





```
In [131]: trips18_clean.ww.groupby.value_counts().index.categories
Out[131]: InDex(['Monday', 'Tuesday', 'Wednesday', 'Thursday', 'Friday', 'Saturday',
          'Sunday'],
          dtype='object')
```

```
In [132]: # Final df clean
trips18_clean.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 1695316 entries, 6 to 1863719
Data columns (total 23 columns):
#   Column                                Dtype
---  --
0   duration_sec                          int64
1   start_time                            datetime64[ns]
2   end_time                              datetime64[ns]
3   start_station_id                      int32
4   start_station_name                    object
5   start_station_latitude                float64
6   start_station_longitude               float64
7   end_station_id                        int32
8   end_station_name                      object
9   end_station_latitude                  float64
10  end_station_longitude                 float64
11  bike_id                               int64
12  user_type                             category
13  member_birth_year                    int32
14  member_gender                         category
15  bike_share_for_all_trip               bool
16  start_petro_area                      category
17  end_metro_area                        category
18  month                                 category
19  weekday                               category
20  hour                                  int32
21  member_age                            int64
22  age_group                             category
dtypes: bool(1), category(8), datetime64[ns](2), float64(4), int32(3), int64(3), object(2)
memory usage: 189.2+ MB
```

```
In [133]: # trips18.info()

In [134]: # # Issue 2: add new columns for trip duration in minute, trip start date in yyyy-mm-dd format, trip s
          # start hour of the day, day of week and end time
trips18_clean['duration_minute'] = trips18_clean['duration_sec']/60
trips18_clean['start_date'] = trips18_clean.start_time.dt.strftime('%Y-%m-%d')
trips18_clean['start_hourofday'] = trips18_clean.start_time.dt.strftime('%H')
trips18_clean['start_dayofweek'] = trips18_clean.start_time.dt.strftime('%A')
trips18_clean['start_month'] = trips18_clean.start_time.dt.strftime('%B')
```

```
Out[134]:
   duration_sec  start_time  end_time  start_station_id  start_station_name  start_station_latitude  start_station_longitude  end_station
98      453  2018-01-31  2018-02-01  110  Street Park (17th St at Folsom St)  37.783708  -122.415204  1
7       180  2018-01-31  2018-01-31  81  Berry St at 4th St  37.775880  -122.393170
8       996  2018-01-31  2018-01-31  134  Valencia St at 24th St  37.752428  -122.420628
9       825  2018-01-31  2018-01-31  305  Ryland Park  37.347275  -121.895617  3
11      432  2018-01-31  2018-01-31  89  Division St at Potrero Ave  37.769218  -122.407646

5 rows x 28 columns
```

```
In [135]: # Issue 3: add a new column calculating riders' age from 'member_birth_year'
trips18_clean['member_age'] = 2019 - trips18_clean['member_birth_year']
trips18_clean.describe()

Out[135]:
   duration_sec  start_station_id  start_station_latitude  start_station_longitude  end_station_id  end_station_latitude  end_station_longitude
count  1.695316e+06  1.695316e+06  1.695316e+06  1.695316e+06  1.695316e+06  1.695316e+06  1.695316e+06
mean    6.734631e+02  1.204515e+02  3.776863e+01  -1.223512e+02  1.186763e+02  3.776873e+01  -1.223507e+02
std     5.418813e+02  1.001461e+02  1.016470e+01  1.194266e+01  1.001396e+02  1.014922e+01  -1.128235e+02
min      6.10000e+01  3.00000e+00  3.726331e+01  -1.224737e+02  3.00000e+00  3.726331e+01  -1.224737e+02
25%    3.43000e+02  3.60000e+01  3.726331e+01  -1.224116e+02  3.00000e+01  3.777106e+01  -1.224094e+02
50%    5.40000e+02  8.90000e+01  3.777106e+01  -1.223974e+02  8.90000e+01  3.777106e+01  -1.223974e+02
75%    8.31000e+02  1.86000e+02  3.77625e+01  -1.223876e+02  1.86000e+02  3.776278e+01  -1.223846e+02
max    6.00000e+03  3.81000e+02  3.780222e+01  -1.218333e+02  3.81000e+02  3.78022e+01  -1.218333e+02
```

```
In [136]: # # Issue 4: filter out outlier ages from visually examination of the distribution above
          # # Issue 5: cast 'member_birth_year' and 'member_age' to integer instead of float type
trips18_clean = trips18_clean.query('member_age <= 70')
trips18_clean['member_birth_year'] = trips18_clean['member_birth_year'].astype('int')
trips18_clean['member_age'] = trips18_clean['member_age'].astype('int')
trips18_clean.info()
Out[136]:
<class 'pandas.core.frame.DataFrame'>
Int64Index: 1695316 entries, 6 to 1863719
Data columns (total 28 columns):
#   Column                                Non-Null Count  Dtype
---  --
0   duration_sec                          1695316 non-null  int64
1   start_time                            1695316 non-null  datetime64[ns]
2   end_time                              1695316 non-null  datetime64[ns]
3   start_station_id                      1695316 non-null  int32
4   start_station_name                    1695316 non-null  object
5   start_station_latitude                1695316 non-null  float64
6   start_station_longitude               1695316 non-null  float64
7   end_station_id                        1695316 non-null  int32
8   end_station_name                      1695316 non-null  object
9   end_station_latitude                  1695316 non-null  float64
10  end_station_longitude                 1695316 non-null  float64
11  bike_id                               1695316 non-null  int64
12  user_type                             1695316 non-null  category
13  member_birth_year                    1695316 non-null  int32
14  member_gender                         1695316 non-null  category
15  bike_share_for_all_trip               1695316 non-null  bool
16  start_petro_area                      1695316 non-null  category
17  end_metro_area                        1695316 non-null  category
18  month                                 1695316 non-null  category
19  weekday                               0 non-null       category
20  hour                                  1695316 non-null  int32
21  member_age                            1695316 non-null  int32
22  age_group                             1695316 non-null  category
23  duration_minute                       1695316 non-null  object
24  start_date                            1695316 non-null  object
25  start_hourofday                       1695316 non-null  object
26  start_dayofweek                       1695316 non-null  object
27  start_month                           1695316 non-null  object
dtypes: bool(1), category(8), datetime64[ns](2), float64(5), int32(4), int64(2), object(6)
memory usage: 247.6+ MB
```

**What is the structure of your dataset?**

The original combined data contains approximately 1,693,719 individual trip records with 16 variables collected. These trips depend on various factors that is Duration of trip, Start Date, Time and End Date time, of the trip, User Type etc.

