

Communicate Data Findings (Ford GoBike System Data)

by (Amal Aljabri)

Preliminary Wrangling

This document explores a dataset that includes information about individual rides made in a bike-sharing system covering the greater San Francisco Bay Area.

```
In [1]: # Import all packages and set plots to be embedded inline
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from matplotlib import cm
import seaborn as sns
import matplotlib inline
```

```
In [2]: from IPython.core.display import display
sns.set(rc={'figure.figsize':(14,6)})
```

```
In [3]: # Load in the dataset into a pandas dataframe
df = pd.read_csv('201802-fordgo-bike-tripdata.csv')
```

Dataset Overview

```
In [4]: # High-level overview of data shape and composition
print(df.shape)
print(df.dtypes)
print(df.info())
df.head(10)
```

```
(183412, 16)
duration_sec      int64
start_time        object
end_time          object
start_station_id  float64
start_station_name object
start_station_latitude float64
start_station_longitude float64
end_station_id    float64
end_station_name  object
end_station_latitude float64
end_station_longitude float64
bike_id           int64
user_type         object
member_birth_year float64
member_gender     object
bike_share_for_all_trip object
dtype: object
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 183412 entries, 0 to 183411
Data columns (total 16 columns):
 #   Column              Non-Null Count  Dtype
---  --
 0   duration_sec        183412 non-null    int64
 1   start_time          183412 non-null    object
 2   end_time            183412 non-null    object
 3   start_station_id    183215 non-null    float64
 4   start_station_name  183215 non-null    object
 5   start_station_latitude 183412 non-null    float64
 6   start_station_longitude 183412 non-null    float64
 7   end_station_id      183215 non-null    float64
 8   end_station_name    183215 non-null    object
 9   end_station_latitude 183412 non-null    float64
10   end_station_longitude 183412 non-null    float64
11   bike_id             183412 non-null    int64
12   user_type           183412 non-null    object
13   member_birth_year   175147 non-null    float64
14   member_gender       175147 non-null    object
15   bike_share_for_all_trip 183412 non-null    object
dtypes: float64(7), int64(2), object(7)
memory usage: 22.4+ MB
None
```

```
Out [4]:
```

	duration_sec	start_time	end_time	start_station_id	start_station_name	start_station_latitude	start_station_longitude	end_station
0	52185	2019-02-28 17:32:10.145	2019-03-01 08:01:55.975	21.0	Montgomery St BART Station (Market St at 2nd St)	37.789625	-122.400811	1
1	42521	2019-02-28 18:53:37.780	2019-03-01 06:24:08.140	23.0	The Embarcadero at Seast St	37.791464	-122.391034	8
2	61854	2019-02-28 17:54:26.010	2019-03-01 04:02:36.842	86.0	Market St at Dolores St	37.789305	-122.426826	
3	36490	2019-02-28 17:54:26.010	2019-03-01 04:02:36.842	375.0	Grove St at Masonic Ave	37.774836	-122.446546	7
4	1585	2019-02-28 23:41:06.780	2019-03-01 00:20:44.070	7.0	Frank H Ogawa Plaza	37.804562	-122.271738	22
5	1793	2019-02-28 23:41:06.780	2019-03-01 00:07:58.715	93.0	4th St at Mission Bay Blvd S	37.770407	-122.391198	32
6	1147	2019-02-28 23:55:35.104	2019-03-01 00:14:42.588	300.0	Palm St at Willow St	37.317298	-121.884995	31
7	1615	2019-02-28 23:41:06.780	2019-03-01 00:08:02.750	10.0	Washington St at Kearny St	37.795393	-122.404770	12
8	1570	2019-02-28 23:41:06.780	2019-03-01 00:07:58.715	10.0	Washington St at Kearny St	37.795393	-122.404770	12
9	1049	2019-02-28 23:49:47.699	2019-03-01 00:07:17.025	19.0	Post St at Kearny St	37.788975	-122.403452	12

```
In [5]: # Descriptive statistics for numeric variables
df.describe()
```

```
Out [5]:
```

	duration_sec	start_station_id	start_station_latitude	start_station_longitude	end_station_id	end_station_latitude	end_station_longitude
count	183412.000000	183215.000000	183412.000000	183412.000000	183215.000000	183412.000000	183412.000000
mean	726.078435	138.590427	37.771223	-122.352664	136.249123	37.771427	-122.352664
std	1794.389780	111.778864	0.099581	0.117097	111.515131	0.096490	0.116146
min	61.000000	4.000000	37.317298	-122.453704	3.000000	37.317298	-122.453704
25%	325.000000	37.000000	-122.412408	-122.412408	44.000000	37.770407	-122.412408
50%	516.000000	104.000000	37.780760	-122.396285	100.000000	37.781010	-122.396285
75%	796.000000	239.000000	37.797280	-122.286533	235.000000	37.797320	-122.286533
max	85444.000000	398.000000	37.880222	-121.874119	398.000000	37.880222	-121.874119

```
In [6]: # Check if duplicates exist
df.duplicated().sum()
```

