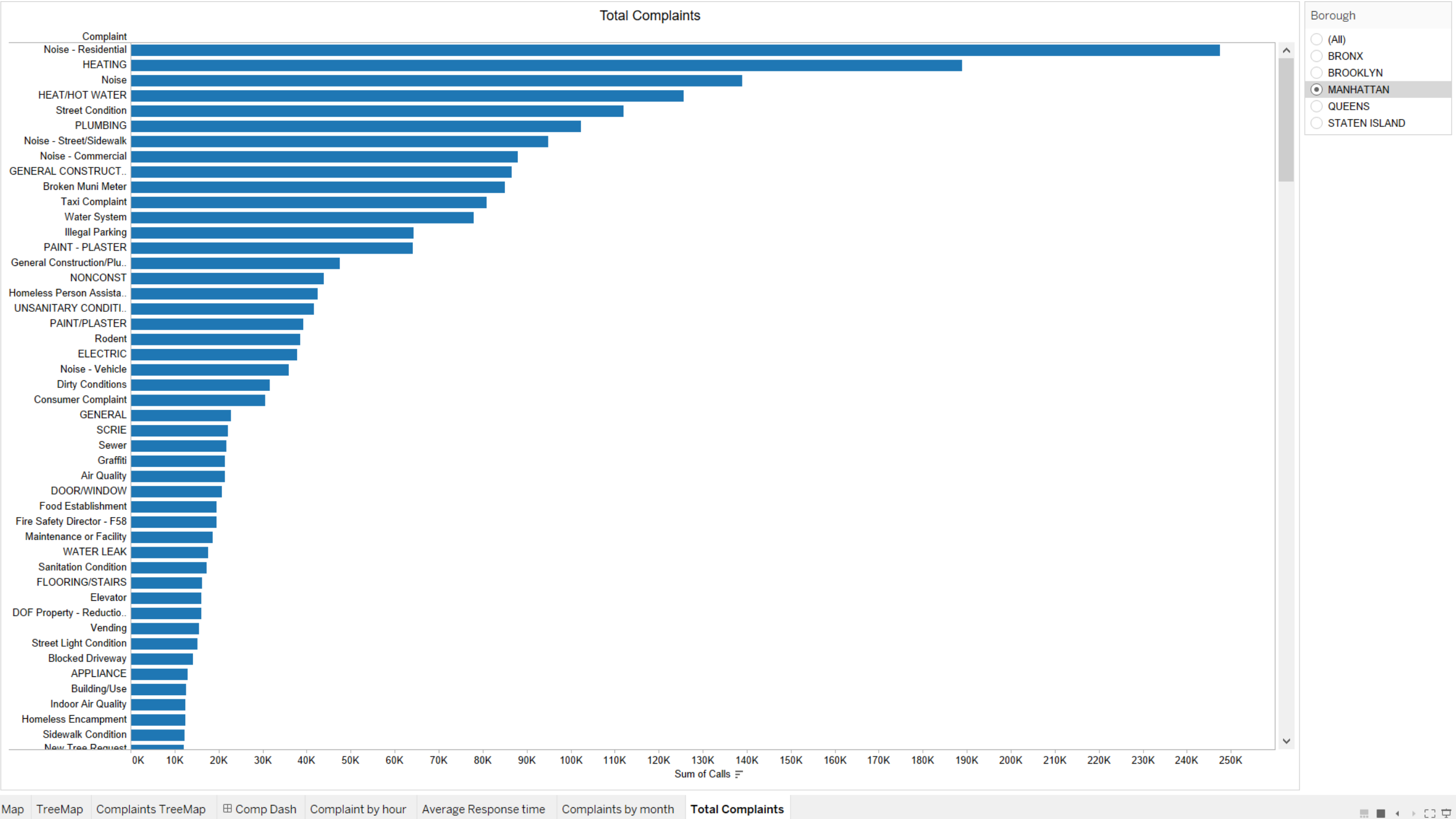
- Filtered dataset to remove noisy records. Some records had missing or incorrect closing date.
- Used Pyspark to filter the data based on keys.
- The data size was reduced from 8.8 gb to approx. 300mb.
- Visual analytics was used to draw information from the resulting dataset.
- Tableau was used for visual analytics.