**Main features of my interest in the Dataset :**

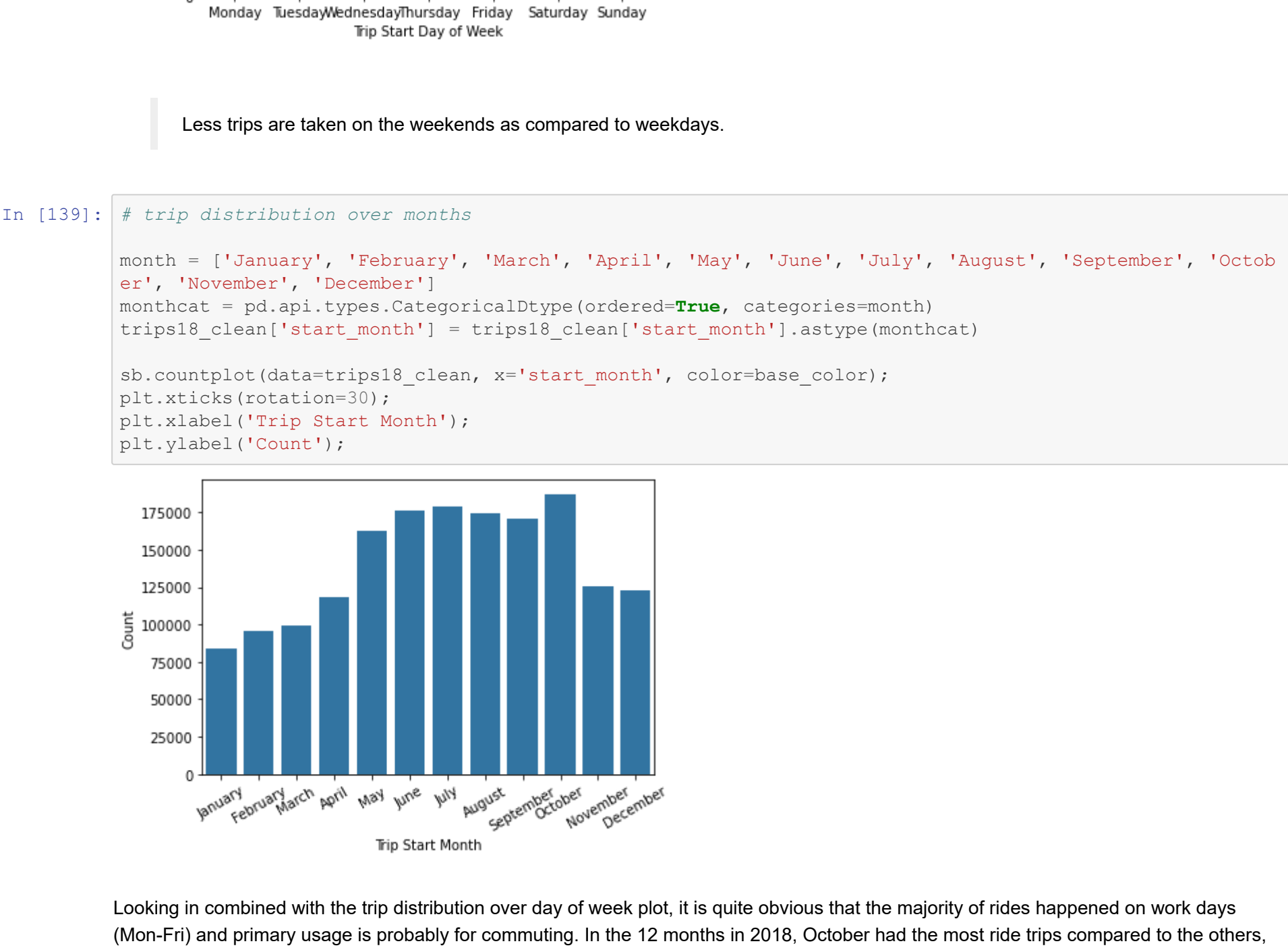
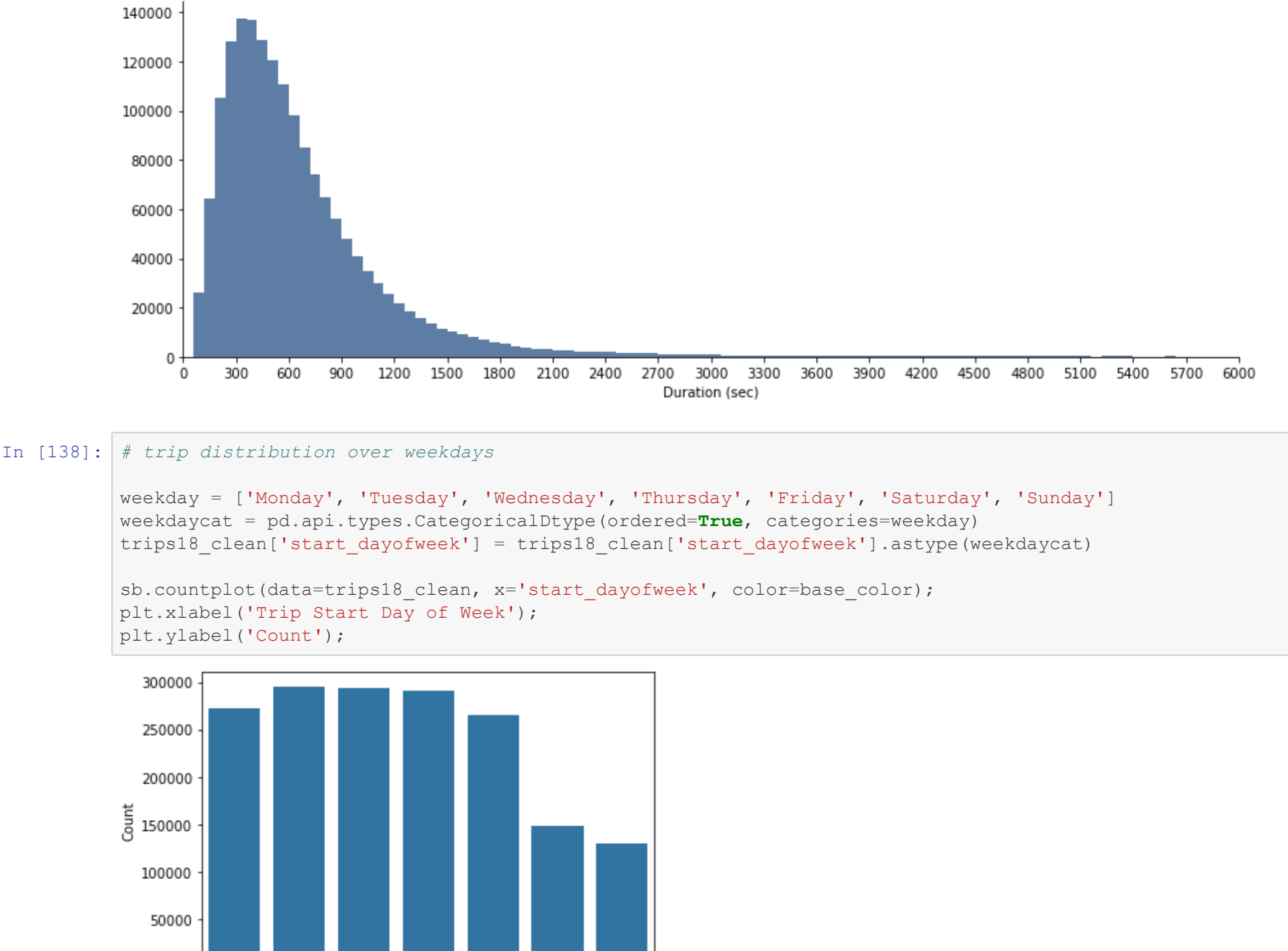
I'm most interested in figuring out what features are best for predicting the more number of bike trips in the dataset.

**Features in the dataset which will help me support our investigation :**

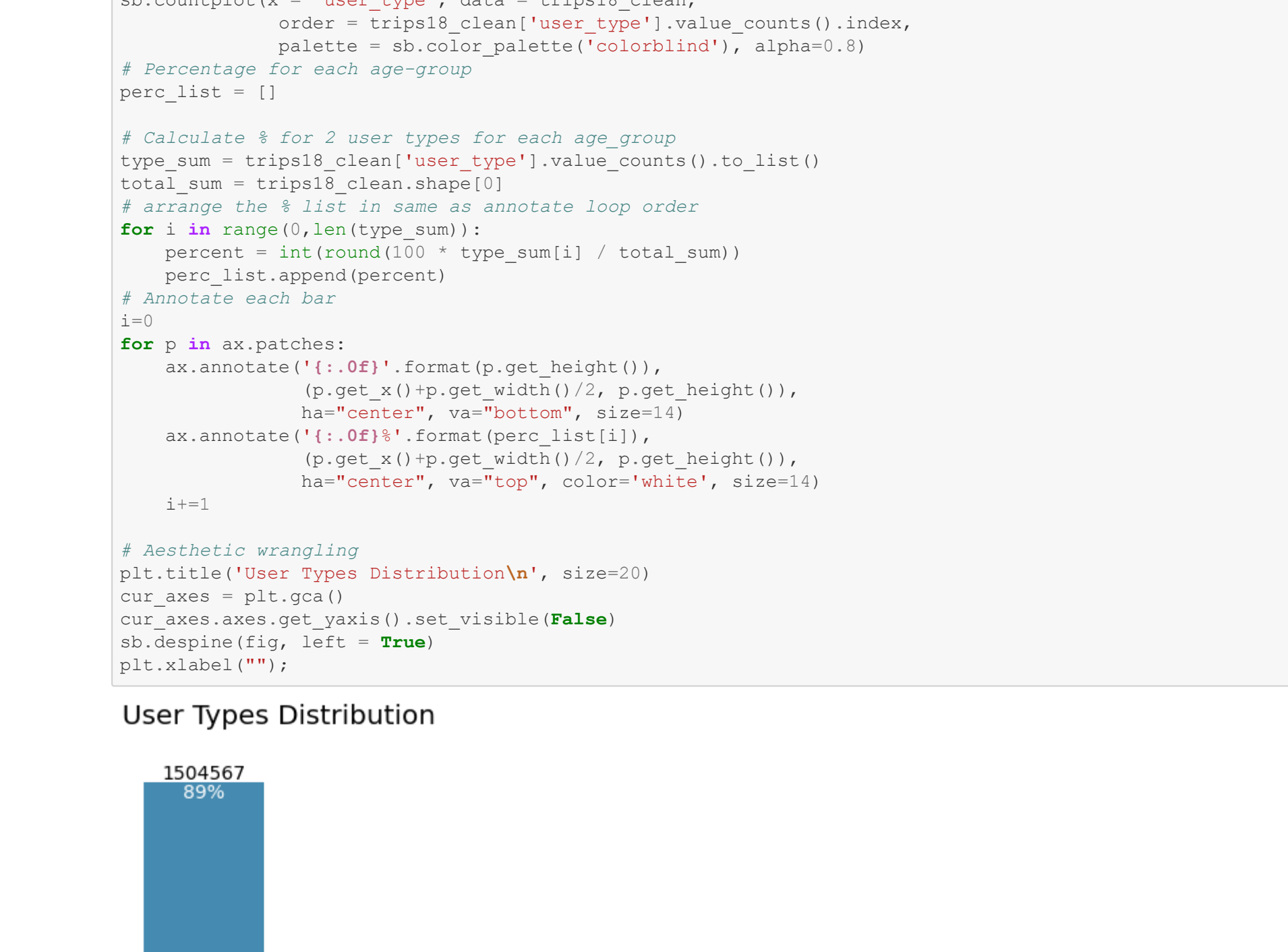
According to the dataset, the features like Start Date, Time, Month, Duration of Trip and User Type will have the strongest effect on Bike Trips.

**Univariate Exploration**

A series of plots to first explore the trips distribution over hour-of-day, day-of-week and month.