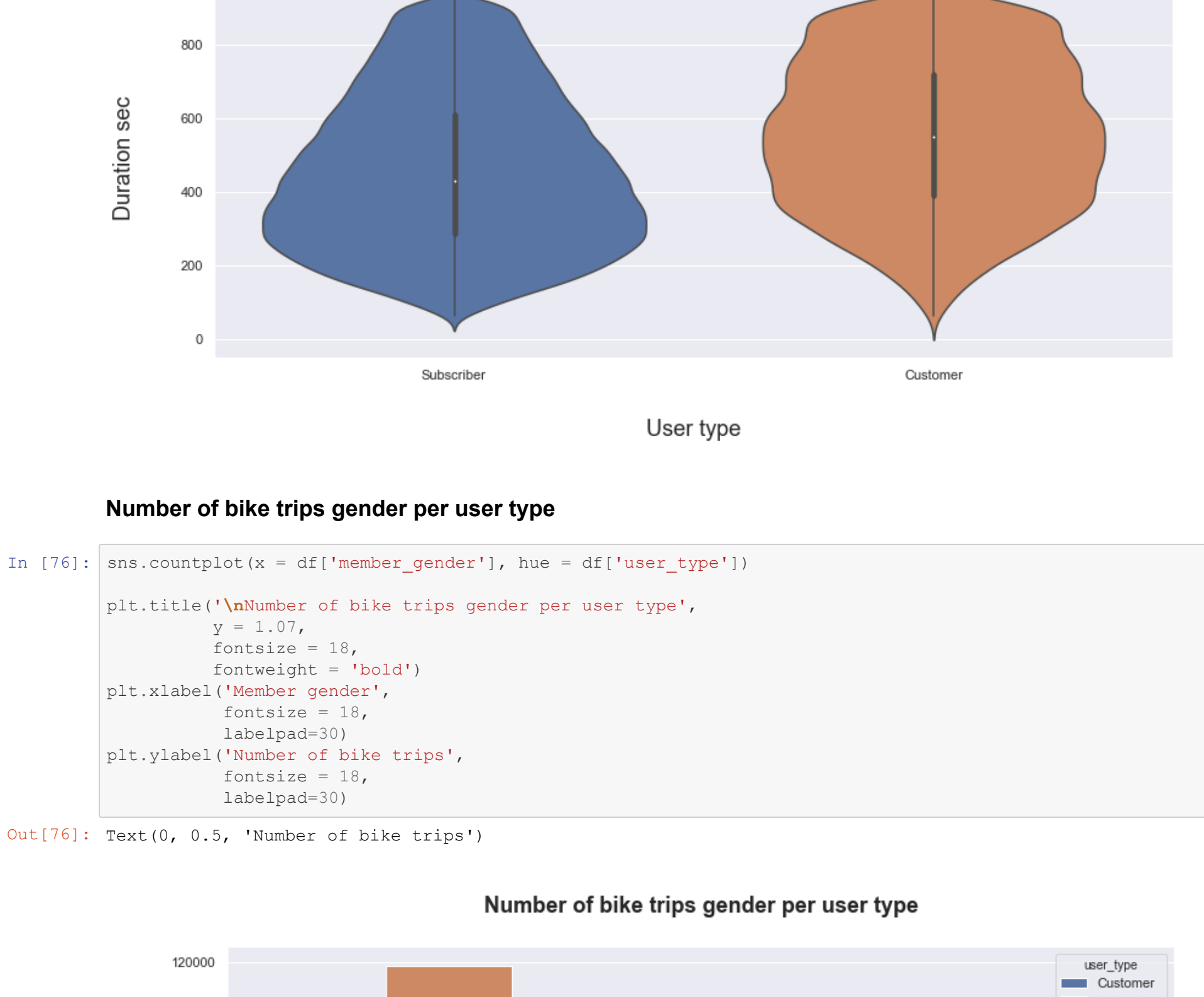
```
Out [6]: 0
```

```
In [7]: # Check if unique values
df.nunique()
```