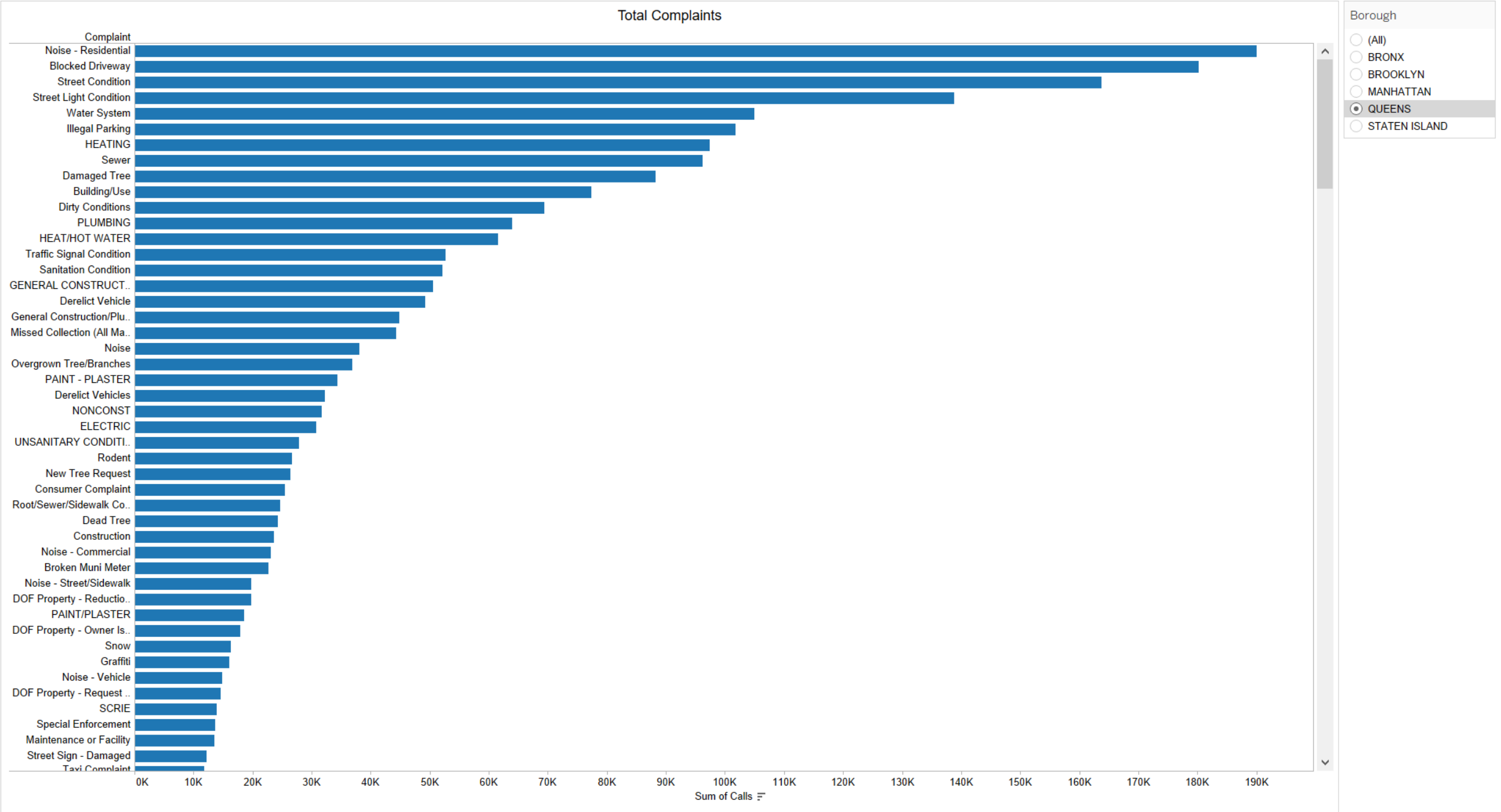
Complaints by location

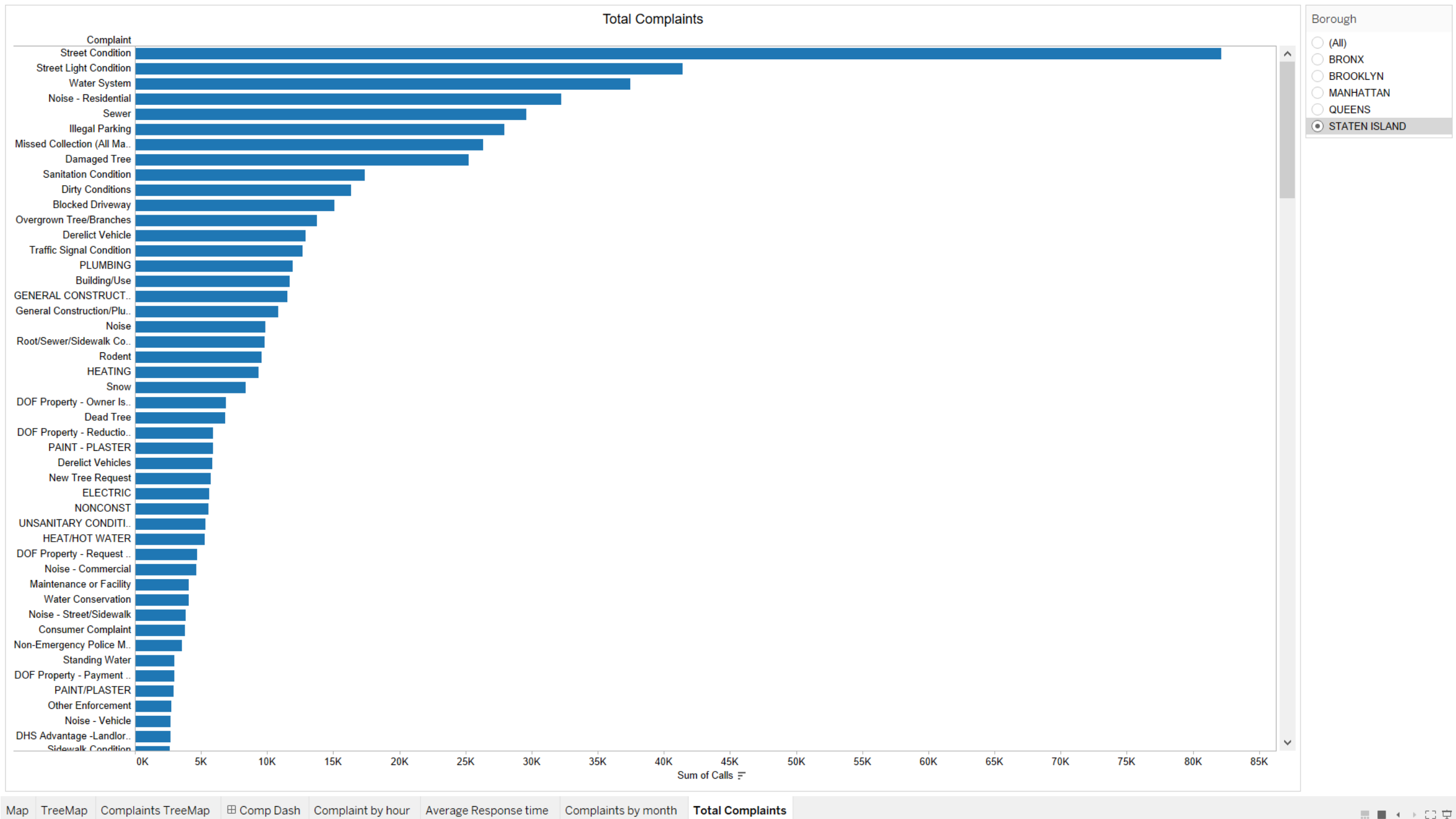
- For Bronx, 'Noise', 'Heating' and 'Street Light Conditions' were the most common complaints.
- For Brooklyn, 'Noise', 'Heating' and 'Street Conditions' were the most common complaints.
- Manhattan also had, 'Noise', 'Heating' and 'Street Conditions' as the most reported complaints.
- Queens had 'Noise', 'Street Light Conditions' and 'Blocked Driveway' as the most reported complaints.
- Staten Island had 'Street Light Conditions', 'Street Conditions' and 'Water System' as the most reported complaints.
- 'Noise', 'Heating' and 'Street Light Conditions' were the overall most reported complaints.