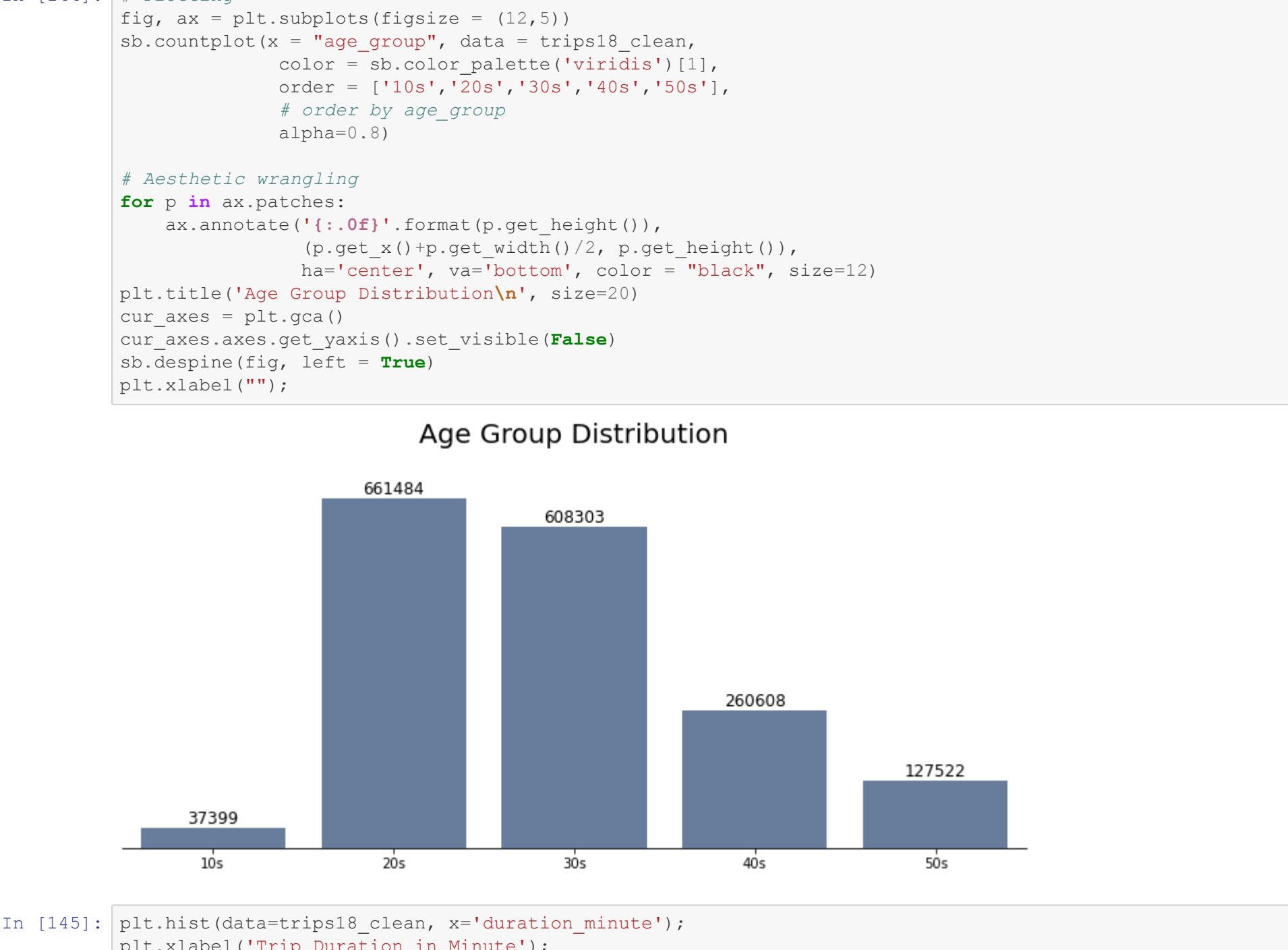
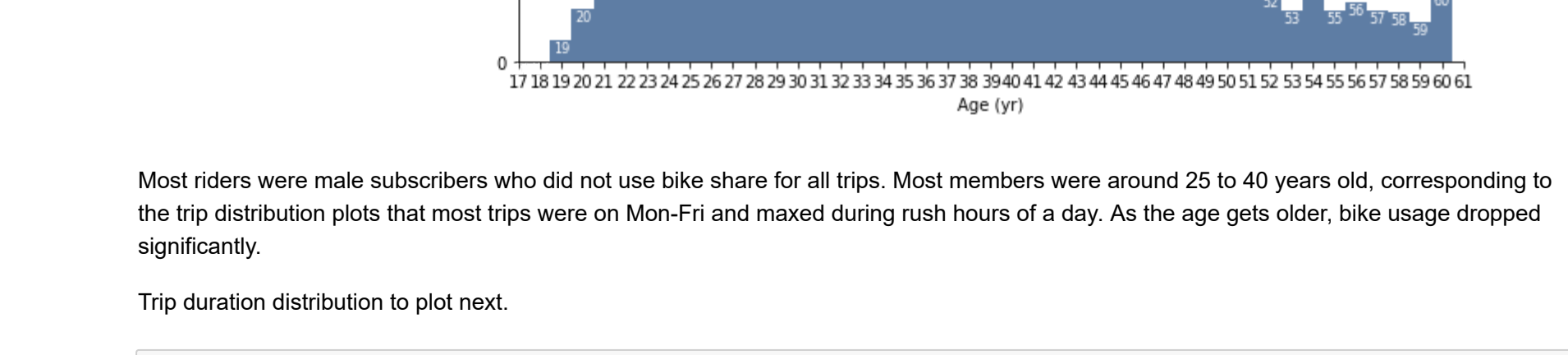
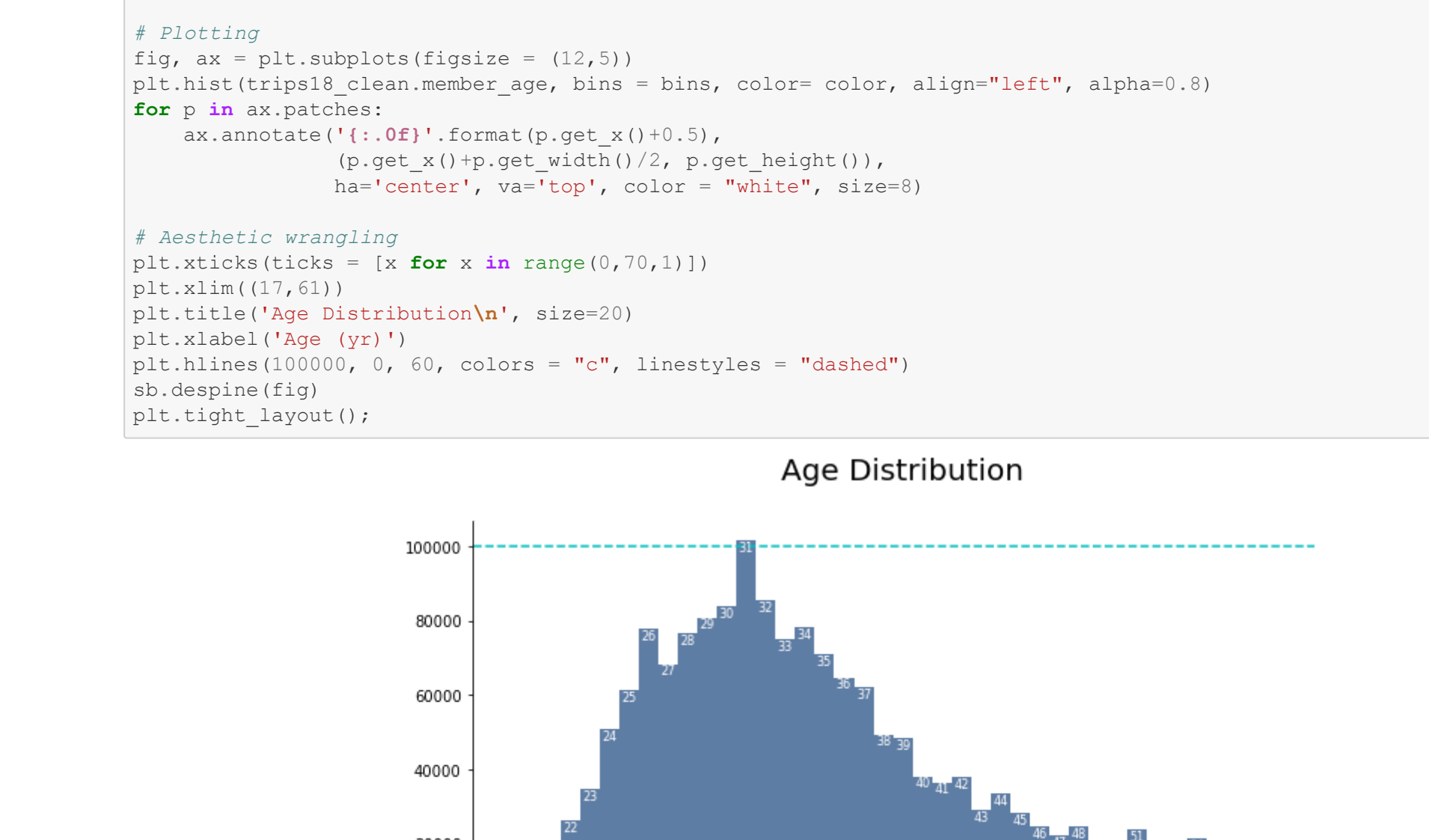
Less trips are taken on the weekends as compared to weekdays.



Looking in combined with the trip distribution over day of week plot, it is quite obvious that the majority of rides happened on work days (Mon-Fri) and primary usage is probably for commuting. In the 12 months in 2018, October had the most ride trips compared to the others, but overall it was the most popular during summer time (May-Sept), probably due to the weather in the area.

