

# DATA 698 Final Research Project

## Data Collection & Processing

Philip Tanofsky

2022-11-13

## Data Collection & Preprocessing

### Data Collection

Citi Bike provides individual bike trip data on a monthly basis available at <https://ride.citibikenyc.com/system-data>. The September 2022 bike trip data for New York City was downloaded and unzipped. The dataset contains 13 variables for each bike trip originating at a NYC-based docking station. A note on the system data page indicates trips taken by staff to service or inspect the system have been removed from the dataset. Also, any trips below 60 seconds have also been omitted. With this preprocessing by the data maintainers, the remaining trips are considered to be valid bike trips.

- **ride.id:** Unique identifier of the bike trip
- **rideable.type:** Factor variable - classic, electric, and docked
- **started.at:** Timestamp of trip departure
- **ended.at:** Timestamp of trip arrival
- **start.station.name:** Name of departure docking station
- **start.station.id:** Unique identifier of departure docking station
- **end.station.name:** Name of arrival docking station
- **end.station.id:** Unique identifier of arrival docking station
- **start.lat:** Latitude of departure location
- **start.lng:** Longitude of departure location
- **end.lat:** Latitude of arrival location
- **end.lng:** Longitude of arrival location
- **member.casual:** Factor variable for user type - member or casual

Based on the **started.at** and **ended.at** variables, four variables are derived for each bike trip.

- **day:** Day of the month
- **start.hour:** Hour of the trip departure
- **weekday:** Day of the week for the trip
- **trip.duration:** Duration of bike trip in minutes.

The NYC Open Data (free public data published by New York City agencies and partners) provides a GeoJSON file for the polygons defining each neighborhood in NYC according for the 2010 Neighborhood Tabulation Areas (NTAs). Each NTA is associated with one of the five NYC boroughs. (<https://data.cityofnewyork.us/City-Government/2010-Neighborhood-Tabulation-Areas-NTAs-/cpf4-rkhq>)

The elevation of each Citi Bike docking station is determined using the R library **elevatr** based on the latitude and longitude of each station. The elevation is defined in meters above sea level.

- **elevation:** Units above sea level
- **elev.units:** Unit of measurement for elevation

## Data Preprocessing

Upon initial inspection of the 3,507,123 bike trips in September 2022, a total of 8,012 entries do not contain an `end.station.id` and `end.station.name` listed. Of that count, 3,838 do not contain an `end.lat` and `end.lng` values. These 8,012 without a defined destination docking station will be defined as ‘Abandoned’ meaning the user did not properly dock the bike. For this purpose, the `end.lat` and `end.lng` will be removed as the abandoned bikes temporarily remove a bike from the bikeshare system. Another rider cannot rent an abandoned bike until the bike is properly docked.

Evaluation of the `end.station.id` for the NYC based bike trips includes docking stations located in New Jersey. The Citi Bike bikeshare system does include docking stations in Jersey City and Hoboken. A number of bike trips end in New Jersey which does remove the bike from the NYC-based docking stations of which this research is focused.

## Surplus Calculation

The individual bike trip information was sorted by timestamp and grouped by docking station for each 15-minute interval to count the number of bikes departing and the number bikes arriving. By subtracting the number of departures from the number of arrivals for each station for each interval, we are able to determine the running increase or decrease of bikes at the docking station. This total is defined as the variable `surplus`. A summation of the `surplus` for each docking station is calculated over the course of the month to determine which docking stations are more likely departure stations or arrival stations.