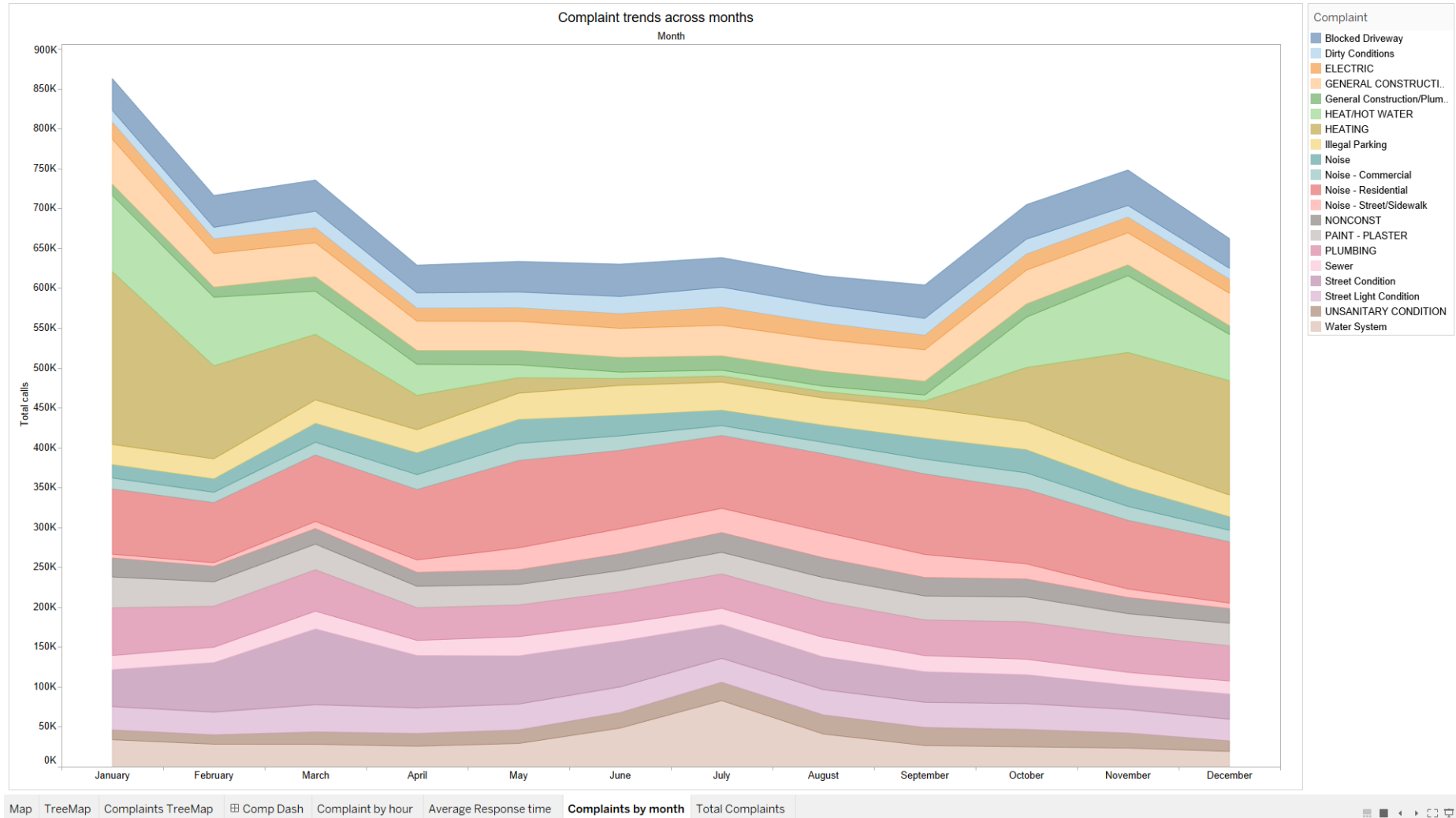




Complaints by month

- Heating related complaints were higher at the start and end of the year.
- Street condition complaints were somewhat higher at the start of the year and show a decline after that.
- Water system complaints show a peak in the middle of the year.
- Illegal parking complaints increase slightly in the second half of the year.
- Noise, Street Light Conditions, Blocked Driveway, etc., are some complaints that follow don't vary too much throughout the year.