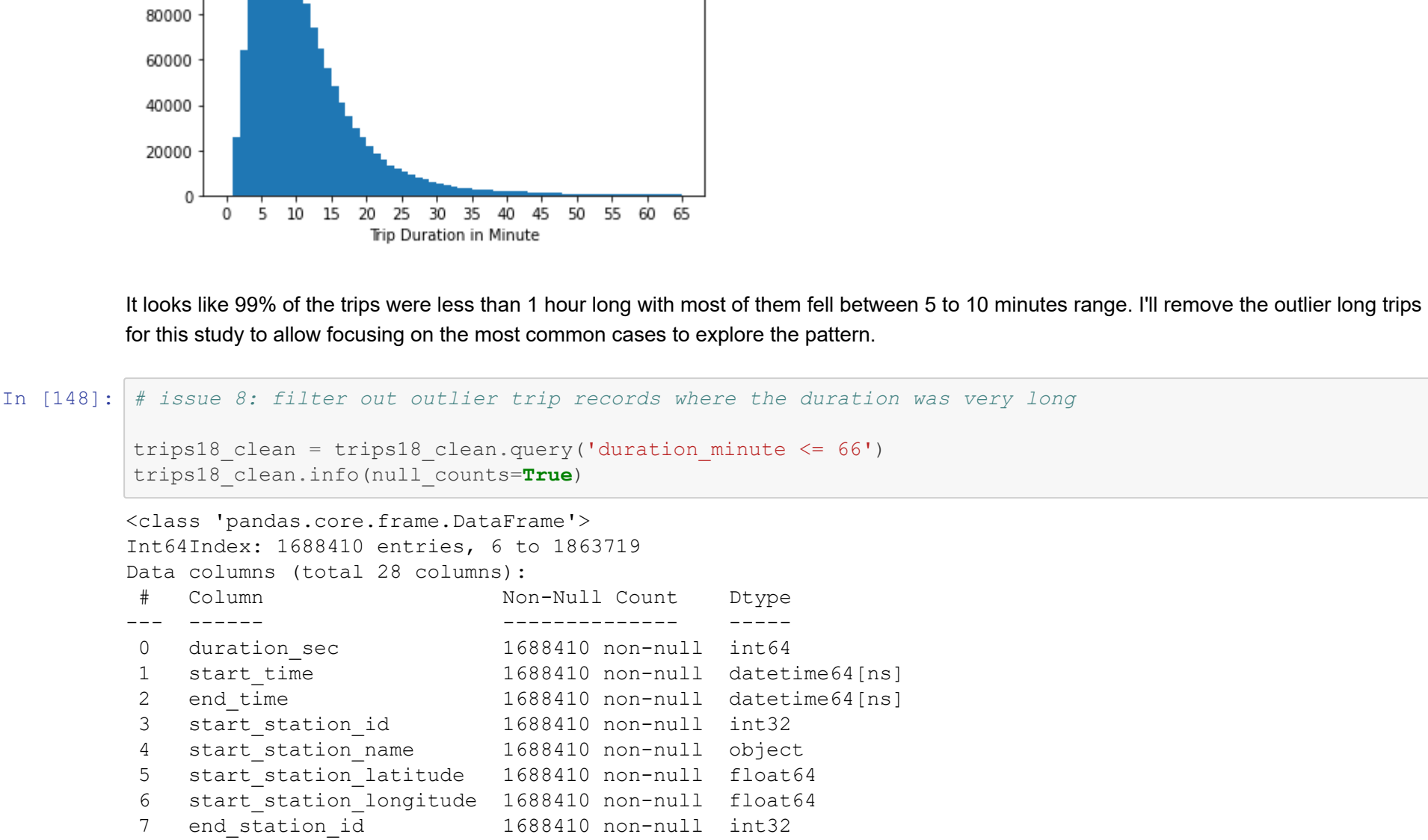
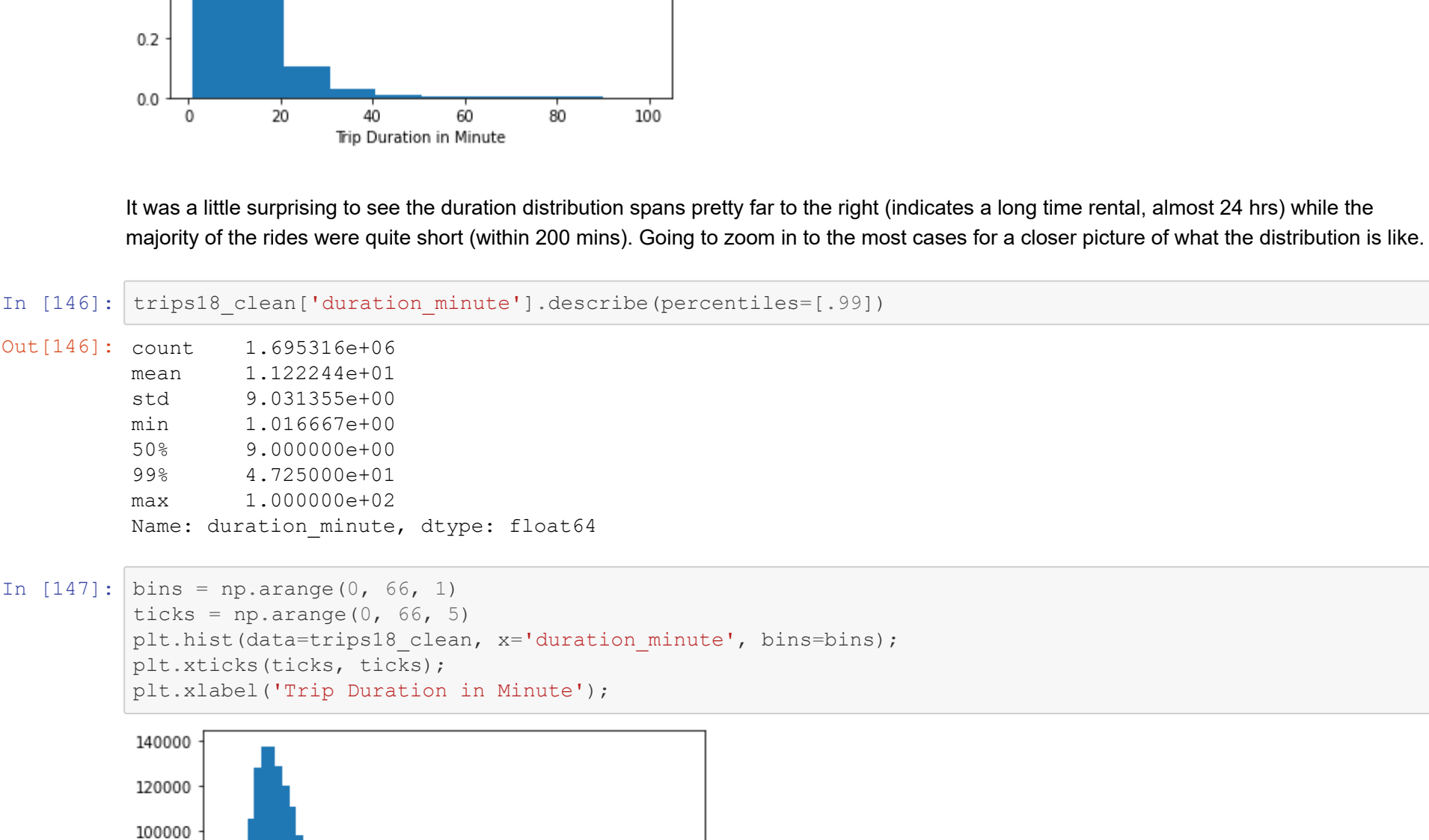


From graph above, most of the users are members which is about 7.6 times as many as casual users.



Most riders were male subscribers who did not use bike share for all trips. Most members were around 25 to 40 years old, corresponding to the trip distribution plots that most trips were on Mon-Fri and biked during rush hours of a day. As the age gets older, bike usage dropped significantly.

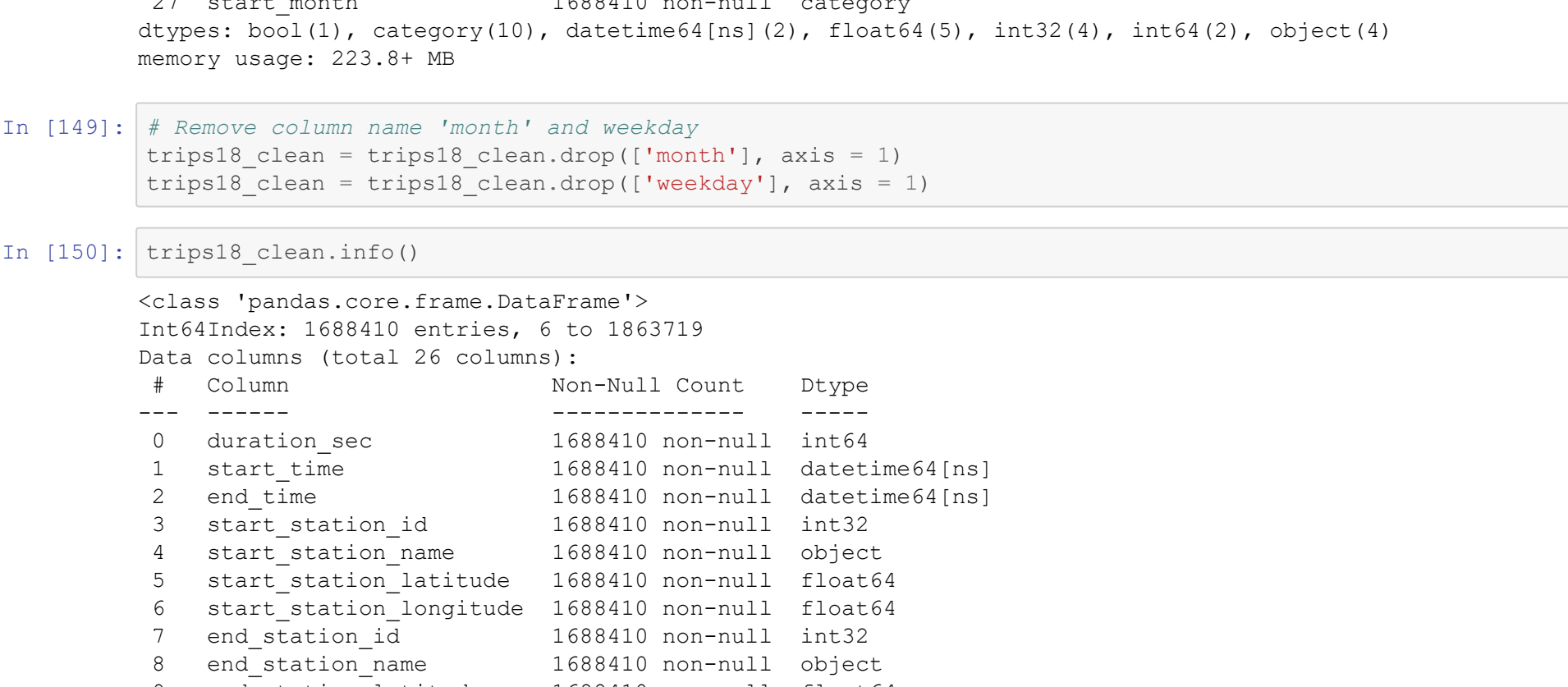
Trip duration distribution to plot next.



It was a little surprising to see the duration distribution spans pretty far to the right (indicates a long time rental, almost 24 hrs) while the majority of the rides were quite short (within 200 minutes). Going to zoom in to the most cases for a closer picture of what the distribution is like.

**trips18\_clean['duration\_minute'].describe(percentiles=[.99])**

```
Out[146]:
count    1.695316e+06
mean      1.122244e+01
std        9.031355e+00
min        1.016667e+00
50%        9.000000e+00
99%        4.725000e+01
max        1.000000e+02
Name: duration_minute, dtype: float64
```



It looks like 99% of the trips were less than 1 hour long with most of them between 5 to 10 minutes range. I'll remove the outlier long trips for the age group of all trip focusing on the most common cases to explore the pattern.