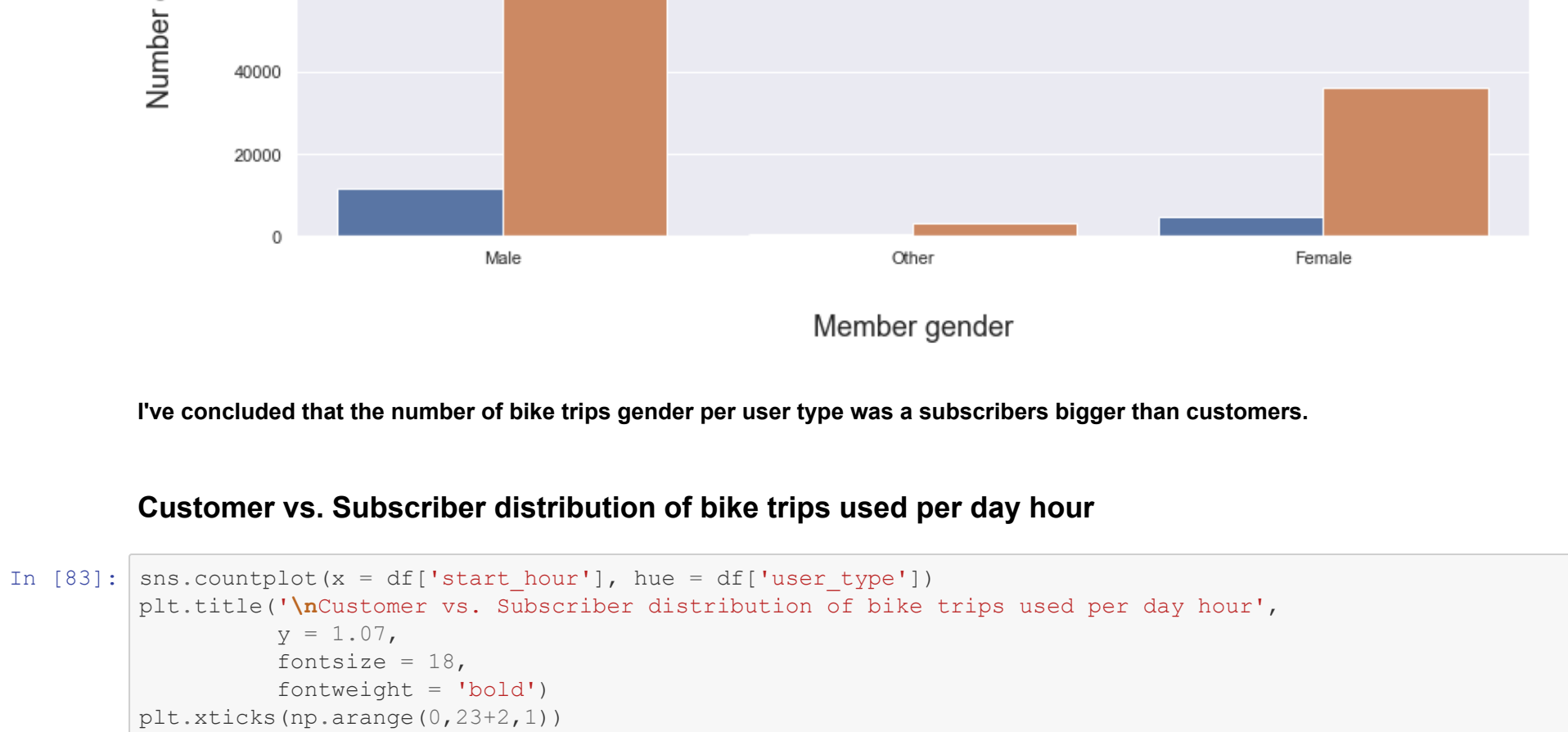
```
Out [7]:
```

	duration_sec	start_time	end_time	start_station_id	start_station_name	start_station_latitude	start_station_longitude	end_station_id	end_station_name	end_station_latitude	end_station_longitude
0	4752	183401	183397	329	329	334	335	329	329	335	335
1	183401	183397	329	334	335	329	335	335	335	335	335
2	183397	329	334	335	335	329	335	335	335	335	335
3	329	334	335	335	335	329	335	335	335	335	335
4	334	335	329	335	335	329	335	335	335	335	335
5	335	335	329	335	335	329	335	335	335	335	335
6	335	335	329	335	335	329	335	335	335	335	335
7	335	335	329	335	335	329	335	335	335	335	335
8	335	335	329	335	335	329	335	335	335	335	335
9	335	335	329	335	335	329	335	335	335	335	335
10	335	335	329	335	335	329	335	335	335	335	335
11	335	335	329	335	335	329	335	335	335	335	335
12	335	335	329	335	335	329	335	335	335	335	335
13	335	335	329	335	335	329	335	335	335	335	335
14	335	335	329	335	335	329	335	335	335	335	335
15	335	335	329	335	335	329	335	335	335	335	335
16	335	335	329	335	335	329	335	335	335	335	335
17	335	335	329	335	335	329	335	335	335	335	335
18	335	335	329	335	335	329	335	335	335	335	335
19	335	335	329	335	335	329	335	335	335	335	335
20	335	335	329	335	335	329	335	335	335	335	335
21	335	335	329	335	335	329	335	335	335	335	335
22	335	335	329	335	335	329	335	335	335	335	335
23	335	335	329	335	335	329	335	335	335	335	335
24	335	335	329	335	335	329	335	335	335	335	335
25	335	335	329	335	335	329	335	335	335	335	335
26	335	335	329	335	335	329	335	335	335	335	335
27	335	335	329	335	335	329	335	335	335	335	335
28	335	335	329	335	335	329	335	335	335	335	335
29	335	335	329	335	335	329	335	335	335	335	335
30	335	335	329	335	335	329	335	335	335	335	335
31	335	335	329	335	335	329	335	335	335	335	335
32	335	335	329	335	335	329	335	335	335	335	335
33	335	335	329	335	335	329	335	335	335	335	335
34	335	335	329	335	335	329	335	335	335	335	335
35	335	335	329	335	335	329	335	335	335	335	335
36	335	335	329	335	335	329	335	335	335	335	335
37	335	335	329	335	335	329	335	335	335	335	335
38	335	335	329	335	335	329	335	335	335	335	335
39	335	335	329	335	335	329	335	335	335	335	335
40	335	335	329	335	335	329	335	335	335	335	335
41	335	335	329	335	335	329	335	335	335	335	335
42	335	335	329	335	335	329	335	335	335	335	335