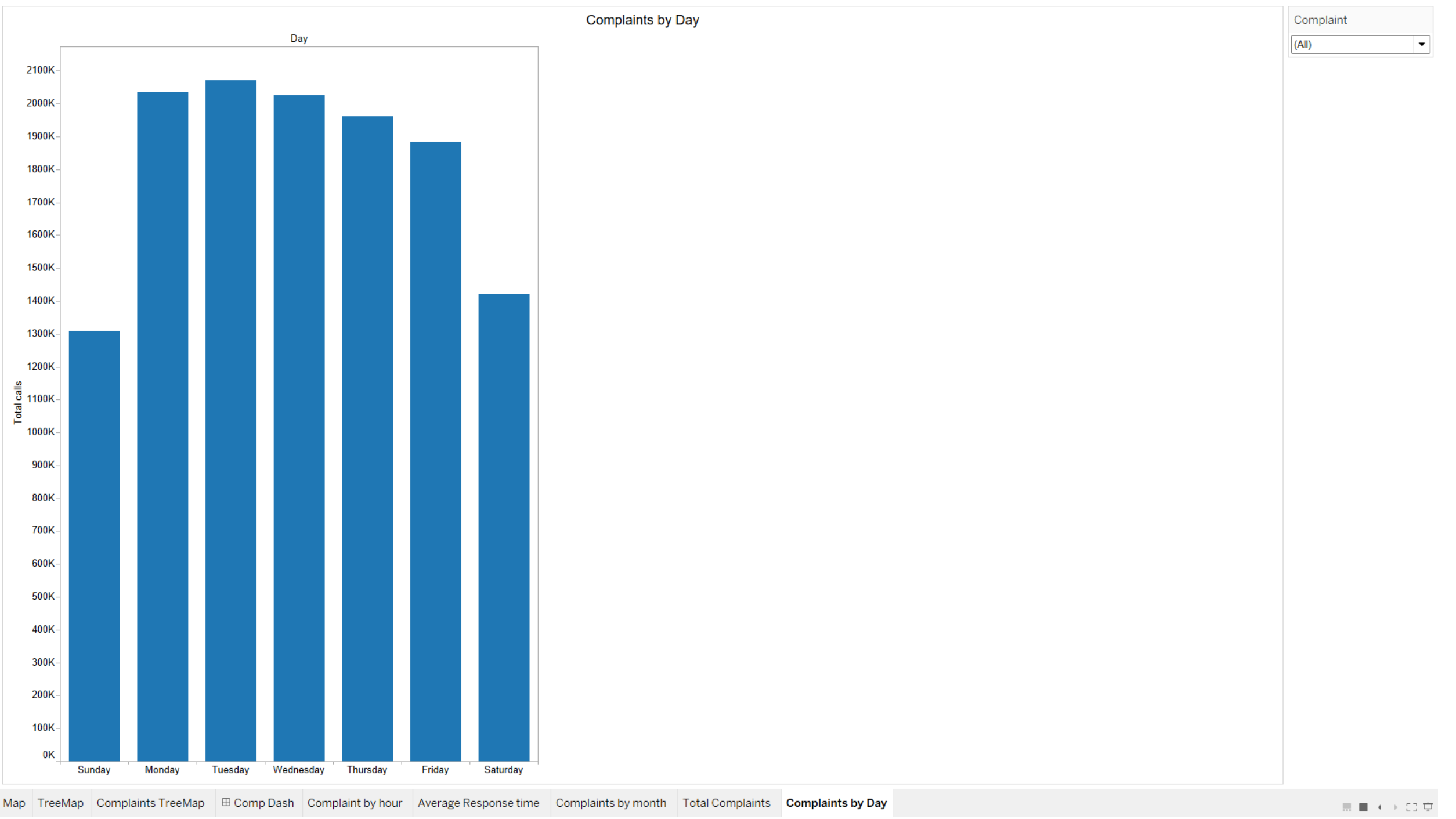
Trends for top 20 Complaints

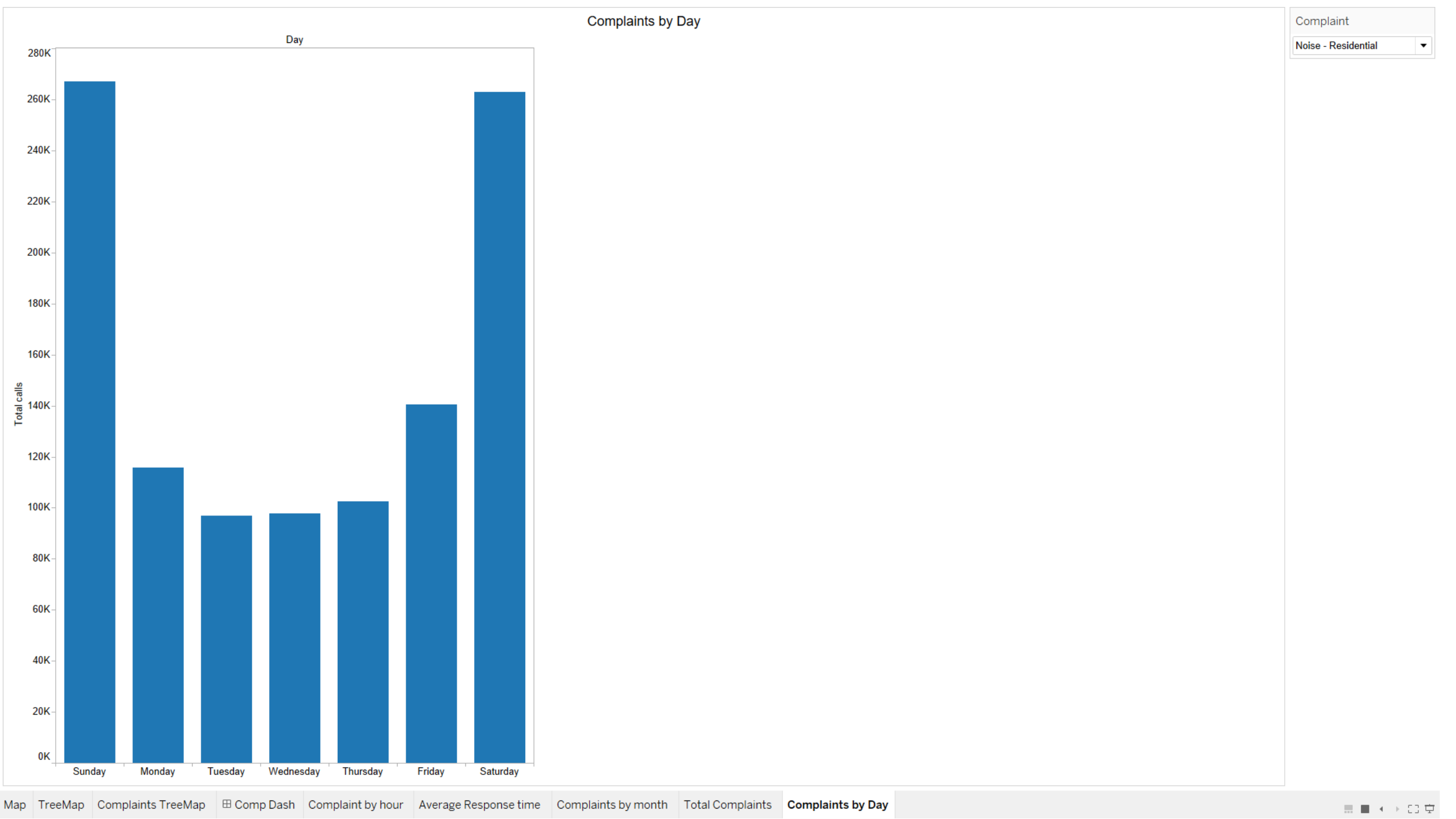


Complaints by day of week

- Complaints were higher in weekdays.
- Noise complaints and Parking Violations were significantly higher in weekends.
- Construction complaints were lower in weekends.
- Blocked driveway complaints were slightly higher in weekends.