**Issue 8: filter out outlier trip records where the duration was very long**

```
trips18_clean = trips18_clean.query('duration_minute <= 66')
trips18_clean.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 1688410 entries, 6 to 1863719
Data columns (total 28 columns):
#   Column                                Non-Null Count  Dtype
---  --
0   duration_sec                          1688410 non-null  int64
1   start_time                            1688410 non-null  datetime64[ns]
2   end_time                              1688410 non-null  datetime64[ns]
3   start_station_id                      1688410 non-null  int32
4   start_station_name                    1688410 non-null  object
5   start_station_latitude                1688410 non-null  float64
6   start_station_longitude               1688410 non-null  float64
7   end_station_id                        1688410 non-null  int32
8   end_station_name                      1688410 non-null  object
9   end_station_latitude                  1688410 non-null  float64
10  end_station_longitude                 1688410 non-null  float64
11  bike_id                               1688410 non-null  int64
12  user_type                             1688410 non-null  category
13  member_birth_year                    1688410 non-null  int32
14  member_gender                         1688410 non-null  category
15  bike_share_for_all_trip               1688410 non-null  bool
16  start_petro_area                      1688410 non-null  category
17  end_metro_area                        1688410 non-null  category
18  month                                 1688410 non-null  category
19  weekday                               0 non-null       category
20  hour                                  1688410 non-null  int32
21  member_age                            1688410 non-null  int32
22  duration_minute                       1688410 non-null  float64
23  start_date                            1688410 non-null  object
24  start_hourofday                       1688410 non-null  object
25  start_dayofweek                       1688410 non-null  object
26  start_month                           1688410 non-null  category
27  start_year                            1688410 non-null  category
dtypes: bool(1), category(10), datetime64[ns](2), float64(5), int32(4), int64(2), object(4)
memory usage: 223.8+ MB
```

**Remove column name 'month' and 'weekday'**

```
trips18_clean = trips18_clean.drop('month', axis = 1)
trips18_clean = trips18_clean.drop('weekday', axis = 1)
```

**trips18\_clean.info()**

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 1688410 entries, 6 to 1863719
Data columns (total 26 columns):
#   Column                                Non-Null Count  Dtype
---  --
0   duration_sec                          1688410 non-null  int64
1   start_time                            1688410 non-null  datetime64[ns]
2   end_time                              1688410 non-null  datetime64[ns]
3   start_station_id                      1688410 non-null  int32
4   start_station_name                    1688410 non-null  object
5   start_station_latitude                1688410 non-null  float64
6   start_station_longitude               1688410 non-null  float64
7   end_station_id                        1688410 non-null  int32
8   end_station_name                      1688410 non-null  object
9   end_station_latitude                  1688410 non-null  float64
10  end_station_longitude                 1688410 non-null  float64
11  bike_id                               1688410 non-null  int64
12  user_type                             1688410 non-null  category
13  member_birth_year                    1688410 non-null  int32
14  member_gender                         1688410 non-null  category
15  bike_share_for_all_trip               1688410 non-null  bool
16  start_petro_area                      1688410 non-null  category
17  end_metro_area                        1688410 non-null  category
18  hour                                  1688410 non-null  int32
19  member_age                            1688410 non-null  int32
20  age_group                             1688410 non-null  category
21  duration_minute                       1688410 non-null  float64
22  start_date                            1688410 non-null  object
23  start_hourofday                       1688410 non-null  object
24  start_dayofweek                       1688410 non-null  category
25  start_month                           1688410 non-null  category
dtypes: bool(1), category(8), datetime64[ns](2), float64(5), int32(4), int64(2), object(4)
memory usage: 220.6+ MB
```

**Discuss the distribution(s) of your variable(s) of interest. Were there any unusual points? Did you need to perform any transformations?**

There were more trips on work days (Mon-Fri) compared to weekends. Summer time was the most popular season of a year, likely due to the weather.

Userside, there were more bike riders than female, and most members were subscribers compared to casual riders. The majority of the member did not use bike share for all of their trips, and most were around 25 to 40 years old.

Most rides were quick and short, lasted between 5 to 10 minutes, though there were some very long outliers like 24hrs.

No transformation was needed luckily due to the straightforwardness of the data.

**Of the features you investigated, were there any unusual distributions? Did you perform any operations on the data to tidy, adjust, or change the form of the data? If so, why did you do this?**

1. From duration graph, we can see the most rides takes within 5-10 minutes (300-600 secs). Durations dramatically higher are outliers probably due to user forget to log out after finish their ride.
2. From age graph, we can see people from 27 to 31 years old use the GoBike over 100,000 times. Customers could have given the wrong birth year information that cause the super old age data outliers.
3. From gender graph, male users are about 3 times as many as female users.
4. From user\_type graph, most of the users are members which is about 7.6 times as many as casual users.
5. From month graph, March and April have high usage over 200,000 and November and December's usage all below 130,000. This probably due to weather reasons. Hot summer season and cold winter season prevent users from riding. I also suspect holiday season could be reason for the dramatic drop in November and December, and school summer vacations could be reason for the drop in May.
6. From weekday graph, weekday usage is higher than weekend usage by almost 2 times. From morning and evening rush hours, there are two peaks of traffic during the day around 8-9am and 5-6pm. The times correspond to the morning and evening rush hours and so I believe people go to and leave from work on using the bike to commute between home, office or Bart/conect to other transportation.
7. From Most Popular Starting/Ending Stations graph, most popular starting stations are in SF area. And the fact that popular GoBike stations are near Bart or Caltrain stations confirms the idea people use bike to commute between public transportation and work/home.
8. From Area Distribution for Bike Activity graph, most bike riding activities happened in San Francisco area. This is make sense because SF is the first location of GoBike. I am curious about if there is any ride that starting from one area but end in another area.
9. Inter-area Bike Activity graph, among the 114592 bike activity in df\_clean, only 65 inter-area rides. All of the rides are between SF and Oakland across the Bay Bridge.

**save the clean data to a .csv file**