43	335	335	329	335	335	329	335	335	335	335	335
44	335	335	329	335	335	329	335	335	335	335	335
45	335	335	329	335	335	329	335	335	335	335	335
46	335	335	329	335	335	329	335	335	335	335	335
47	335	335	329	335	335	329	335	335	335	335	335
48	335	335	329	335	335	329	335	335	335	335	335
49	335	335	329	335	335	329	335	335	335	335	335
50	335	335	329	335	335	329	335	335	335	335	335
51	335	335	329	335	335	329	335	335	335	335	335
52	335	335	329	335	335	329	335	335	335	335	335
53	335	335	329	335	335	329	335	335	335	335	335
54	335	335	329	335	335	329	335	335	335	335	335
55	335	335	329	335	335	329	335	335	335	335	335
56	335	335	329	335	335	329	335	335	335	335	335
57	335	335	329	335	335	329	335	335	335	335	335
58	335	335	329	335	335	329	335	335	335	335	335
59	335	335	329	335	335	329	335	335	335	335	335
60	335	335	329	335	335	329	335	335	335	335	335
61	335	335	329	335	335	329	335	335	335	335	335
62	335	335	329	335	335	329	335	335	335	335	335
63	335	335	329	335	335	329	335	335	335	335	335
64	335	335	329	335	335	329	335	335	335	335	335
65	335	335	329	335	335	329	335	335	335	335	335
66	335	335	329	335	335	329	335	335	335	335	335
67	335	335	329	335	335	329	335	335	335	335	335
68	335	335	329	335	335	329	335	335	335	335	335
69	335	335	329	335	335	329	335	335	335	335	335
70	335	335	329	335	335	329	335	335	335	335	335
71	335	335	329	335	335	329	335	335	335	335	335
72	335	335	329	335	335	329	335	335	335	335	335
73	335	335	329	335	335	329	335	335	335	335	335
74	335	335	329	335	335	329	335	335	335	335	335
75	335	335	329	335	335	329	335	335	335	335	335
76	335	335	329	335	335	329	335	335	335	335	335
77	335	335	329	335	335	329	335	335	335	335	335
78	335	335	329	335	335	329	335	335	335	335	335
79	335	335	329	335	335	329	335	335	335	335	335
80	335	335	329	335	335	329	335	335	335	335	335
81	335	335	329	335	335	329	335	335	335	335	335
82	335	335	329	335	335	329	335	335	335	335	335
83	335	335	329	335	335	329	335	335	335	335	335
84	335	335	329	335	335	329	335	335	335	335	335
85	335	335	329	335	335	329	335	335	335	335	335
86	335	335	329	335	335						


```
In [117]: sns.violinplot(data= df[df['duration_sec'] <= 923],x='user_type',y='duration_sec')
plt.xlabel('User type',
           fontsize = 18,
           labelpad=30)
plt.ylabel('Duration sec',
           fontsize = 18,
           labelpad=30)
```