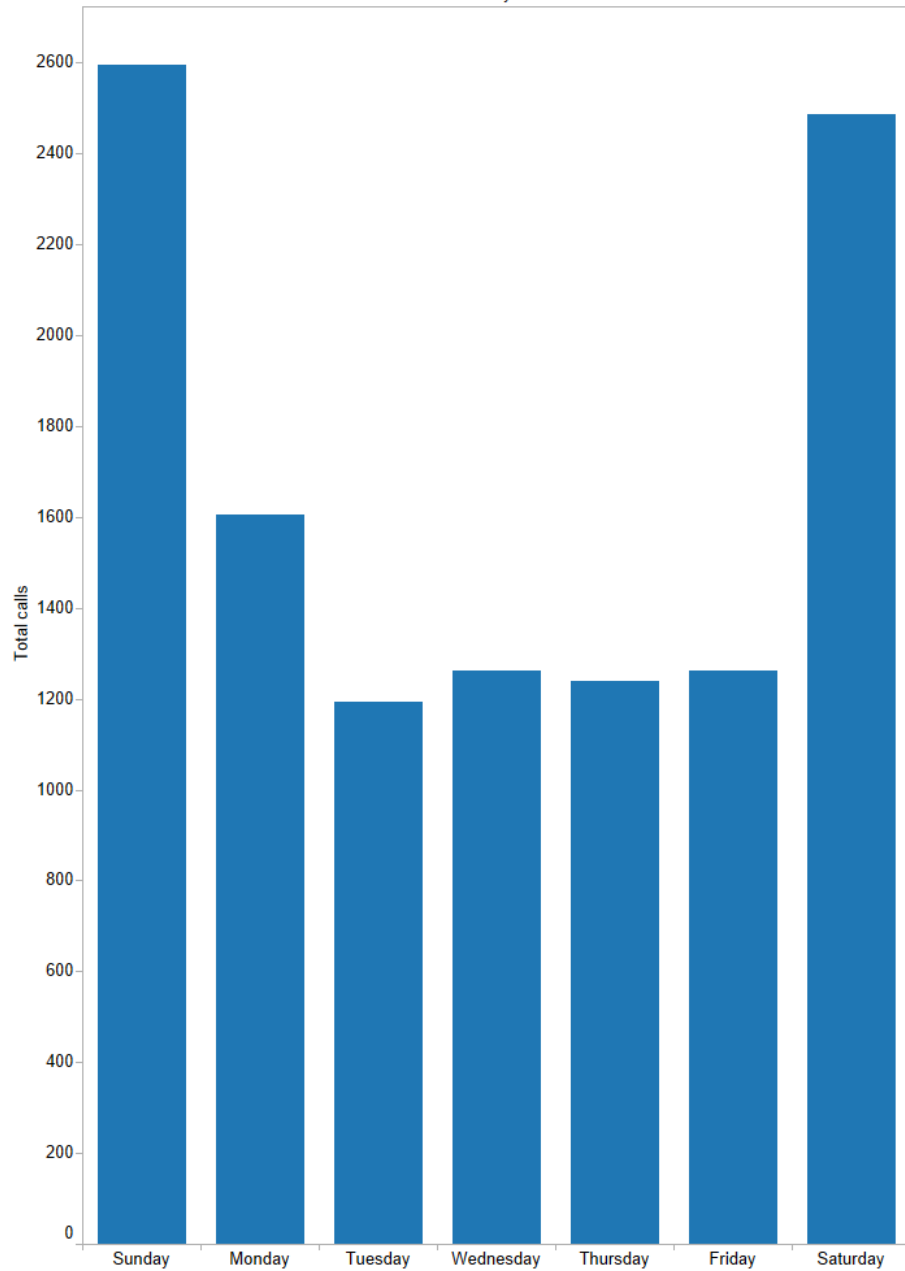
Complaints by Day





Complaints by Day

Day



Complaint

Violation of Park Rules



Map

TreeMap

Complaints TreeMap

Comp Dash

Complaint by hour

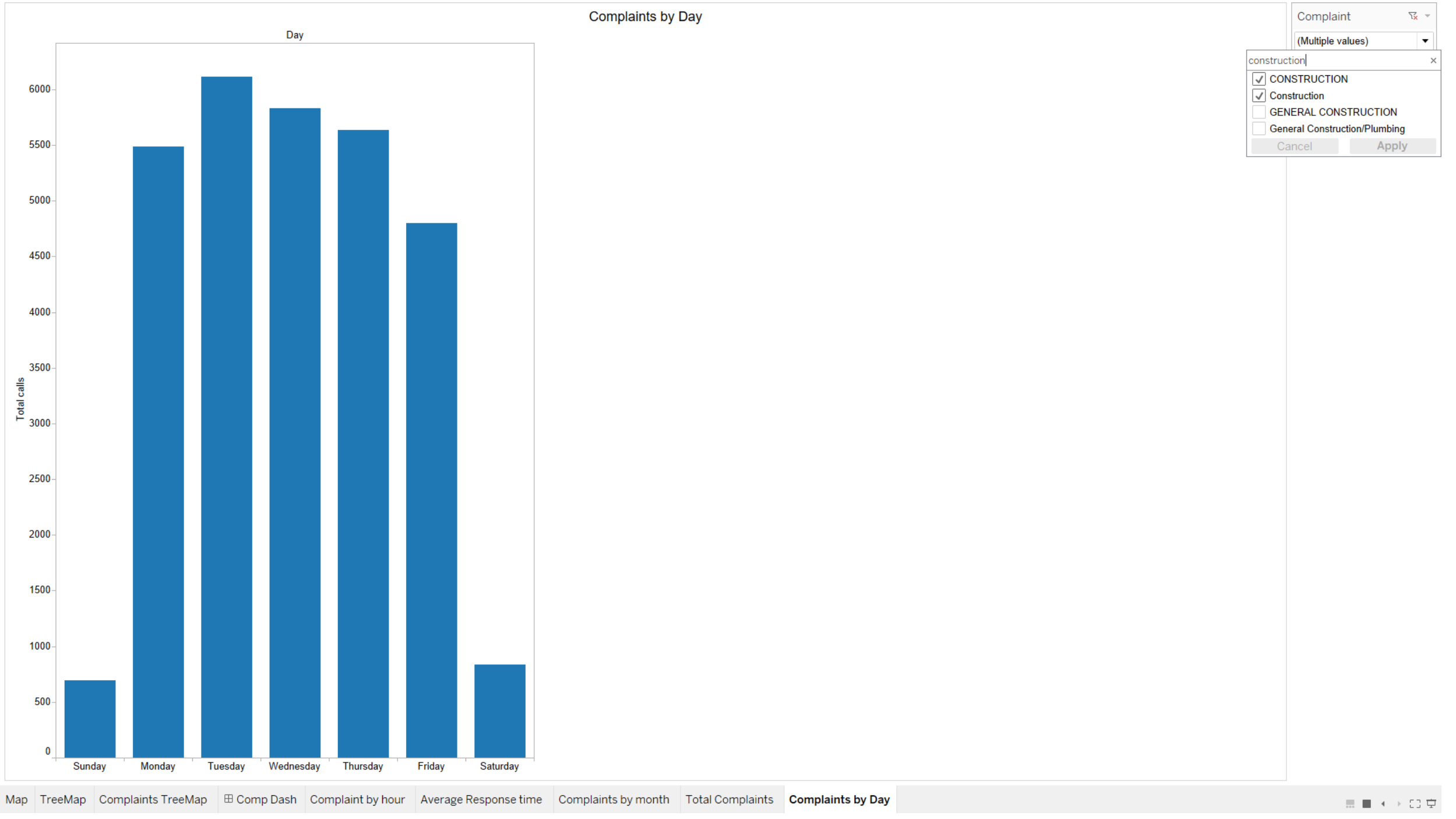
Average Response time

Complaints by month

Total Complaints

Complaints by Day





Map

TreeMap

Complaints TreeMap

Comp Dash

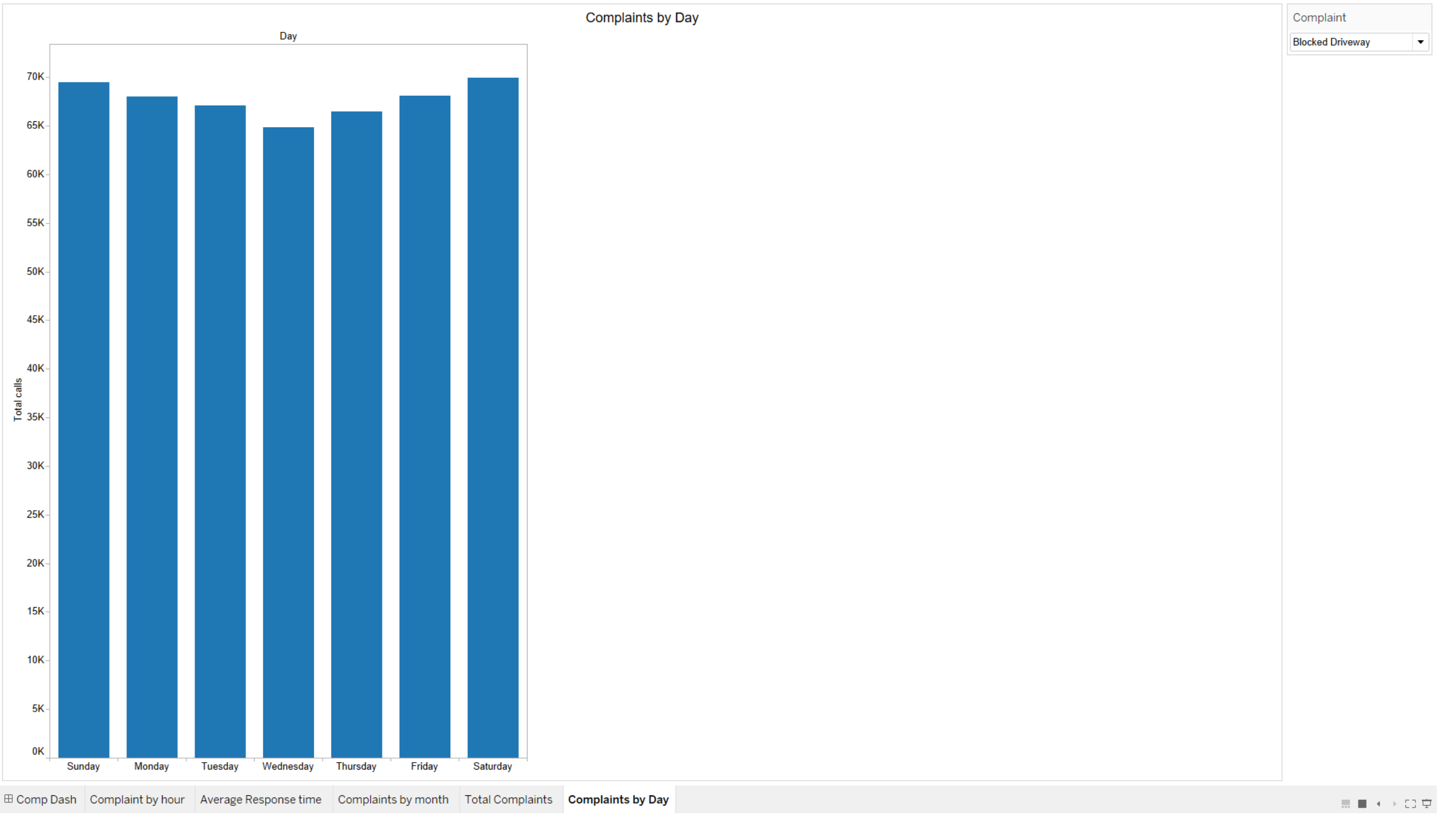
Complaint by hour

Average Response time

Complaints by month