---

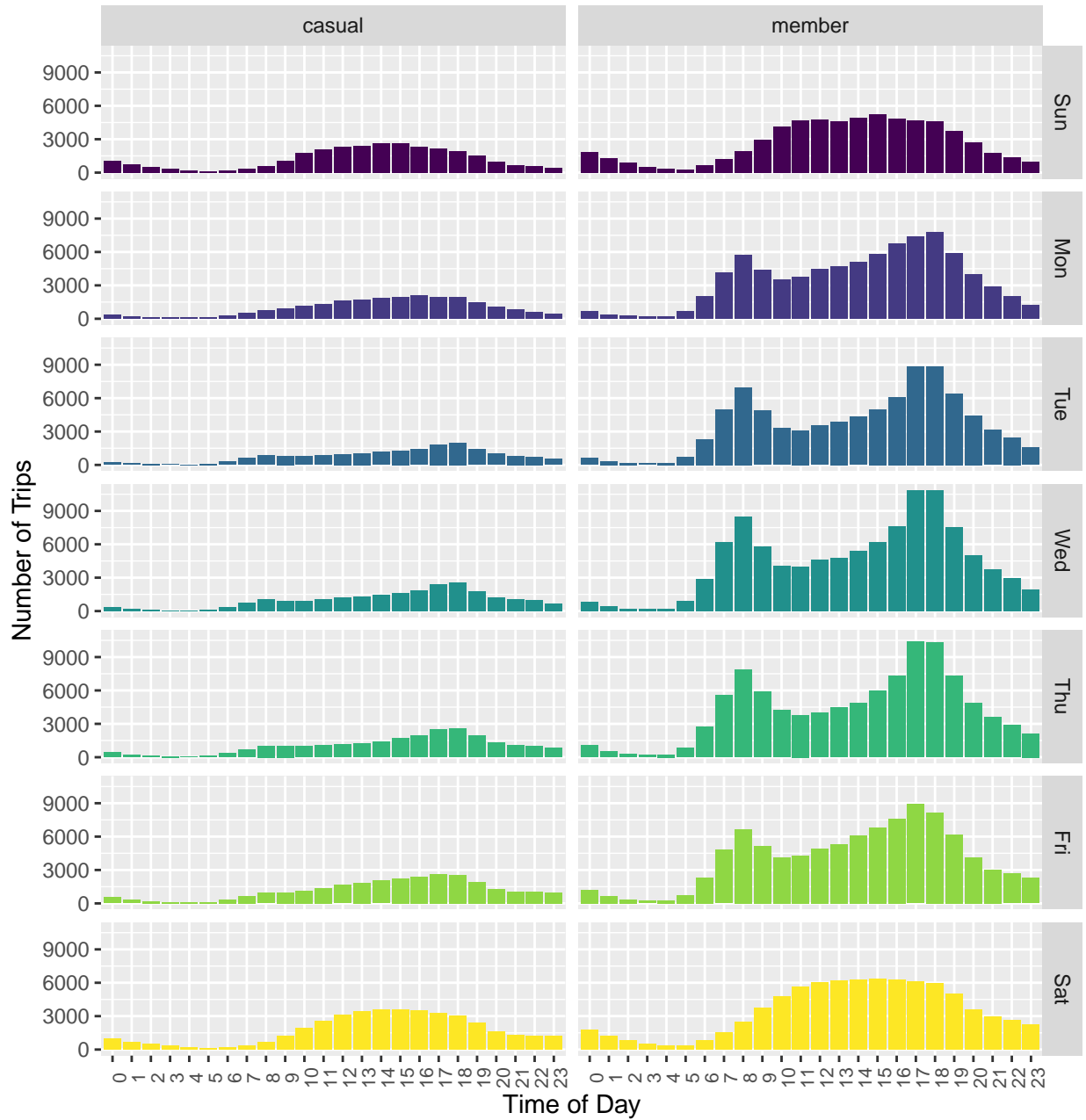
## Data Exploration & Analysis

Exploratory data analysis is performed on the Citi Bike trips for September 2022 to evaluate the patterns of bike use and identify docking stations with surplus. First, the count of bike trips are assessed to find the high volume days of the week and time of day. Next, the duration of bike trips are evaluated to assess when bikes are individually likely to be unavailable longer. Finally, the surplus of bikes by docking station and borough are analyzed to determine which areas of the New York City are more prone to having lower bike availability.

## Count of Bike Trips

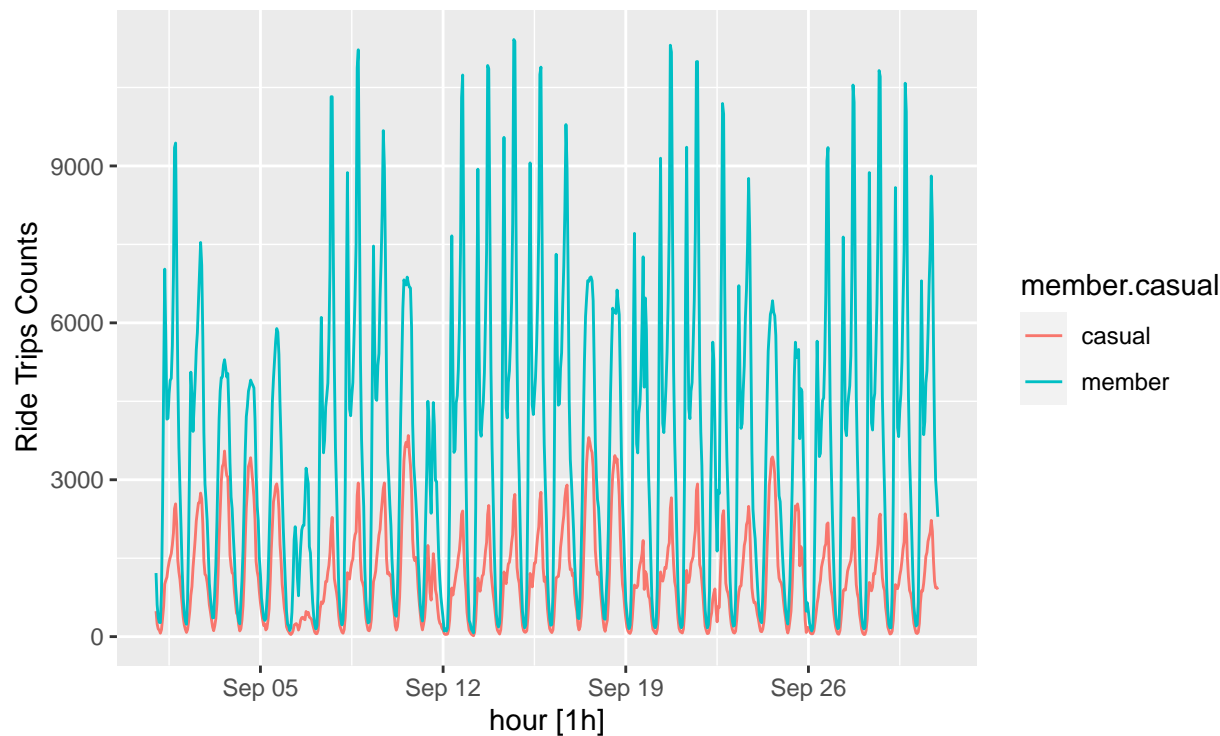
The count of bike trips are separated by user type - member and casual. With the separation by user type, two distinct patterns emerge of bike use over the course of the day and over the course of a week. The member trips are more likely to follow the workday pattern of spikes in the morning and evening as individuals are traveling to work or from work. The casual users show a pattern not necessarily indicative of the workday but instead of tourist or recreational use. The member trips account for an overall higher volume of trips.