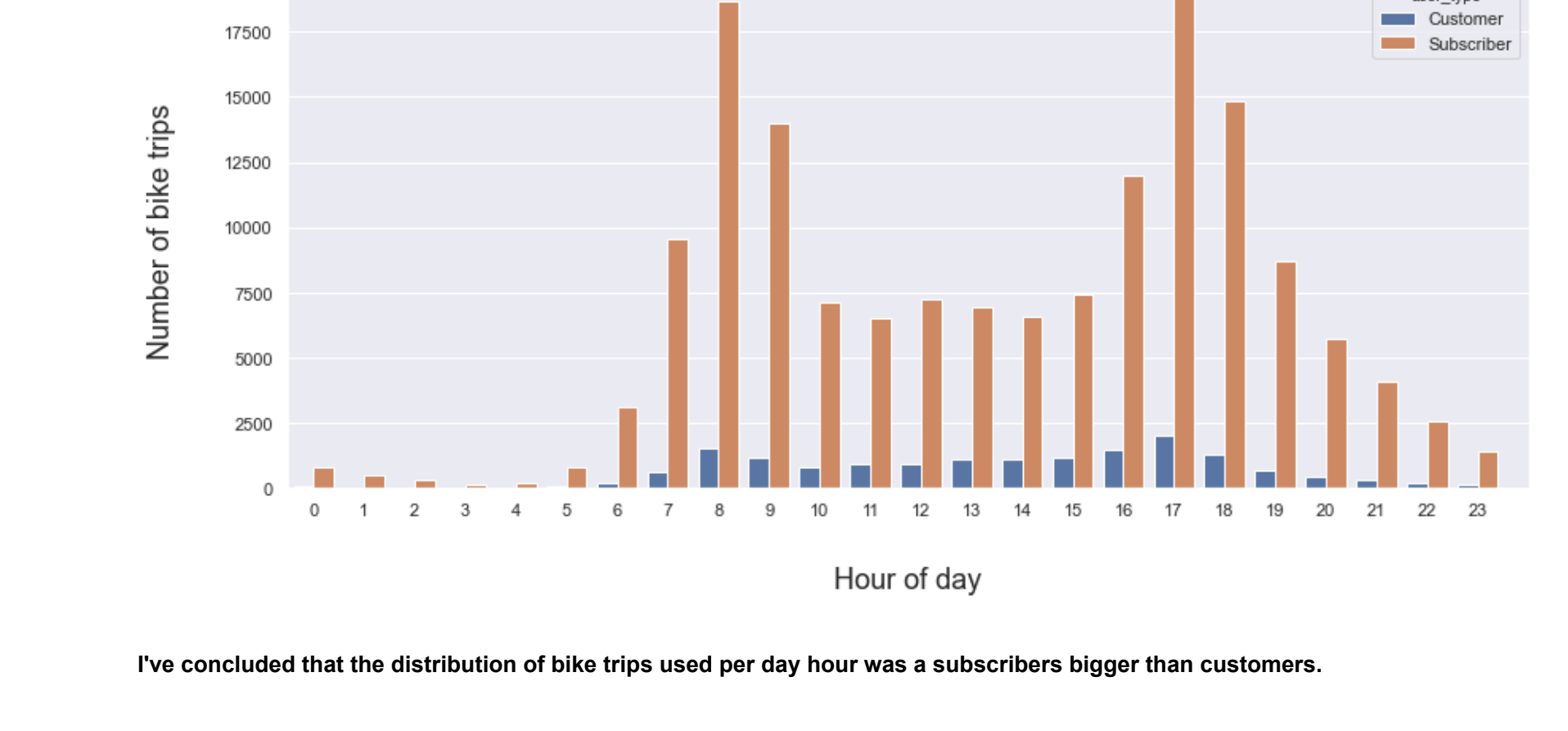


Number of bike trips gender per user type



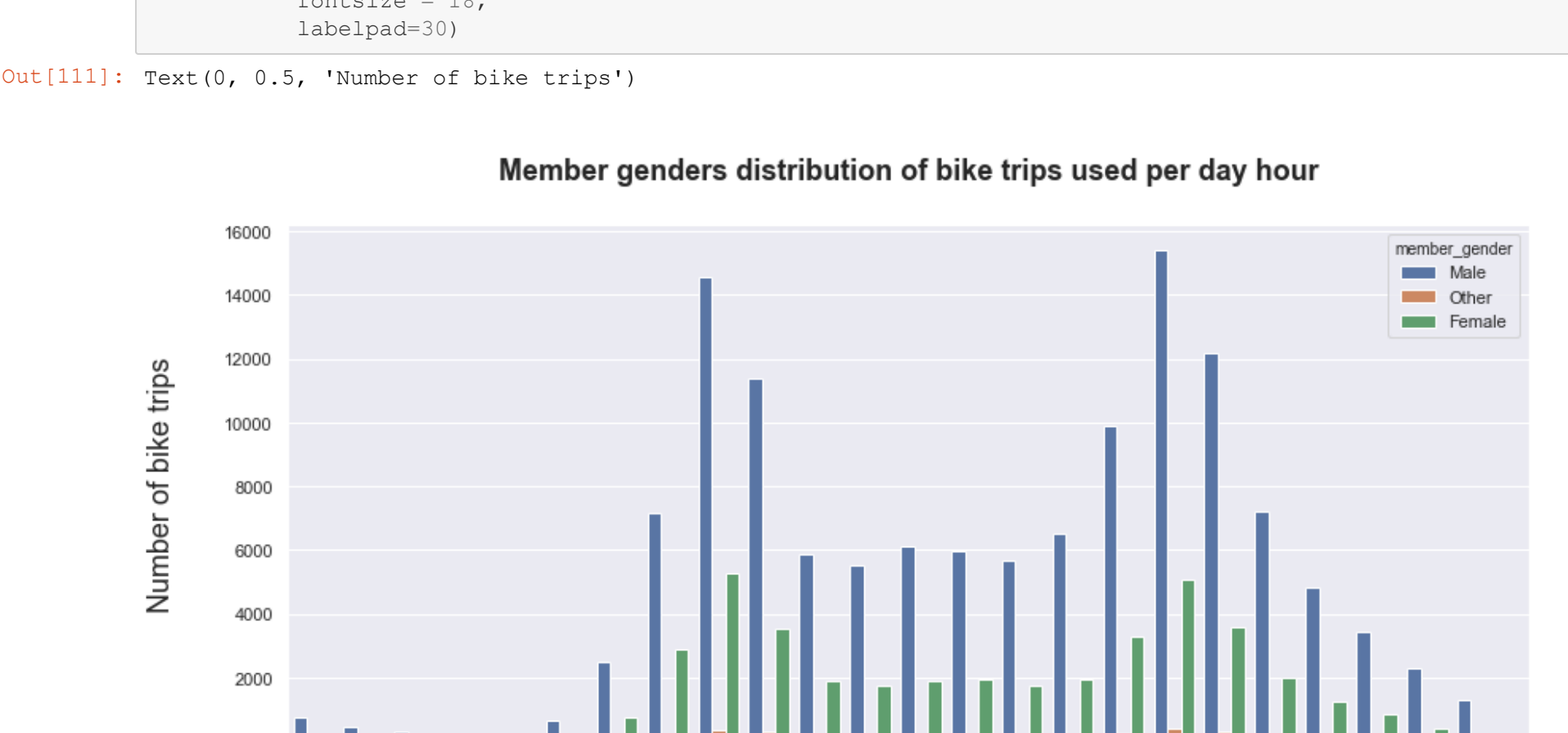
I've concluded that the number of bike trips gender per user type was a subscribers bigger than customers.