Total Complaints

Complaints by Day

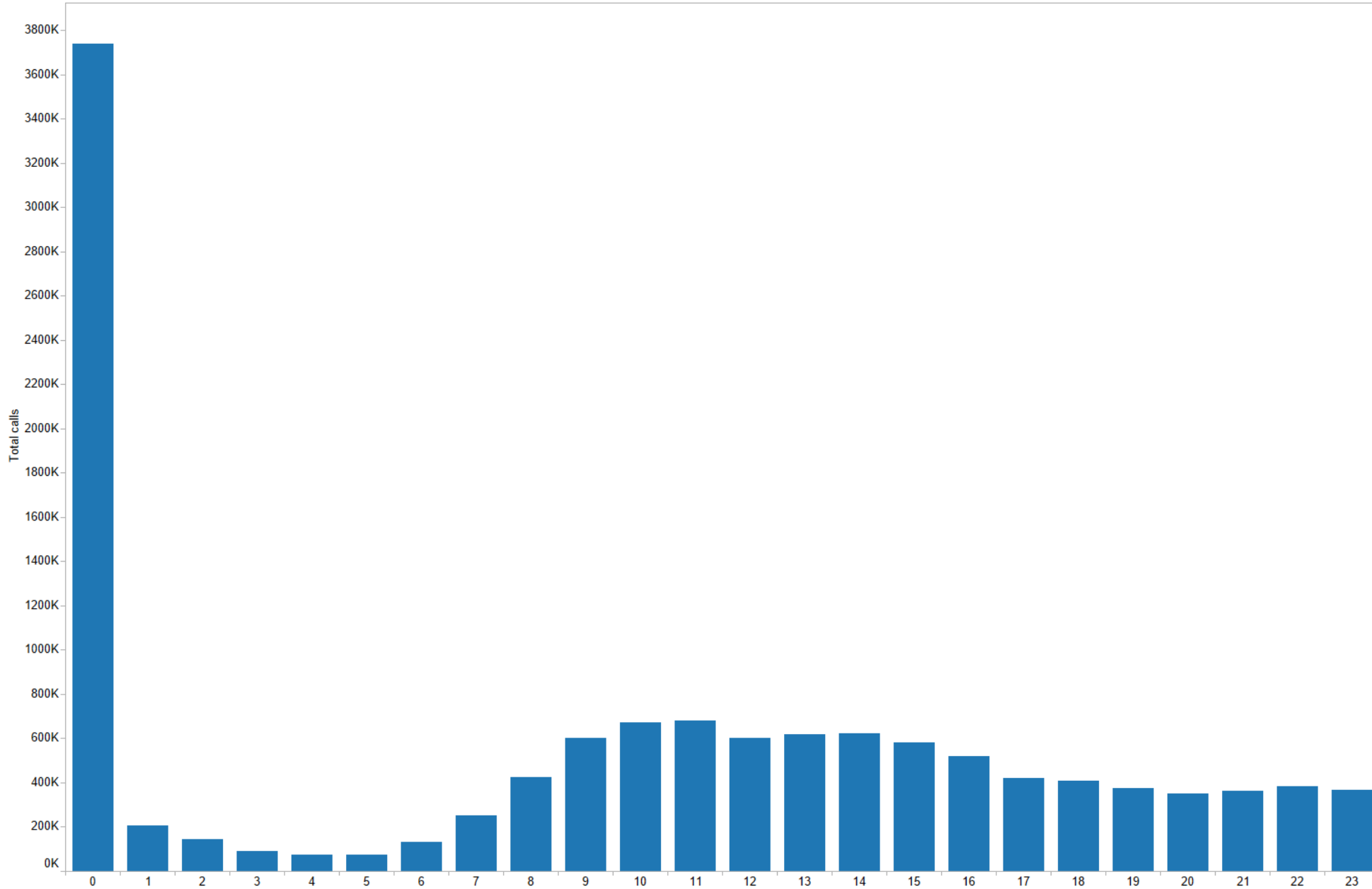


Complaints by time of the day

- Most complaints calls were made close to midnight.
- Illegal parking complaints were high throughout the daytime and fall out after midnight.
- Noise complaints reach a peak close to midnight.
- Parking violation complaints are maximum in the after noon, show a decline after that.
- Blocked driveway complaints are maximum in early morning and night periods.