Average Number of Bike Trips by Time, Day, and User

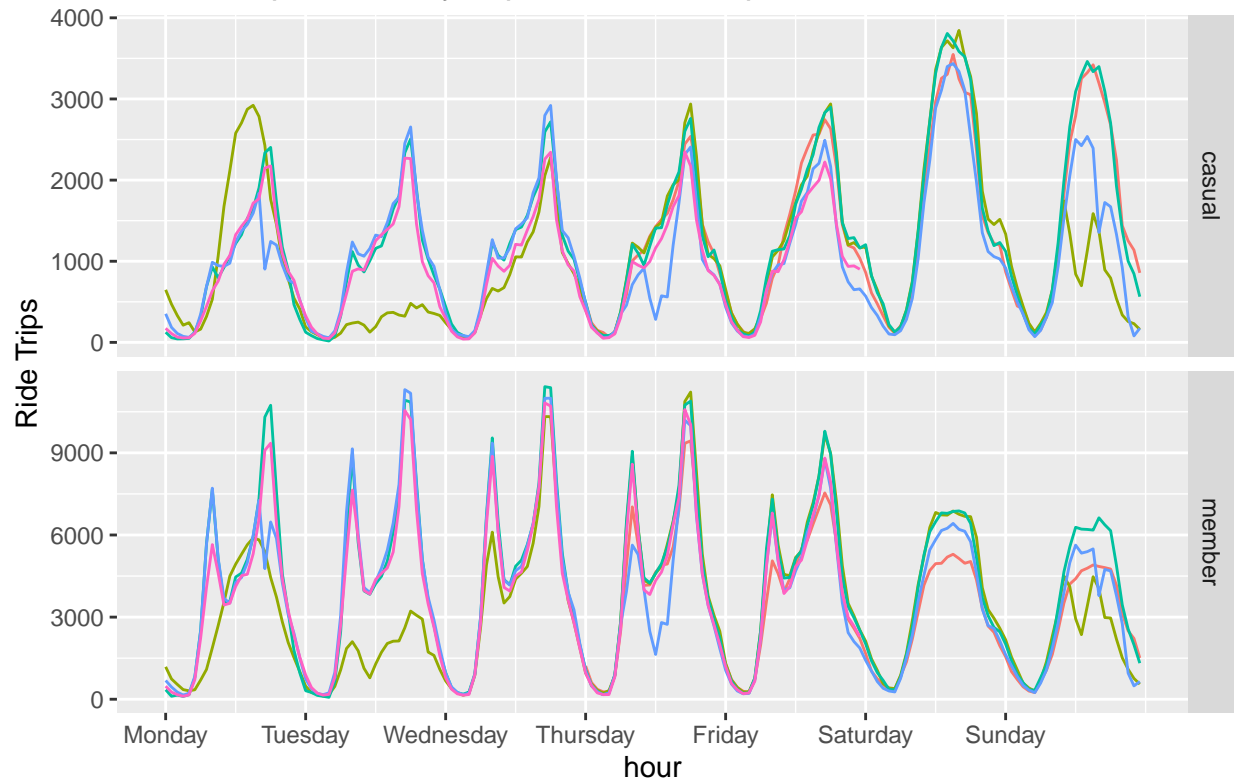


The member user counts show rush-hour spikes for members on weekdays whereas casual users do show higher counts around 5pm and 6pm on weekdays. The comparison of the two user groups also points to greater usage overall by member users. For everyday of the week, the member average member counts per hour are greater than the casual users. On weekend days, both user groups show a pattern indicative of recreational use with plateau use during the middle of the day and without distinct spikes found on the weekdays. Also, the higher usage of by member users throughout the middle of the day may indicate even if someone is a member the primary reason may not be transportation to and from work.