```
trips18_clean.to_csv('fordgoBike_trips_2018_clean.csv', index=False)
```

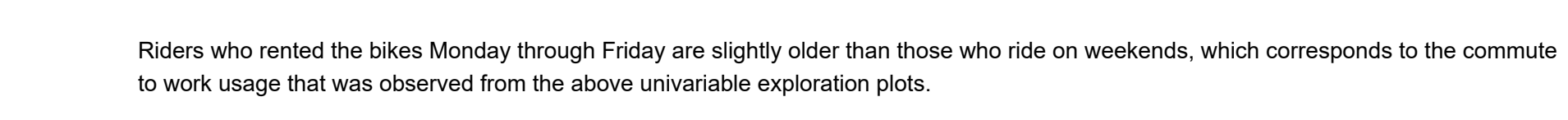
**Bivariate Exploration**

**Average Trip Duration on Weekdays**



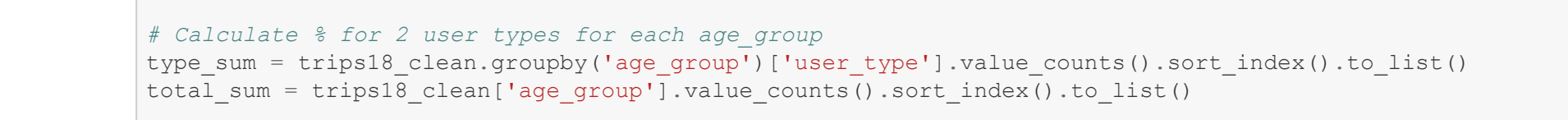
The riding trips are much shorter on Monday through Friday compared to weekends. It indicates a pretty stable and efficient usage of the sharing system on normal work days, while more casual flexible use on weekends.

**Average Trip Duration by month**



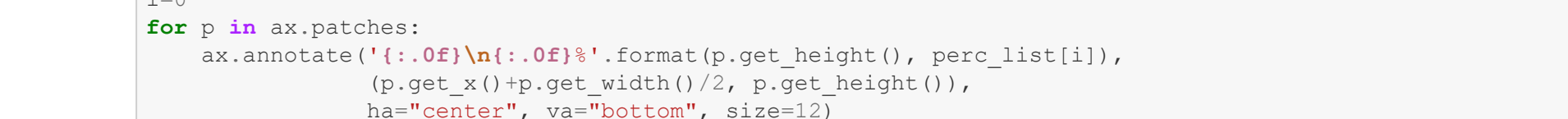
Starting from March, the average usage time shown an increasing trend almost all the way till October before it down turned. The average longest trips happened during the summer months, around June, July and September, which probably has a lot to do with the weather in the area.

**Member Ages by weekdays**



Riders who rented the bikes Monday through Friday are slightly older than those who ride on weekends, which corresponds to the commute to work usage that was observed from the above univariable exploration plots.

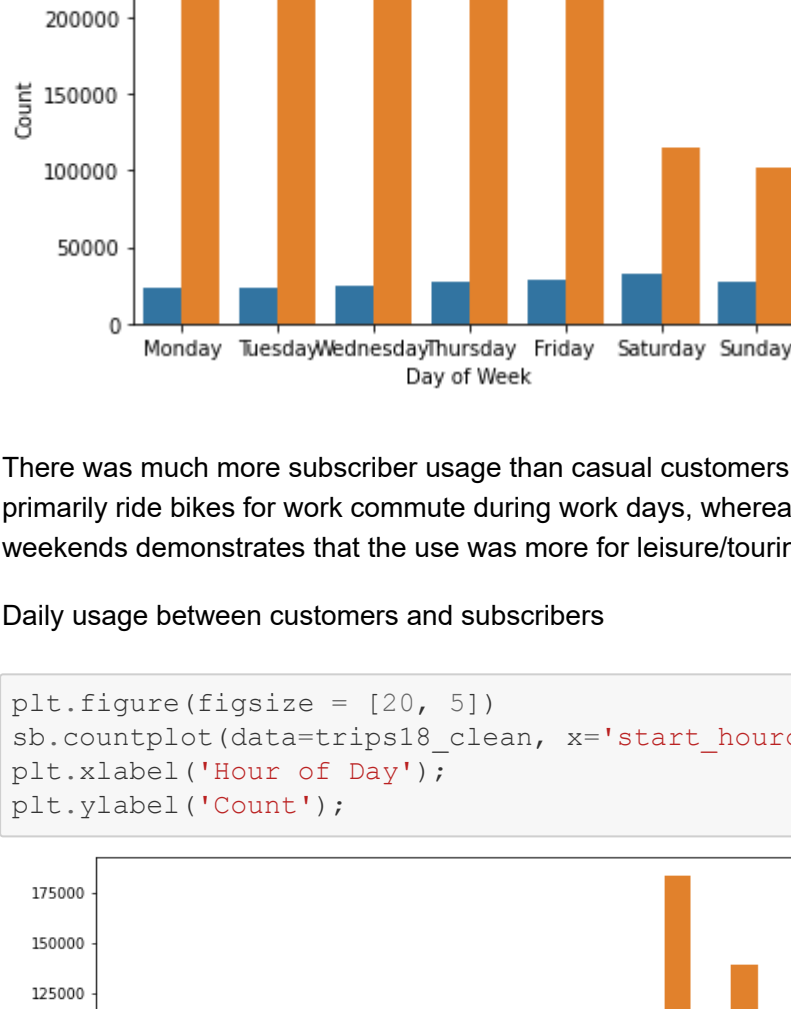
**Weekly usage between customers and subscribers**



From graph above, user's subscription ratio increases as the age increase. Younger user tend to use the service but do not subscribe.



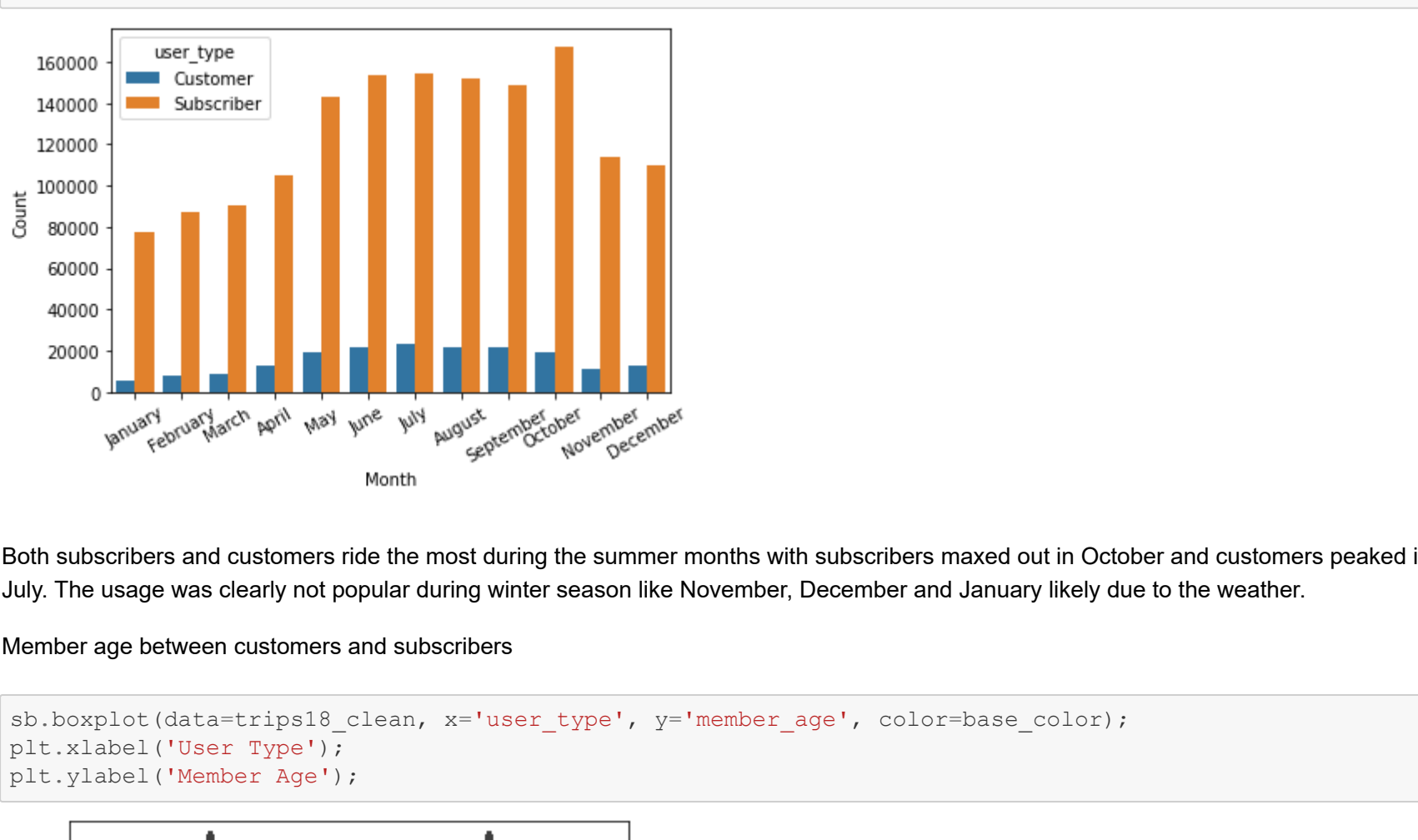
```
In [156]: sb.countplot(data=trips18_clean, x="start_dayofweek", hue="user_type");
plt.title('User type wise Start Week');
plt.xlabel('Day of Week');
plt.ylabel('Count');
```



There was much more subscriber usage than casual customers overall. The drop of volume on weekends for subscribers indicates that they primarily ride bikes for work commute during work days, whereas almost the opposite pattern of a slight increase of use for customers on weekends demonstrates that the use was more for leisure/touring and relaxing purposes.

Daily usage between customers and subscribers

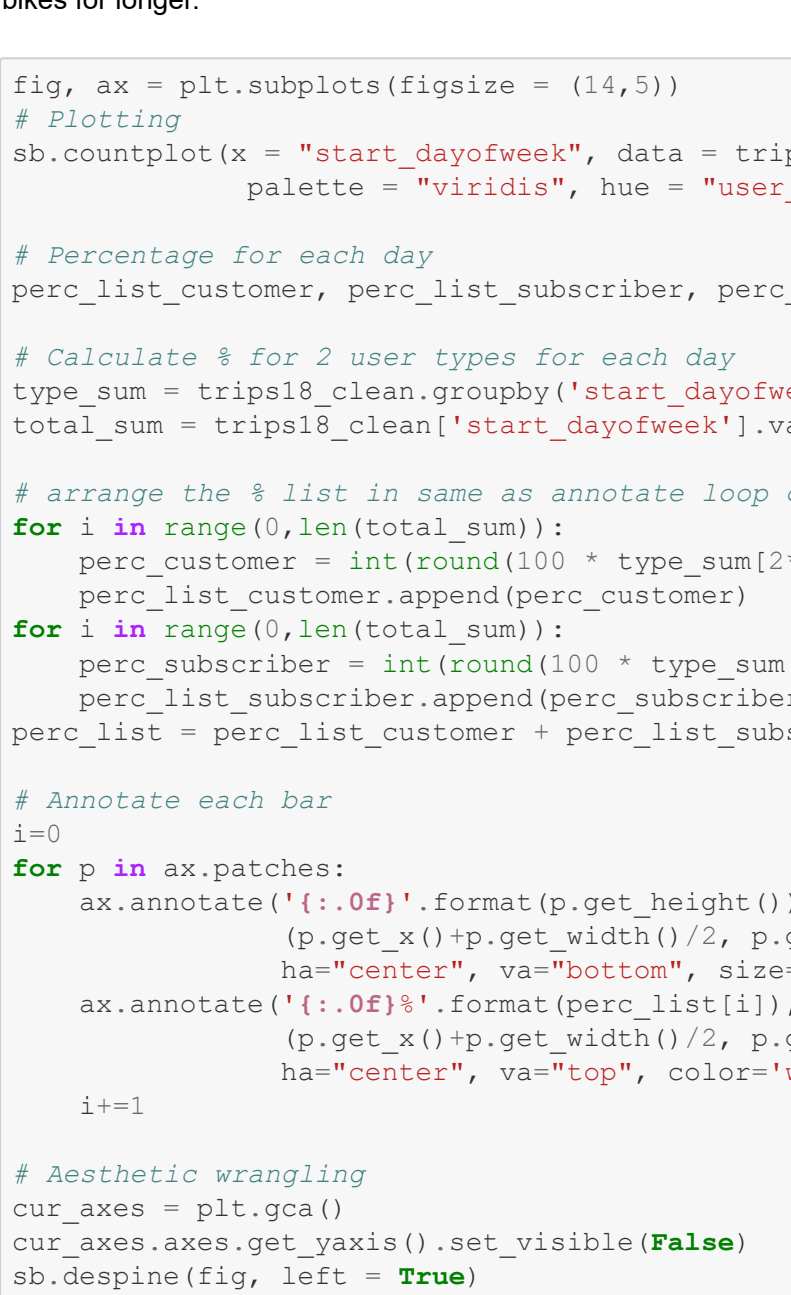
```
In [157]: plt.figure(figsize = (20, 5))
sb.countplot(data=trips18_clean, x="start_hourofday", hue="user_type");
plt.xlabel('Hour of Day');
plt.ylabel('Count');
```



Subscriber usage clearly peaks out on typical rush hours when people go to work in the morning and getting off work in the afternoon, double confirmed their usage purpose and goal of riding. Similar pattern was not observed among customers who tend to ride most in the afternoon or early evening as for a different purpose than the subscribers riders.

Yearly usage between customers and subscribers

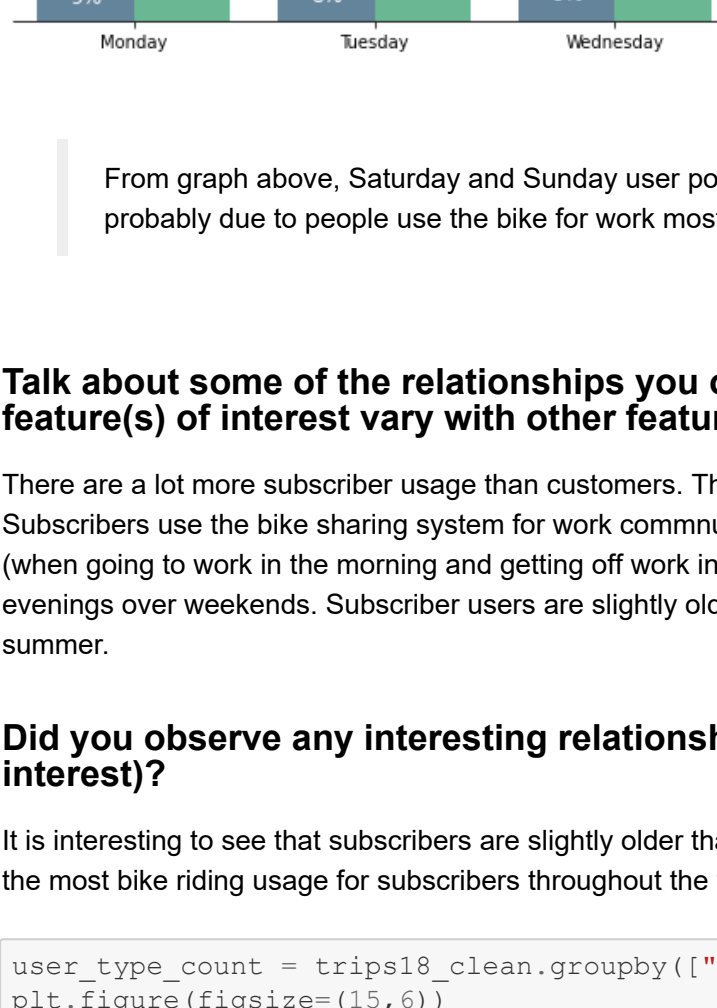
```
In [158]: sb.countplot(data=trips18_clean, x="start_month", hue="user_type");
plt.xticks(rotation=30);
plt.xlabel('Month');
plt.ylabel('Count');
```



Both subscribers and customers ride the most during the summer months with subscribers maxed out in October and customers peaked in July. The usage was clearly not popular during winter season like November, December and January likely due to the weather.

Member age between customers and subscribers