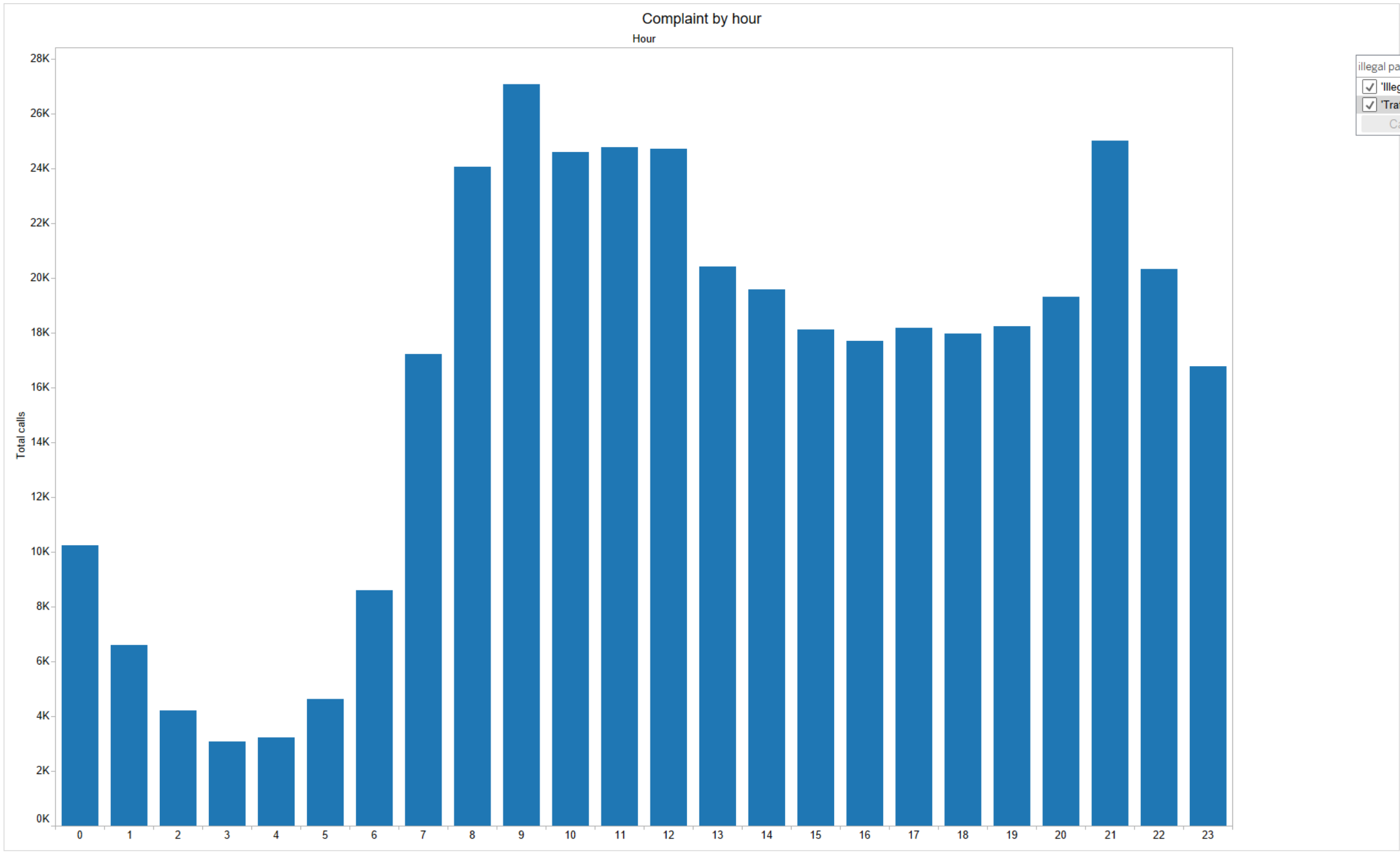
Complaint by hour


Hour




Complaint type

(All) ▼



Complaint type 

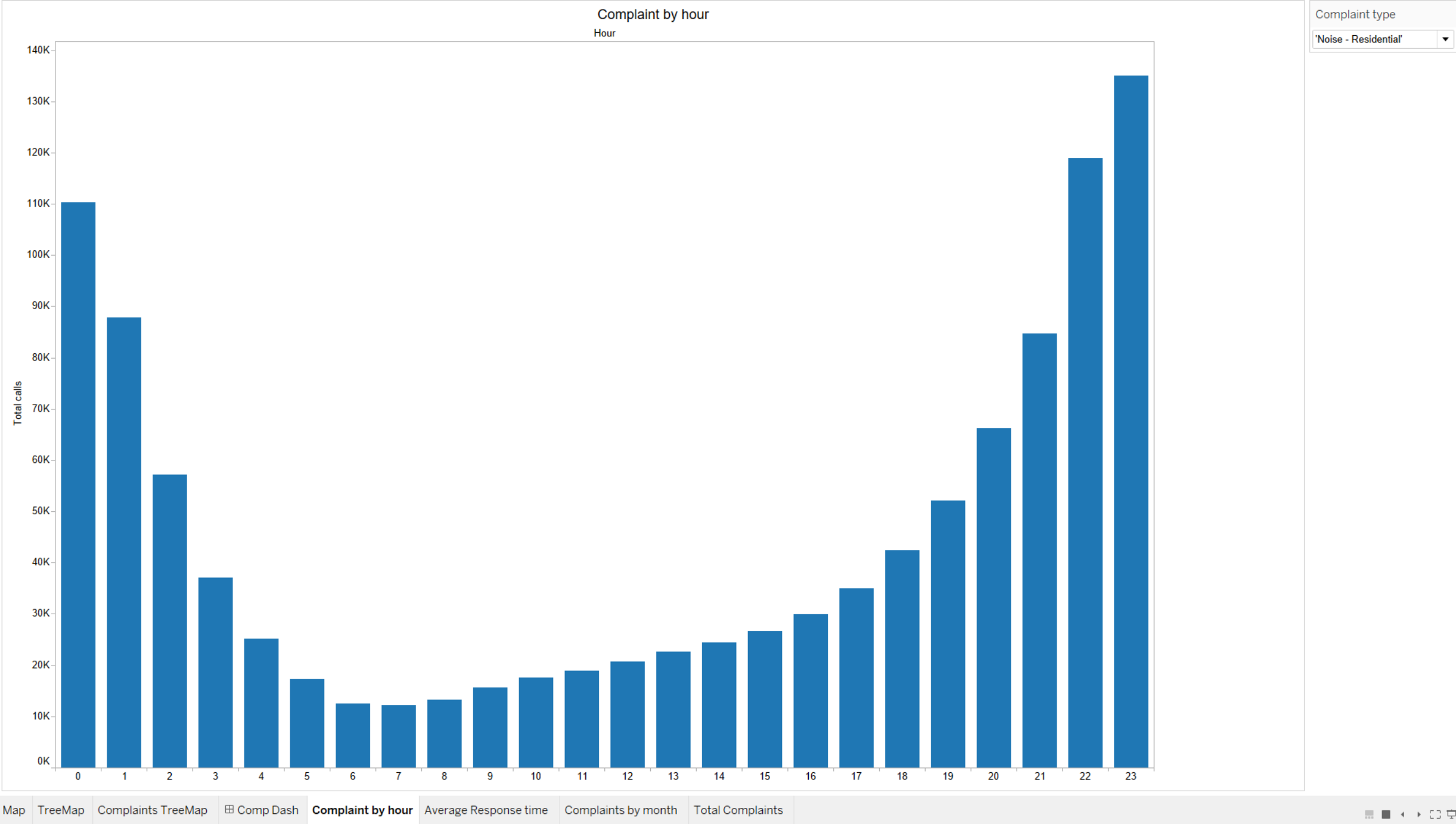
(Multiple values) ▼

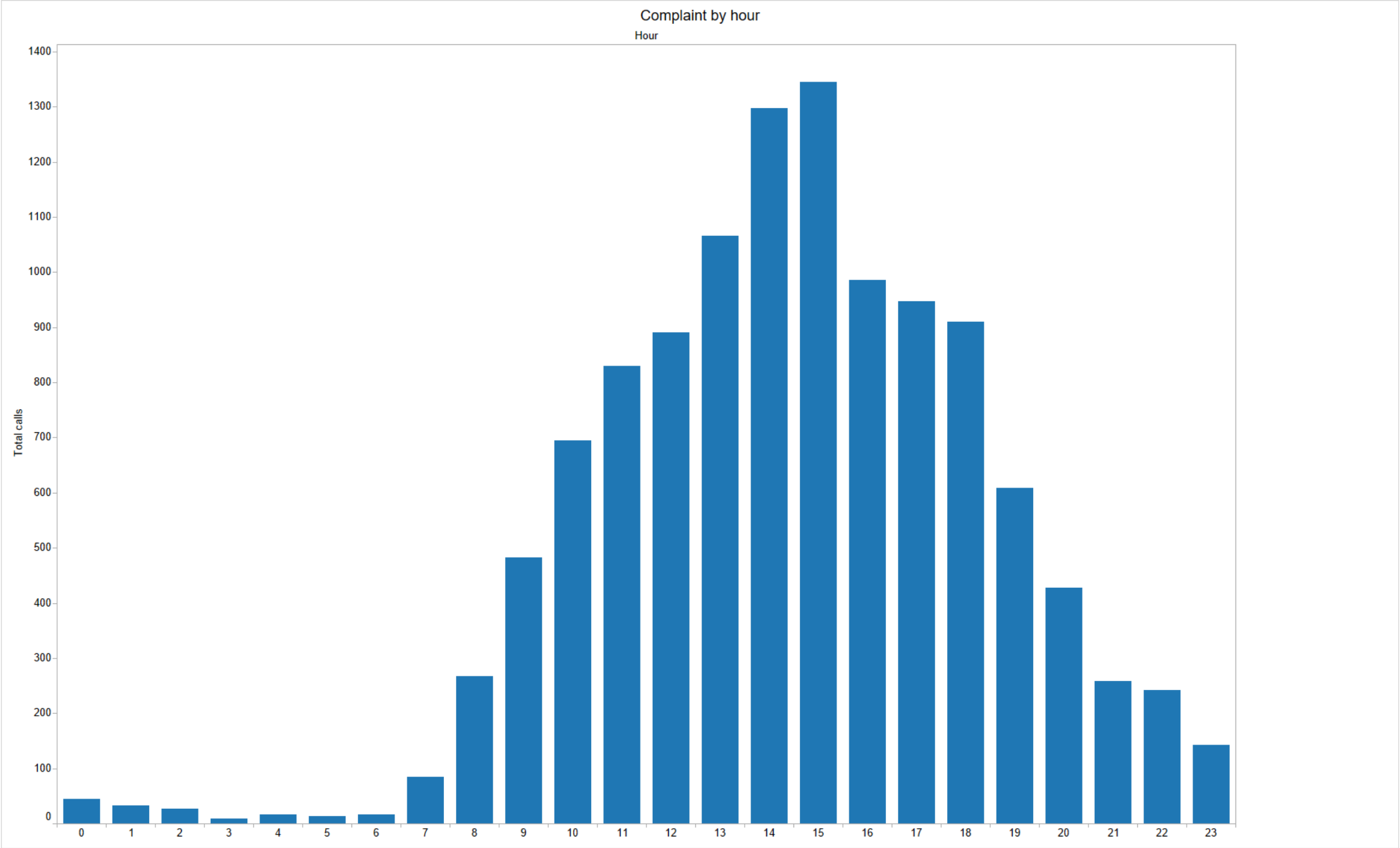
illegal parking 

☒ 'Illegal Parking'

☒ 'Traffic/Illegal Parking'