Ride Trips by Hour – Sept. 2022  
Citi Bike NYC



Seasonal plot: Weekly Trip Counts for Sept. 2022



The time-series chart ‘Ride Trips by Hour’ confirms the higher usage by members with a clear pattern of weekday volumes corresponding to transportation for work purposes. The casual users follow a daily pattern with increases toward the end of the week - Thursday through Saturday.

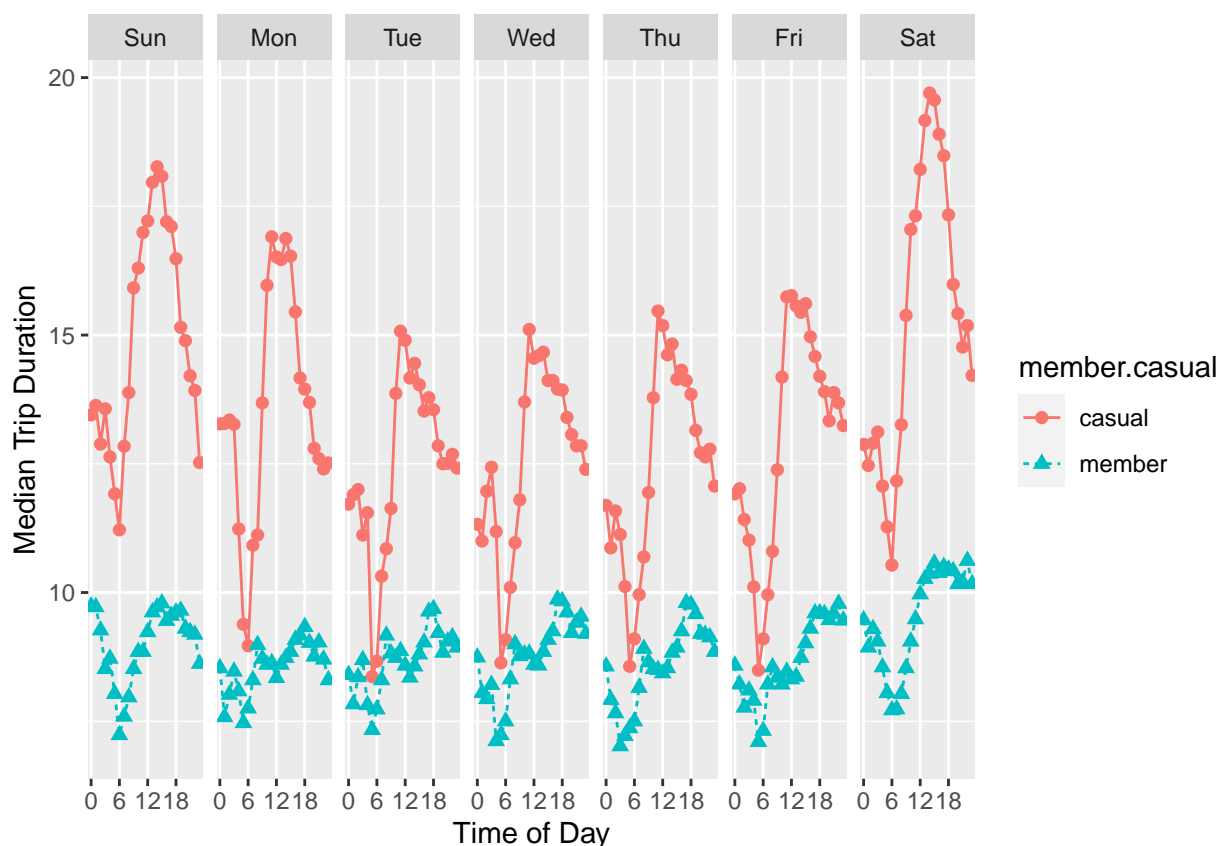
The weekly seasonal plot ‘Seasonal plot: Weekly Trip Counts’ denotes the same pattern week over week. Based on the plot, the bike trips show consistency from week to week based on member users utilizing the bikes for work and casual users renting the bikes for recreational purposes.

By count of trips per day of the week, Wednesday appears the highest volume day. Given the data comes from September 2022, we believe this observation is a result of pandemic return-to-office policies in which individuals more likely to return to office during the middle of the week instead of Monday or Friday for those with hybrid schedules.

Note: On Tuesday, Sept, 6 2022, the weather consisted of rain the entire day and thus a clear drop in use is evident in both user groups.