```
In [159]: sb.boxplot(data=trips18_clean, x="user_type", y="member_age", color=base_color);
plt.xlabel('Member Age');
```



Similar to the Member age by weekdays plot, subscribers who ride most often Monday through Friday are slightly older than customers, with a wider range of ages as well.

The trip duration distribution is much narrower for subscribers compared to casual riders on the shorter/quicker trip end overall. It seems like subscribers have a more specific usage or targeted goal riding the bikes compared to customers who vary more and generally rented the bikes for longer.

```
In [160]: fig, ax = plt.subplots(figsize = (14,5))
# Plotting
sb.countplot(x = "start_dayofweek", data = trips18_clean,
             palette = "viridis", hue = "user_type", alpha = 0.8)

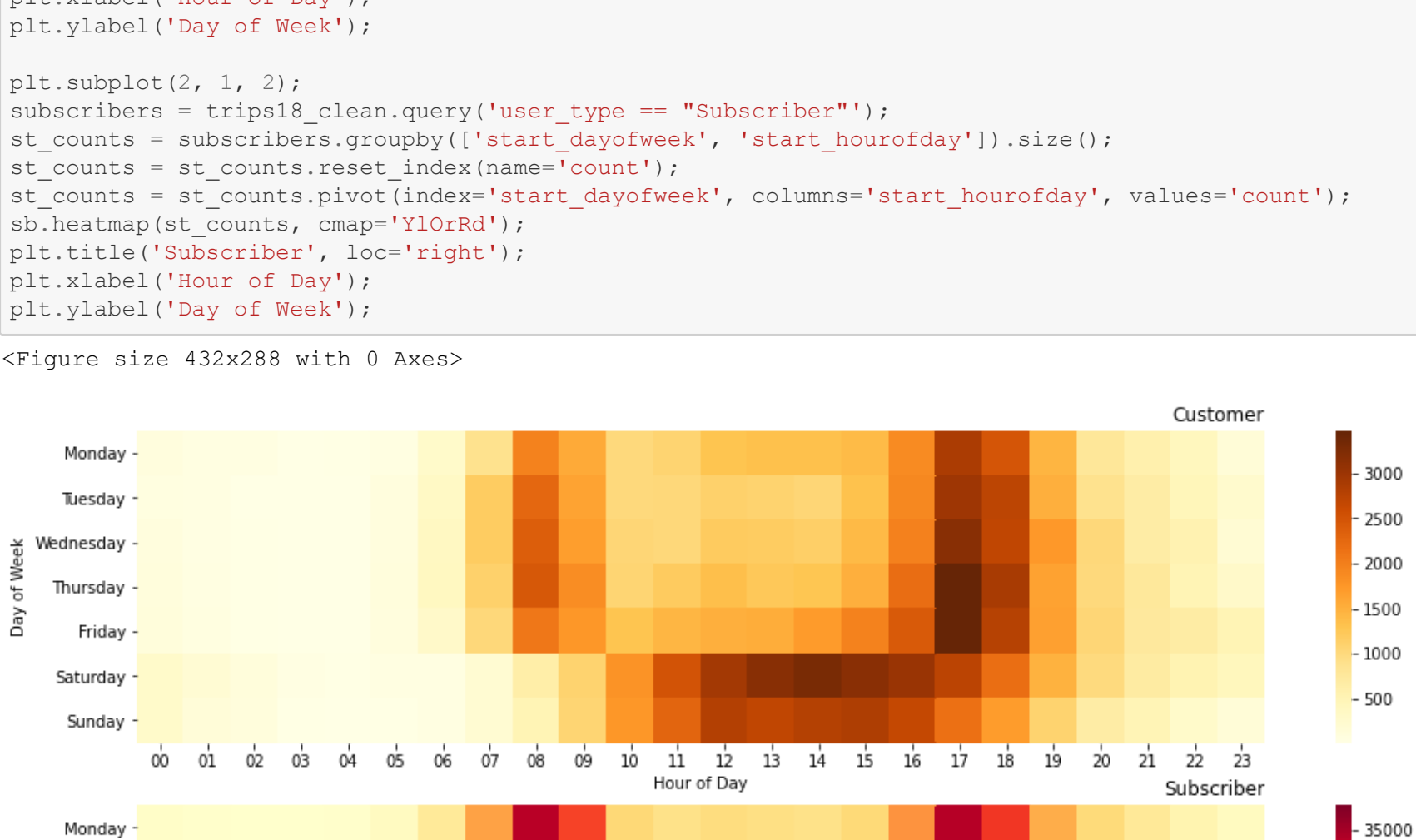
# Percentage for each day
perc_list_customer, perc_list_subscriber, perc_list = [], [], []

# Calculate % for 2 user types for each day
type_sum = trips18_clean.groupby(['start_dayofweek'])['user_type'].value_counts().sort_index().to_list()
total_sum = trips18_clean['start_dayofweek'].value_counts().sort_index().to_list()

# Arrange the % list in same as annotate loop order
for i in range(0,len(total_sum)):
    perc_customer = int(round(100 * type_sum[2*i] / total_sum[i]))
    perc_list_customer.append(perc_customer)
for i in range(0,len(total_sum)):
    perc_subscriber = int(round(100 * type_sum[2*i+1] / total_sum[i]))
    perc_list_subscriber.append(perc_subscriber)
perc_list = perc_list_customer + perc_list_subscriber

# Annotate each bar
i=0
for p in ax.patches:
    ax.annotate('{:.0f}%'.format(p.get_height()),
                (p.get_x()+p.get_width()/2, p.get_height()),
                ha="center", va="bottom", size=12)
    ax.annotate('{:.0f}%'.format(perc_list[i]),
                (p.get_x()+p.get_width()/2, p.get_height()),
                ha="center", va="top", color="white", size=12)
    i+=1

# Aesthetic wrangling
cur_ax = plt.gca()
cur_ax.set_ylabel('Count')
cur_ax.set_xlabel('Day of Week')
sb.despine(fig, left = True)
plt.title('Users Type by Weekday', fontweight = 'bold', fontstyle = 'italic', fontsize = 20)
plt.ylabel('Count')
plt.legend()
```



From graph above, Saturday and Sunday user populations have unusual low subscription ratio and total rides. This is probably due to people use the bike for work mostly.

**Talk about some of the relationships you observed in this part of the investigation. How did the feature(s) of interest vary with other features in the dataset?**

There are a lot more subscriber usage than customers. The riding habit/pattern varies a lot between subscribers and customers. Subscribers use the bike sharing system for work commute thus most trips were on work days (Mon-Fri) and especially during rush hours (when going to work in the morning and getting off work in the afternoon), whereas customers tend to ride for fun in the afternoon or early evenings over weekends. Subscriber users are slightly older than customer users who tend to take longer rides overall especially during the summer.

**Did you observe any interesting relationships between the other features (not the main feature(s) of interest)?**

It is interesting to see that subscribers are slightly older than customers on average but take much shorter/quicker rides. And October had the most bike riding usage for subscribers throughout the year.

```
In [161]: user_type_count = trips18_clean.groupby(['start_month', "user_type"]).size().reset_index()
# Plotting
plt.figure(figsize=(15, 6))
color = ('Subscriber':'orange', 'Customer':'black')
axis = sb.pointplot(x="start_month", y=0, hue="user_type", palette=color, scale=7, data=user_type_count)