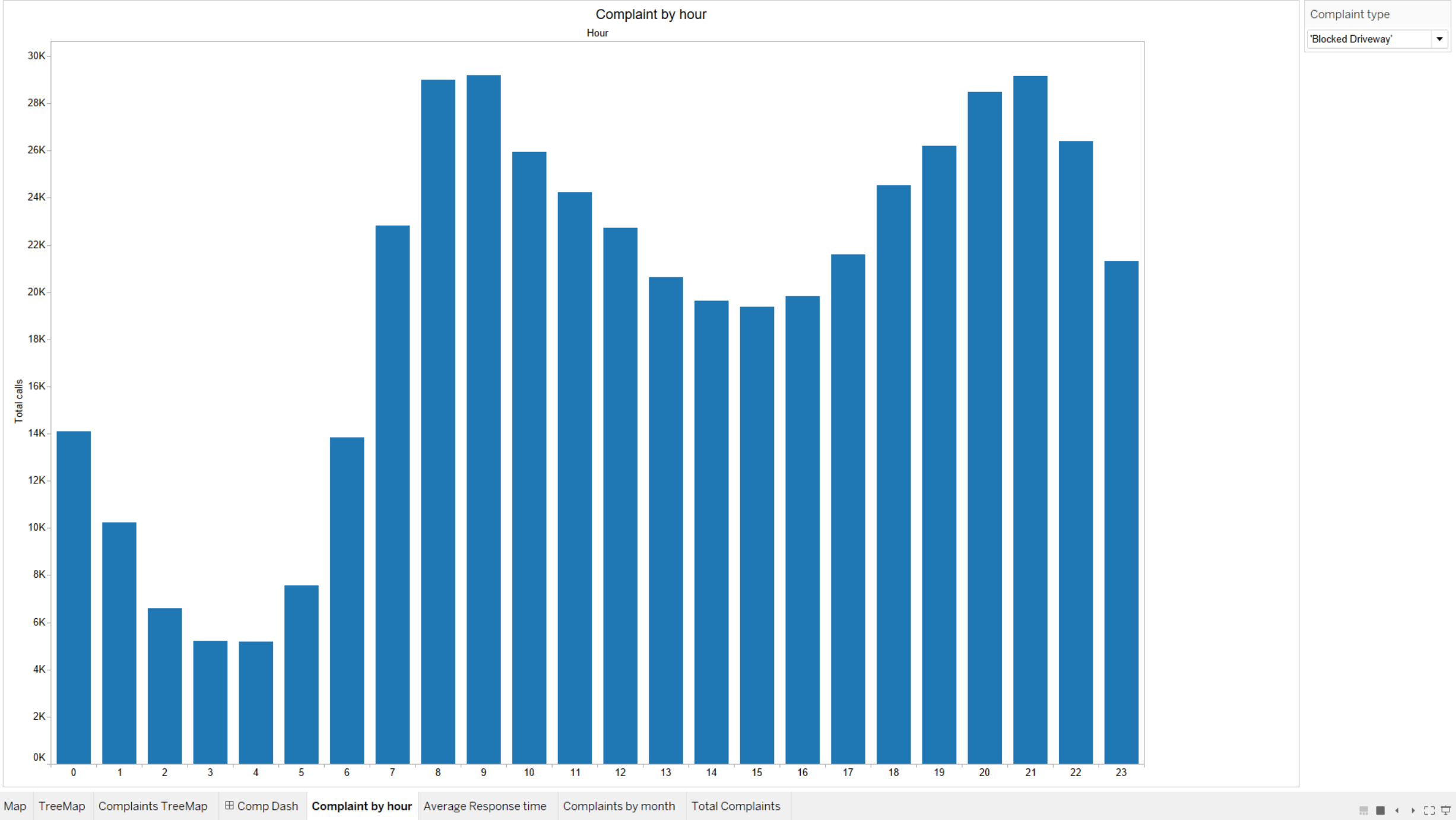
Cancel Apply





Complaint type

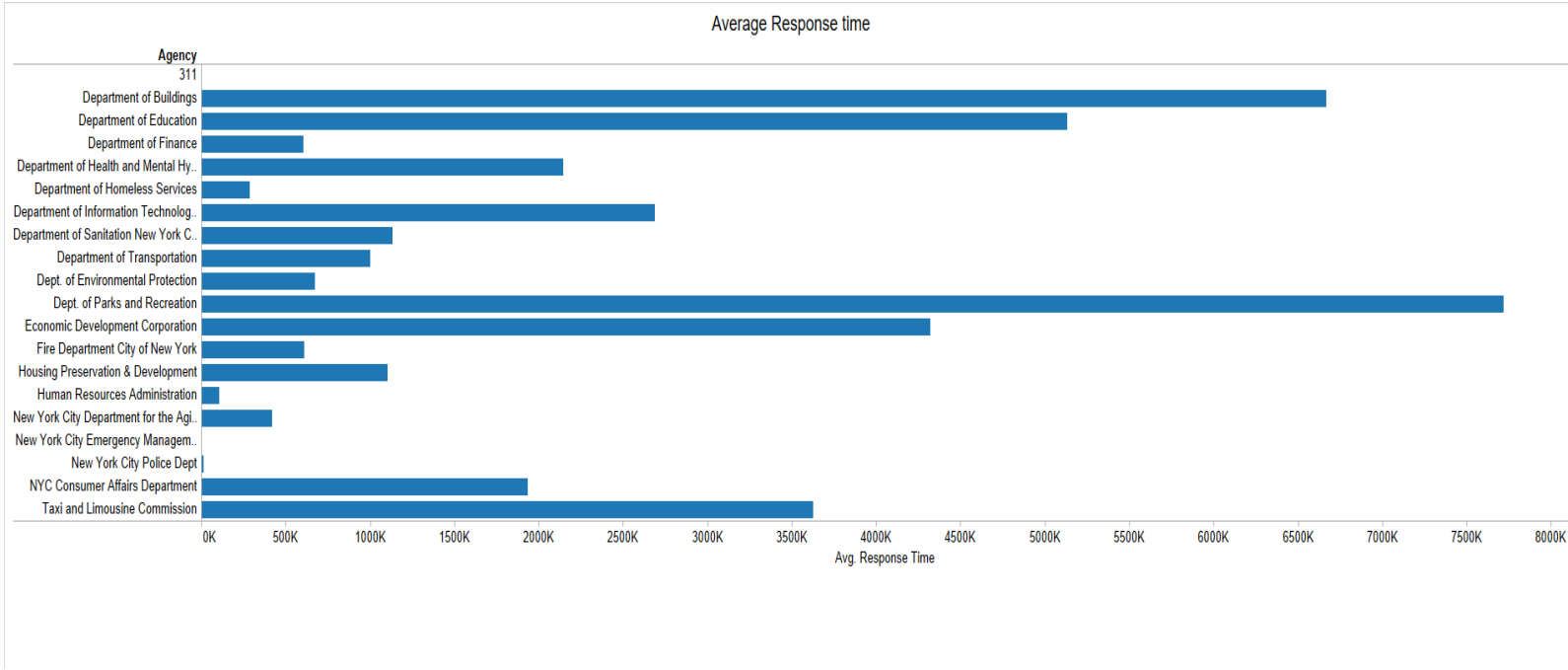
'Violation of Park Rules' ▼



Correlating number of calls with average response time

- For each agency the average response time and number of calls was computed using Pyspark.
- Then using the Statistics module the Pearson's R was computed.
- The correlation was weak and negative.
- The Pearson's R was -0.015.
- Hence there is no relationship between the response time and number of calls.

Response times



Predicting response time using random forest

- Response time for the complaints (“Closed Date” – “Created Date”) varies significantly for each agency.
- We use random forest regression algorithm to predict the response time for a complaint.
- The data was given a 70:30 (train : test) split.
- The depth was fixed as 10.
- Features considered were : “Agency”, “Agency Name”, “Complaint Type”, “Descriptor”, “Incident Zip”, “Borough”, “X-coordinate” and “Y-coordinate”.
- The feature “Agency Name” has 1644 different values and hence max bins were 1644.
- The number of trees was set to 3,7 and 16.

Predicting response time using random forest

- The accuracy was computed in terms of the RMSE (root mean square error).
- The results were as follows:
 - 3 trees : 1560 hours (65 days)
 - 7 trees : 1791 hours (74.5 days)
 - 16 trees : 1782 hours (74.3 days).
- 3 trees had the lowest RMSE of 65 days.
- The features considered didn't give very good prediction results.

Conclusion

- Complaint types varied with location, time, day and month.
- Some locations had really high complaint calls.
- Noise was the most reported complaint across all records, followed by heating.
- Response time was not correlated with the number of calls.
- Some agencies have really long response time.
- Random forest regression was used to predict the response time for a complaint.
- However, there was a high error associated with the prediction.

THANK YOU