Customer vs. Subscriber distribution of bike trips used per day hour

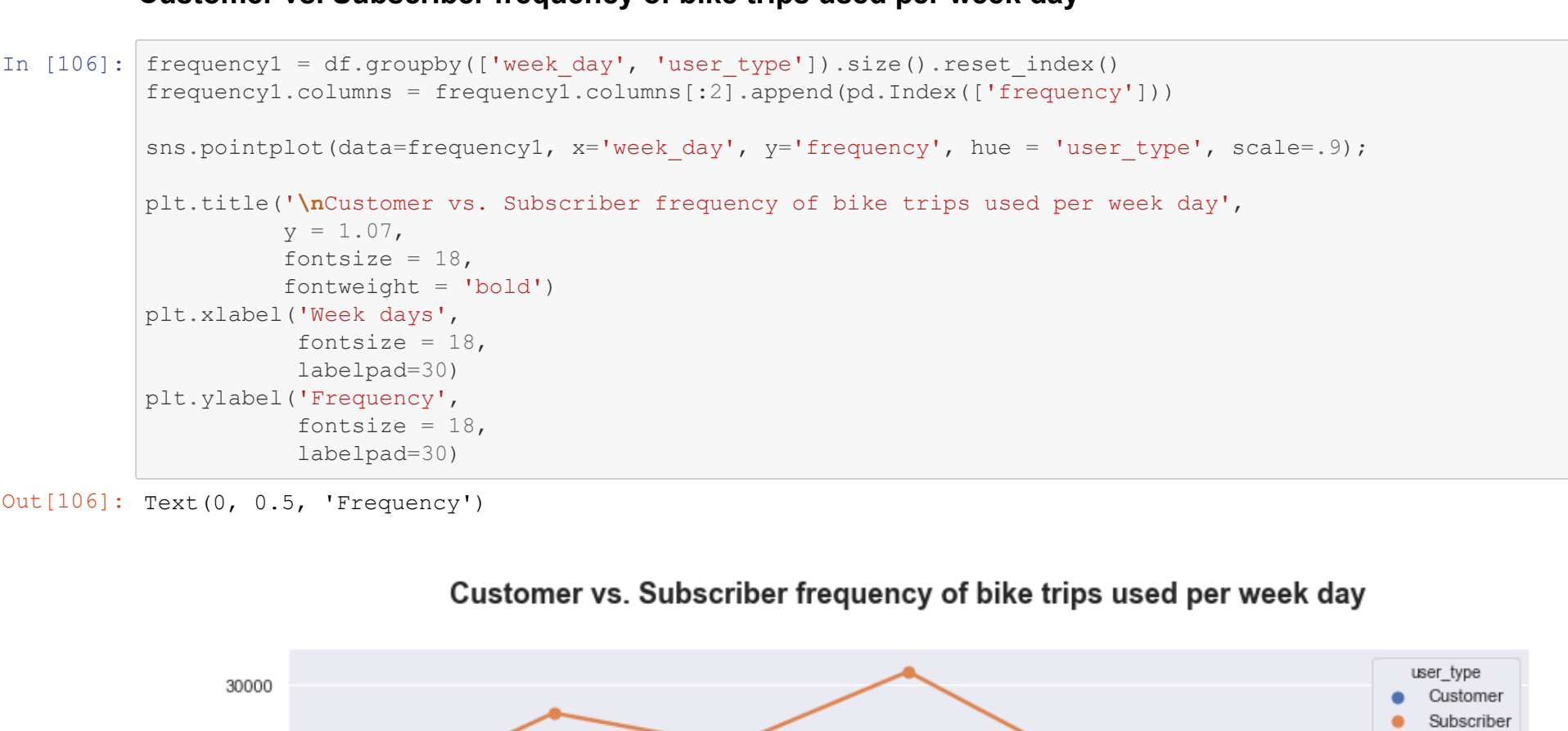


I've concluded that the distribution of bike trips used per day hour was a subscribers bigger than customers.

Member genders distribution of bike trips used per day hour

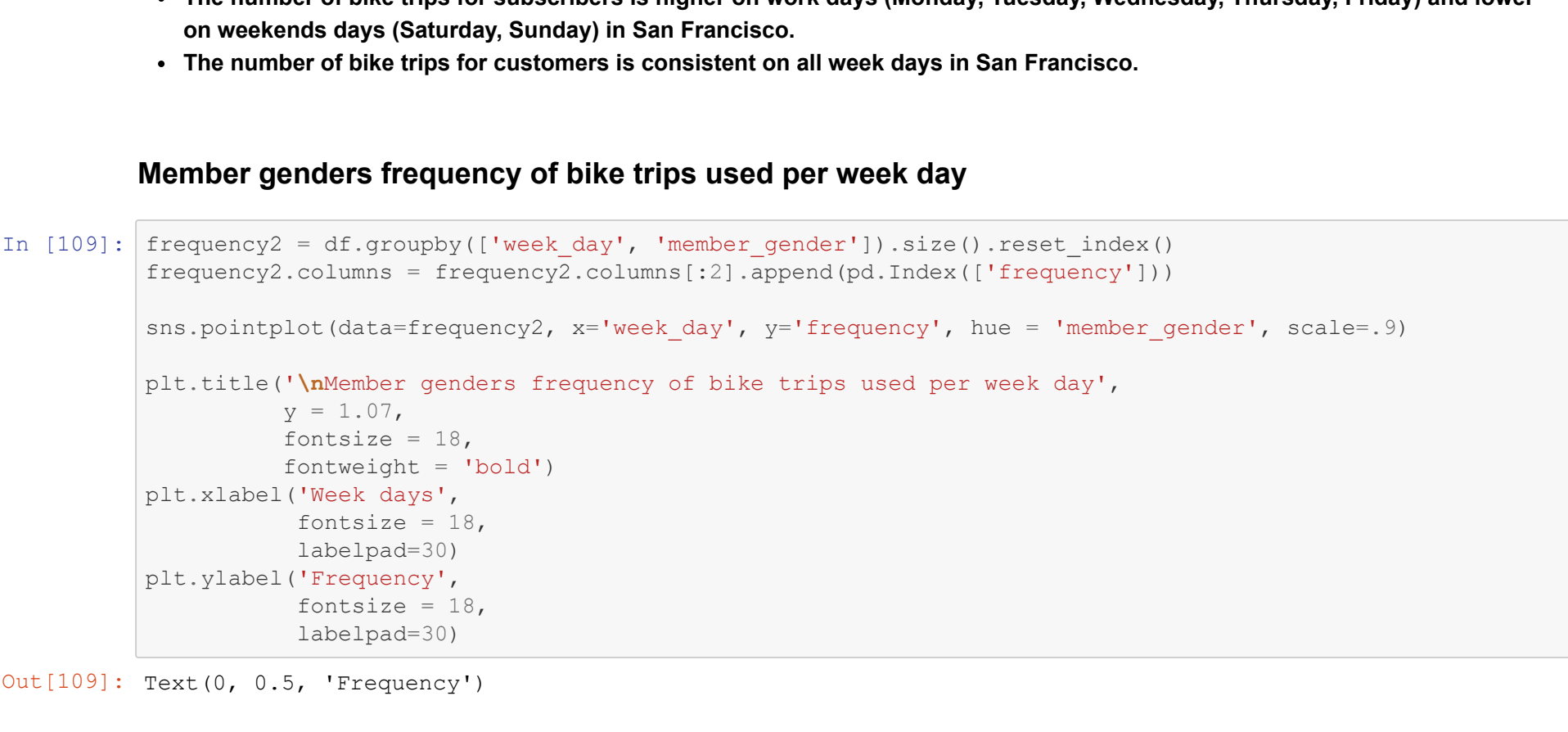


Customer vs. Subscriber frequency of bike trips used per week day



- The number of bike trips for subscribers is higher on work days (Monday, Tuesday, Wednesday, Thursday, Friday) and lower on weekends days (Saturday, Sunday) in San Francisco.
- The number of bike trips for customers is consistent on all week days in San Francisco.