plt.title('The monthly bike rides per user type')
plt.xlabel('Date')
plt.ylabel('Count')
plt.legend()
```



## Multivariate Exploration

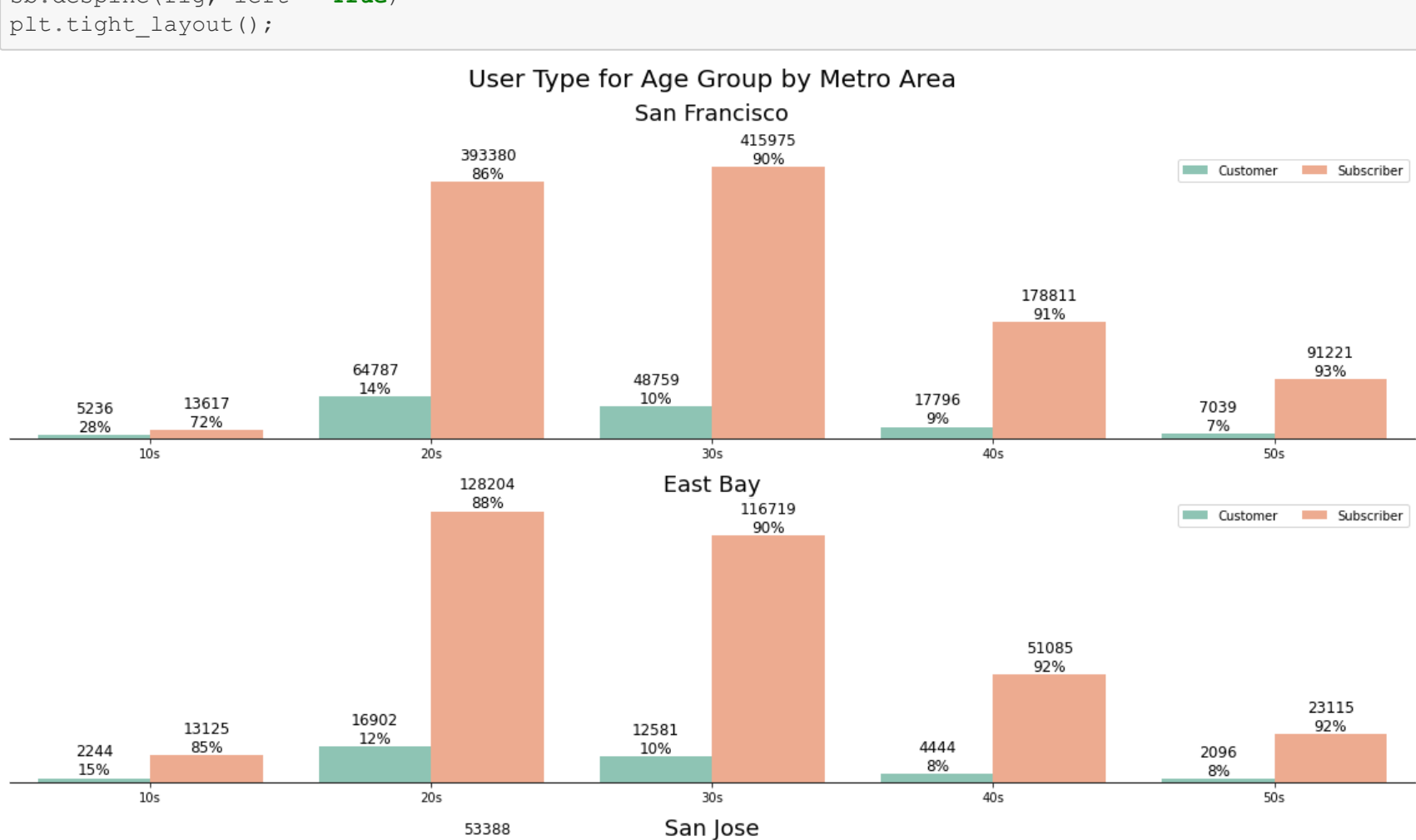
How does the hourly usage vary during weekdays for customers and subscribers?

```
In [162]: plt.subplot(2, 1, 1);
plt.figure(figsize=(5,5))

customers = trips18_clean.query('user_type == "Customer"');
ct_counts = customers.groupby(['start_dayofweek', 'start_hourofday']).size();
ct_counts = ct_counts.reset_index(name='count');
sb.heatmap(ct_counts, cmap='inferno');
plt.title('Customer', loc='right');
plt.xlabel('Hour of Day');

subscribers = trips18_clean.query('user_type == "Subscriber"');
st_counts = subscribers.groupby(['start_dayofweek', 'start_hourofday']).size();
st_counts = st_counts.reset_index(name='count');
sb.heatmap(st_counts, cmap='inferno');
plt.title('Subscriber', loc='right');
plt.xlabel('Hour of Day');

<Figure size 432x288 with 0 Axes>
```



By the above heat map, we can see clearly the trips taken by customers and subscribers on weekdays and weekends.

user\_type for age\_group by metro\_area

```
In [163]: # Since there's only three subplots to create, using the full data should be fine.
fig = plt.figure(figsize = (16, 12))

# Subplot 1: San Francisco, ride count for age_group by user_type
#####
ax1 = plt.subplot(3, 1, 1)
ax1.set_title("San Francisco", size=18)
sb.countplot(data = trips18_clean.query('start_metro_area == "San Francisco"'),
             x = 'age_group', hue = 'user_type', palette = 'Set2', alpha=0.8)
plt.xlabel('')
plt.ylabel('')
cur_ax = plt.gca()
cur_ax.set_ylabel('Count')
cur_ax.set_xlabel('Age Group')
cur_ax.legend(loc = 1, ncol = 2)

# Percentage for each age_group
perc_list_customer, perc_list_subscriber, perc_list = [], [], []

# Calculate % for 2 user types for each age_group
type_sum = trips18_clean.groupby(['start_metro_area == "San Francisco"'])['user_type'].value_counts().sort_index().to_list()
total_sum = trips18_clean.query('start_metro_area == "San Francisco"')['age_group'].value_counts().sort_index().to_list()

# Arrange the % list in same as annotate loop order
for i in range(0,len(total_sum)):
    perc_customer = int(round(100 * type_sum[2*i] / total_sum[i]))
    perc_list_customer.append(perc_customer)
for i in range(0,len(total_sum)):
    perc_subscriber = int(round(100 * type_sum[2*i+1] / total_sum[i]))
    perc_list_subscriber.append(perc_subscriber)
perc_list = perc_list_customer + perc_list_subscriber

# Add annotations
i=0
for p in ax1.patches:
    ax1.annotate('{:.0f}%\n{:.0f}%'.format(p.get_height(), perc_list[i]),
                (p.get_x()+p.get_width()/2, p.get_height()),
                ha="center", va="bottom", size=12)
    i+=1

# Subplot 2: East Bay, ride count for age_group by user_type
#####
ax2 = plt.subplot(3, 1, 2)
ax2.set_title("East Bay", size=18)
sb.countplot(data = trips18_clean.query('start_metro_area == "East Bay"'),
             x = 'age_group', hue = 'user_type', palette = 'Set2', alpha=0.8)
plt.xlabel('')
plt.ylabel('')
cur_ax = plt.gca()
cur_ax.set_ylabel('Count')
cur_ax.set_xlabel('Age Group')
cur_ax.legend(loc = 1, ncol = 2)

# Percentage for each age_group
perc_list_customer, perc_list_subscriber, perc_list = [], [], []

# Calculate % for 2 user types for each age_group
type_sum = trips18_clean.query('start_metro_area == "East Bay"')['user_type'].value_counts().sort_index().to_list()
total_sum = trips18_clean.query('start_metro_area == "East Bay"')['age_group'].value_counts().sort_index().to_list()

# Arrange the % list in same as annotate loop order
for i in range(0,len(total_sum)):
    perc_customer = int(round(100 * type_sum[2*i] / total_sum[i]))
    perc_list_customer.append(perc_customer)
for i in range(0,len(total_sum)):
    perc_subscriber = int(round(100 * type_sum[2*i+1] / total_sum[i]))
    perc_list_subscriber.append(perc_subscriber)
perc_list = perc_list_customer + perc_list_subscriber

# Add annotations
i=0
for p in ax2.patches:
    ax2.annotate('{:.0f}%\n{:.0f}%'.format(p.get_height(), perc_list[i]),
                (p.get_x()+p.get_width()/2, p.get_height()),
                ha="center", va="bottom", size=12)
    i+=1

# Subplot 3: East Bay, ride count for age_group by user_type
#####
ax3 = plt.subplot(3, 1, 3)
ax3.set_title("San Jose", size=18)
sb.countplot(data = trips18_clean.query('start_metro_area == "San Jose"'),
             x = 'age_group', hue = 'user_type', palette = 'Set2', alpha=0.8)
plt.xlabel('')
plt.ylabel('')
cur_ax = plt.gca()
cur_ax.set_ylabel('Count')
cur_ax.set_xlabel('Age Group')
cur_ax.legend(loc = 1, ncol = 2)

# Percentage for each age_group
perc_list_customer, perc_list_subscriber, perc_list = [], [], []

# Calculate % for 2 user types for each age_group
type_sum = trips18_clean.query('start_metro_area == "San Jose"')['user_type'].value_counts().sort_index().to_list()
total_sum = trips18_clean.query('start_metro_area == "San Jose"')['age_group'].value_counts().sort_index().to_list()

# Arrange the % list in same as annotate loop order
for i in range(0,len(total_sum)):
    perc_customer = int(round(100 * type_sum[2*i] / total_sum[i]))
    perc_list_customer.append(perc_customer)
for i in range(0,len(total_sum)):
    perc_subscriber = int(round(100 * type_sum[2*i+1] / total_sum[i]))
    perc_list_subscriber.append(perc_subscriber)
perc_list = perc_list_customer + perc_list_subscriber

# Add annotations
i=0
for p in ax3.patches:
    ax3.annotate('{:.0f}%\n{:.0f}%'.format(p.get_height(), perc_list[i]),
                (p.get_x()+p.get_width()/2, p.get_height()),
                ha="center", va="bottom", size=12)
    i+=1

# Aesthetic Wrangling
fig.suptitle('User Type for Age Group by Metro Area', size=20, y=1.02)
sb.despine(fig, left = True)
plt.tight_layout()
```



From graph above, majority of users is from age groups 20s and 30s. Although San Francisco metro area has most rides but the two largest user populations 20s and 30s age groups have the lowest subscription ratio.

## Conclusion

The relationship between the multiple variables plotted are visualized altogether and information are presented. The short period of usage for subscribers corresponds to their high concentration on rush hours Monday through Friday, which is for work. The customer use shows that they're taking advantage of the bike sharing system quite differently from the subscribers, heavily over weekends and in the afternoon, for city tour or leisure purpose mostly.

The interactions between features are all supplementing each other and quite make sense when looked at combined, there's no big surprise observed.

- We observed both users prefer to use a bike to go to work and to leave work and also in the afternoon they use bikes to go to restaurants or cafes.
- It has observed that on weekdays the usage level is higher than at weekends for subscribers but customer usage shows that approximately the same level at week.
- It observed that the average trip time is approximately 15 minutes.
- It has been observed customer usage trend increases towards the end of the year. On the other hand, subscriber usage trend gets the least level in the winter year.
- It can be seen that subscribers use bikes for a short time period than customers.
- It does not seem that any regular trends in a year. Subscriber usage seems that they are the least ride at the end of the year this situation maybe because of the Christmas holiday, they have a holiday and do not have to go to work. On the other hand, customer usage shows that the highest level in December for them.
- It can be seen that subscribers ride much shorter/quicker trips compared to customers on each day of the week. Both user types have an obvious increase in trip duration on Saturdays and Sundays over weekends. Subscribers usage seems to maintain a very consistent average duration Monday through Friday.

```
In [ ] :
```