## Duration of Bike Trips

Next, we evaluate the duration of bike trips to better understand longevity of individual bike unavailability along with reason for trip. As expected, the average bike trip for all users and all times are below 30 minutes as the rental defined time is 30 minutes with users incurring additional fees beyond the base time limit.



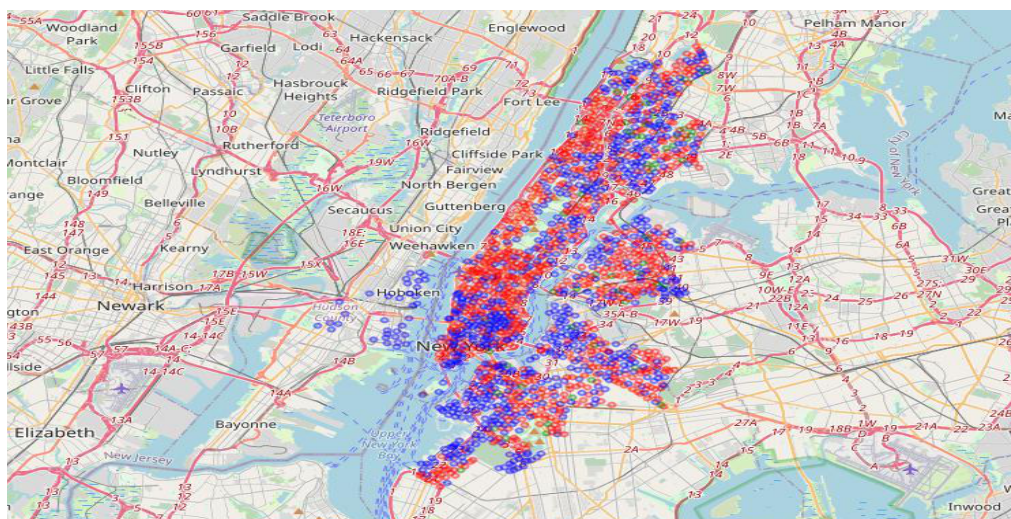
The chart of average trip duration by user type shows a clear distinction between the user types. Member users tend to average bike trips of 10 minutes or less throughout the week, whereas the casual users have a greater variance in average duration based on time of day and day of the week. The member users only average trips greater than 10 minutes during the afternoon and evening on Saturday while casual users typically average even longer trips of greater than 17 minutes on the weekend. Casual users tend to reach of peak of greater than 15 minutes on average everyday of the week, particularly around lunch on weekdays.

We observe that member users tend to have longer average trips during the morning and evening rush hours on weekdays with a slight deviation on Friday evening as the weekend starts. The afternoon and evening hours of Saturday and Sunday show longer trips for members as likely the result of recreational trips.

## Docking Station Surplus

After constructing a matrix of bike arrivals and departures for each NYC-based docking station for every 15-minute interval of September 2022, we calculate the overall surplus or shortage of bikes. The researchers note, the shortage of bikes can logically not fall below zero without the introduction of rebalancing throughout the bikeshare system. The following analysis confirms the rebalancing of bikes across the docking stations occurs to ensure popular departure docking stations have bikes available despite the dearth of bike trip arrivals.

### Overall Monthly Surplus by Docking Station



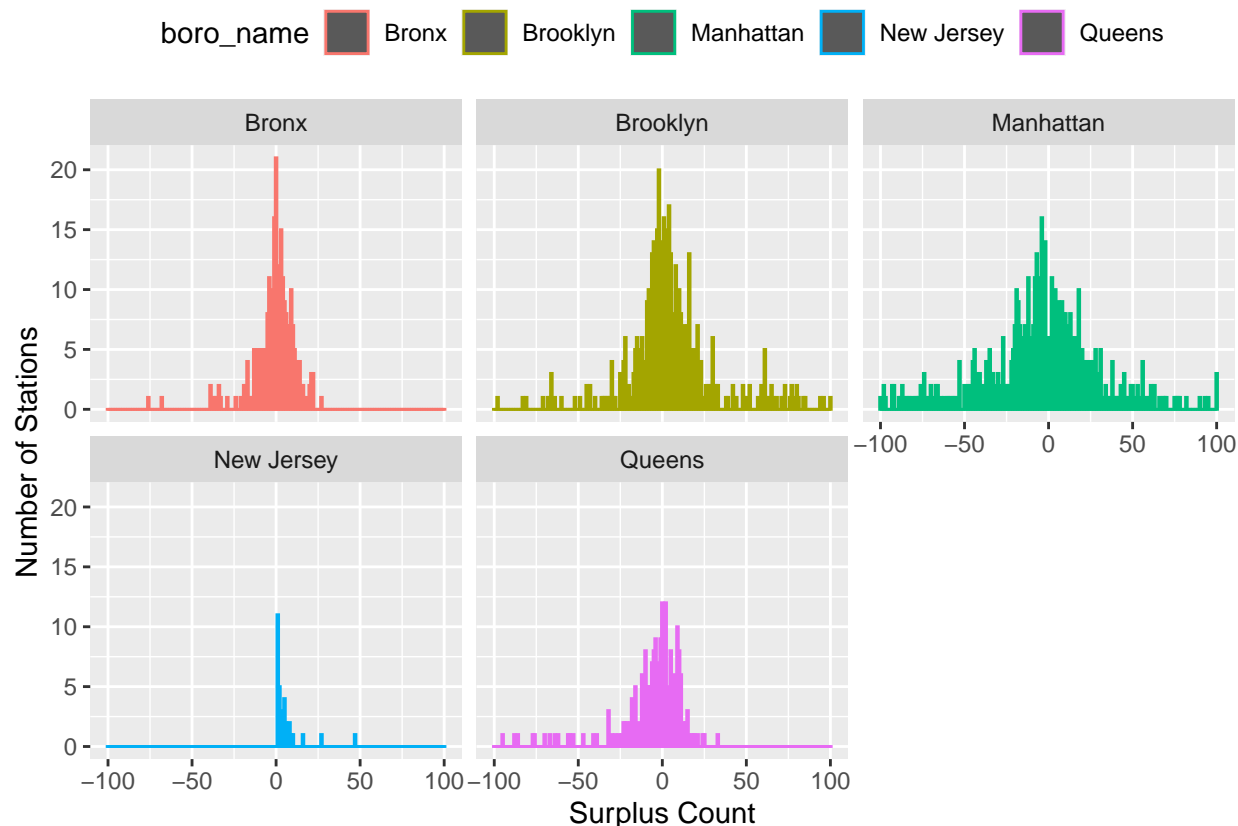
The plot of docking stations across the New York City maps the location of every docking station along with a color to indicate a net positive, negative or even amount of bikes for the given month. The blue stations are net positive, and red stations indicate net negative while green stations are even for the entire month. Of the 1656 docking stations represented, 799 ended with a surplus, 804 with a shortage, and 53 with an even count.

The plot does indicate several blue-colored docking stations in New Jersey. As the dataset does not contain any bike trips originating in NYC, a number of trips have an ending station in Jersey City or Hoboken. The New Jersey based stations are guaranteed to be blue based on that dataset without any departing trips from those docking stations.

The dataset does encompass 472,920 valid combinations of departure and arrival docking stations. Of the valid combinations 15 of the top 20 combinations are the same departure and arrival docking stations, indicative of recreational bike trips. Of the aforementioned valid combinations, 471,313 denote travel between two different docking stations.

## Surplus by borough

With almost an even number of docking stations with a net shortage and a net surplus, we evaluate the shortage and surplus by borough.



The plot denotes a near normal distribution for the four boroughs included in the dataset - the Bronx, Brooklyn, Manhattan, and Queens. As noted previously, the dataset contains only trips originating in New York City, so the plot for New Jersey only indicates docking stations surpluses. The balanced distribution across the four NYC boroughs likely demonstrates the bike trips are contained within a borough. With average duration less than 10 minutes for member users and less than 20 minutes for casual users, the distance traveled for all bike trips is likely less than three miles. The 30-minute base rental rate dictates a limited travel distance which inhibits users from traveling across boroughs via Citi Bike.

## Closing Remarks

The exploratory analysis of the Citi Bike dataset for NYC-based trips in September 2022 shows a consistent pattern of use dependent on time of day and day of the week. The pattern persists week over week for the given month. The majority of users are member which indicates the availability of bikes will be dependent on work schedules during weekdays. The surplus and shortage count of bikes by docking station denotes the uneven direction of bikes in some sections of New York City. The shortage (and surplus) counts confirm the practice of rebalancing by Citi Bike to ensure availability of bikes. The duration of bike trips may not play as large a factor in bike availability compared to time of day, day of week, and directional flow of bikes throughout the city.

## Appendix with Code

```
# Required packages
library(tidyverse)
library(ggplot2)
library(skimr)
library(lubridate)
library(fpp3)
library(assertthat)
library(igraph)
library(ggraph)
library(ggmap)
library(leaflet)
library(rgdal)
library(RColorBrewer)
library(jpeg)
library(data.table)

# Create additional columns of pertinence
# weekday, day of the month, trip duration in minutes, start hour
citibike <- fread("data/202209-citibike-tripdata.csv", data.table=FALSE, check.names=TRUE) %>%
#citibike <- read.csv("data/202209-citibike-tripdata.csv", check.names=TRUE) %>%
  mutate(day = factor(mday(ymd_hms(started_at))),
         start.hour=factor(hour(ymd_hms(started_at))),
         weekday = lubridate::wday(ymd_hms(started_at), label=TRUE, abbr=TRUE),
         trip.duration = as.numeric(difftime(ended_at,started_at,units="mins")),
         member_casual = factor(member_casual)) %>%
  rename(ride.id=ride_id, rideable.type=rideable_type, started.at=started_at,
         ended.at=ended_at, start.station.name=start_station_name, start.station.id=start_station_id,
         end.station.name=end_station_name, end.station.id=end_station_id, start.lat=start_lat,
         start.lng=start_lng, end.lat=end_lat, end.lng=end_lng, member.casual=member_casual)

# Figure out abandoned and label them
citibike$end.station.name[citibike$end.station.name == ''] <- "Abandoned"
citibike$end.station.id[citibike$end.station.id == ''] <- "ABAN"

# Output DF if needed
head(citibike)

# Trip by weekday by segment by time of day
citibike %>%
  group_by(day, member.casual, start.hour) %>%
  summarize(n=n(),
            weekday=weekday[1]) %>%
  group_by(weekday, member.casual, start.hour) %>%
  summarize(n.m=mean(n)) %>%
  ggplot(aes(x=start.hour, y=n.m, fill=weekday)) +
  geom_bar(stat='identity') +
  labs(x='Time of Day',
       y='Number of Trips',
       title='Average Number of Bike Trips by Time, Day, and User') +
  facet_grid(weekday~member.casual) +
  theme(axis.text.x = element_text(size=8, angle=90),
```



```

    legend.position = 'none')

citibike$started.at.ts <- as_datetime(as.character(citibike$started.at))

citibike_by_hour <- citibike %>%
  mutate(hour=lubridate::floor_date(started.at.ts, "1 hour")) %>%
  group_by(hour, member.casual) %>%
  summarize(cnt=n())
#citibike_by_hour

citibike_by_hour_ts <- citibike_by_hour %>%
  as_tsibble(index=hour, key=c(member.casual))
#citibike_by_hour_ts

autoplot(citibike_by_hour_ts, cnt) +
  labs(title = "Ride Trips by Hour - Sept. 2022",
        subtitle = "Citi Bike NYC",
        y = "Ride Trips Counts")
citibike_by_hour_ts %>%
  gg_season(cnt, period = "week") +
  labs(y = "Ride Trips",
        title = "Seasonal plot: Weekly Trip Counts for Sept. 2022")
citibike %>%
  filter(member.casual %in% c('member', 'casual')) %>%
  group_by(weekday, start.hour, member.casual) %>%
  summarize(med.duration=median(trip.duration)) %>%
  ggplot(aes(x=start.hour, y=med.duration, group=member.casual,
             color=member.casual, linetype=member.casual, shape=member.casual)) +
  geom_point(size=2) +
  geom_line(size=0.5) +
  facet_wrap(~weekday, nrow=1) +
  labs(x='Time of Day',
        y='Median Trip Duration') +
  scale_x_discrete(breaks=c(0,6,12,18))
# Create table of start to end station IDs
stations_cols <- citibike %>%
  select(start.station.id, end.station.id)
stations_table <- as.data.frame((table(stations_cols)))

stations_table <- stations_table %>%
  filter(Freq > 0)

stations_table_order <- stations_table[order(-stations_table$Freq),]

stations_table_dif <- stations_table_order %>%
  filter(as.character(start.station.id) != as.character(end.station.id))

# Create table of start to end station IDs
stations_cols <- citibike %>%
  select(start.station.id, end.station.id)
stations_table <- as.data.frame((table(stations_cols)))

stations_table <- stations_table %>%

```

```

filter(Freq > 0)

stations_table_order <- stations_table[order(-stations_table$Freq),]

stations_table_dif <- stations_table_order %>%
  filter(as.character(start.station.id) != as.character(end.station.id))

stations_table_dif1 <- stations_table_dif
stations_table_dif2 <- stations_table_dif

# Added all=TRUE to account for one-sided counts
stations_table_dif_merge <-
  merge(stations_table_dif1,
        stations_table_dif1,
        by.x=c('start.station.id','end.station.id'),
        by.y=c('end.station.id','start.station.id'), all = TRUE)

# Set the NA (one-sided trips) to count of 0
stations_table_dif_merge[is.na(stations_table_dif_merge)] <- 0

stations_table_dif_merge$surplus <- stations_table_dif_merge$Freq.y - stations_table_dif_merge$Freq.x
stations_table_dif_merge <- stations_table_dif_merge[order(-stations_table_dif_merge$surplus),]

# Remove rows with start station id equal to ABAN for abandoned, those are a result of the merge all=TRUE
stations_table_dif_merge <- stations_table_dif_merge %>% filter(start.station.id != 'ABAN')

# Want to identify the surplus (or not) by station for the month
station_surplus_count <- stations_table_dif_merge %>%
  group_by(start.station.id) %>%
  summarize(surplus.sum=sum(surplus))

station_surplus_count <- station_surplus_count[order(-station_surplus_count$surplus.sum),]
colnames(station_surplus_count)[1] <- "station.id"

# Extract just the end station Id, lat, log
# because this had the higher count from the initial dataset, going with end_station_id
end_station_info <- citibike %>%
  select(end.station.id, end.lng, end.lat)

end_station_info <- end_station_info[!duplicated(end_station_info$end.station.id),]

# Now I want the coordinates of all those station_ids
station_surplus_count_coords <- merge(x = station_surplus_count, y = end_station_info, by.x = 'station.id', by.y = 'end.station.id')

colnames(station_surplus_count_coords)[3] <- "lng"
colnames(station_surplus_count_coords)[4] <- "lat"

station_surplus_count_coords <- station_surplus_count_coords %>%
  add_row(station.id = "ABAN", surplus.sum=1363, lng=-73.99, lat=40.67)

# Using basemaps for NYC
m <- leaflet(data=station_surplus_count_coords) %>%

```

```

# setView(zoom=12) %>%
addTiles() %>%
addCircleMarkers(
  ~lng, ~lat,
  popup=~as.character(station.id),
  label=~as.character(station.id),
  radius=.5,
  color = ~ifelse(surplus.sum >= 1, 'blue',
                  ifelse(surplus.sum == 0, 'green', 'red'))
)

# Display map
#m
img <- readJPEG("stations_surplus_sum.jpg")
plot(1:10,ty="n", axes = 0, xlab='', ylab='', main='Overall Monthly Surplus by Docking Station')
rasterImage(img,-1,-1,12,12)
stations_with_elevation <- read.csv('stations_with_elevation.csv', row.names = 1, header= TRUE)

stations_with_boro_hood <- read.csv('stations_with_boro_and_hood.csv', row.names = 1, header= TRUE)

# Combine the elevation, borough, neighborhood, and September surplus
# First let's trim the DF for elevation
stations_with_elevation_trim <- stations_with_elevation %>%
  select(short_name, station_id, elevation, elev_units)

stations_with_boro_hood_trim <- stations_with_boro_hood %>%
  select(short_name, name, station_id, capacity, ntaname, boro_name, lon, lat)

stations_attr_trim <-
  merge(stations_with_elevation_trim,
        stations_with_boro_hood_trim,
        by.x=c('short_name'),
        by.y=c('short_name'), all = TRUE)

stations_attr_trim <- stations_attr_trim %>%
  select(-station_id.y)

colnames(stations_attr_trim)[colnames(stations_attr_trim) == 'station_id.x'] <- 'station_id'
stations_attr_trim[c("boro_name")][is.na(stations_attr_trim[c("boro_name")])] <- "New Jersey"

stations_attr_trim <-
  stations_attr_trim %>%
  mutate(ntaname = ifelse(startsWith(short_name, "JC"), "Jersey City", ntaname))

stations_attr_trim <-
  stations_attr_trim %>%
  mutate(ntaname = ifelse(startsWith(short_name, "HB"), "Hoboken", ntaname))

stations_attr_trim_sur <-
  merge(stations_attr_trim,
        station_surplus_count,

```

```

    by.x=c('short_name'),
    by.y=c('station.id'), all.x = TRUE)

stations_attrs_trim_sur <- stations_attrs_trim_sur %>%
  filter(!is.na(surplus.sum))

stations_attrs_trim_sur %>%
  filter(surplus.sum >= -100 & surplus.sum <= 100) %>%
  ggplot(aes(x=surplus.sum, color=boro_name)) +
  theme(legend.position="top") +
  geom_histogram(bins=201) +
  facet_wrap(~boro_name) +
  labs(x = 'Surplus Count',
       y = 'Number of Stations')
# This removed the duplicate rows by transposing the start and end station ids
stations_table_dif_merge_temp <- stations_table_dif_merge %>% select(start.station.id, end.station.id)

stations_for_graph <- stations_table_dif_merge_temp[!duplicated(lapply(as.data.frame(t(stations_table_dif_merge_temp[,1:2])), FUN=function(x) {
  x[1]
}))

g_stations <- graph_from_data_frame(stations_for_graph, directed=FALSE, vertices=station_surplus_count_coords)

edges_for_plot <- stations_for_graph %>%
  inner_join(station_surplus_count_coords %>% select(station.id, lng, lat), by=c('start.station.id' = 'station.id', 'end.station.id' = 'station.id')) %>%
  rename(x=lng, y=lat) %>%
  inner_join(station_surplus_count_coords %>% select(station.id, lng, lat), by=c('end.station.id' = 'station.id', 'start.station.id' = 'station.id')) %>%
  rename(xend=lng, yend=lat)

#assert_that(nrow(edges_for_plot) == nrow(stations_for_graph))

# https://yihui.org/en/2018/09/code-appendix/

```

LaTeX help <https://tex.stackexchange.com/questions/10684/vertical-space-in-lists>