Member genders frequency of bike trips used per week day



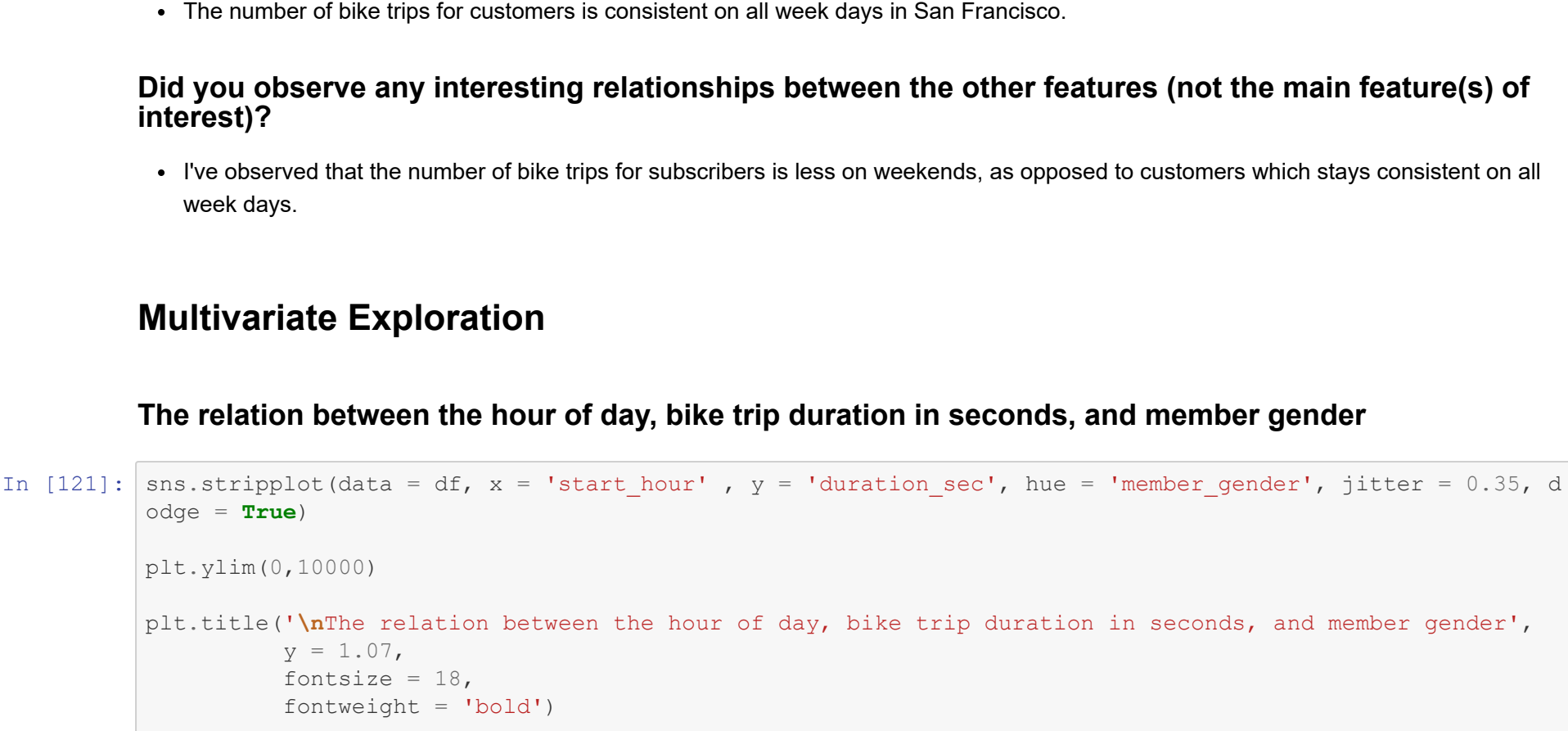
- The importance of including user type in data analysis revealed that there are some differences in the behavior of customers and subscribers
- The average bike trip duration in minutes per user type was (10.67922) are subscriber, while customer are (21.847536).
- The distribution of bike trip duration in minutes per user type was the distribution of customer trip time shows more variability than subscribers.
- The number of bike trips gender per user type was a subscribers bigger than customers.
- The distribution of bike trips used per day hour was a subscribers bigger than customers.
- The number of bike trips for subscribers is higher on work days (Monday, Tuesday, Wednesday, Thursday, Friday) and lower on weekends days (Saturday, Sunday) in San Francisco.
- The number of bike trips for customers is consistent on all week days in San Francisco.

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

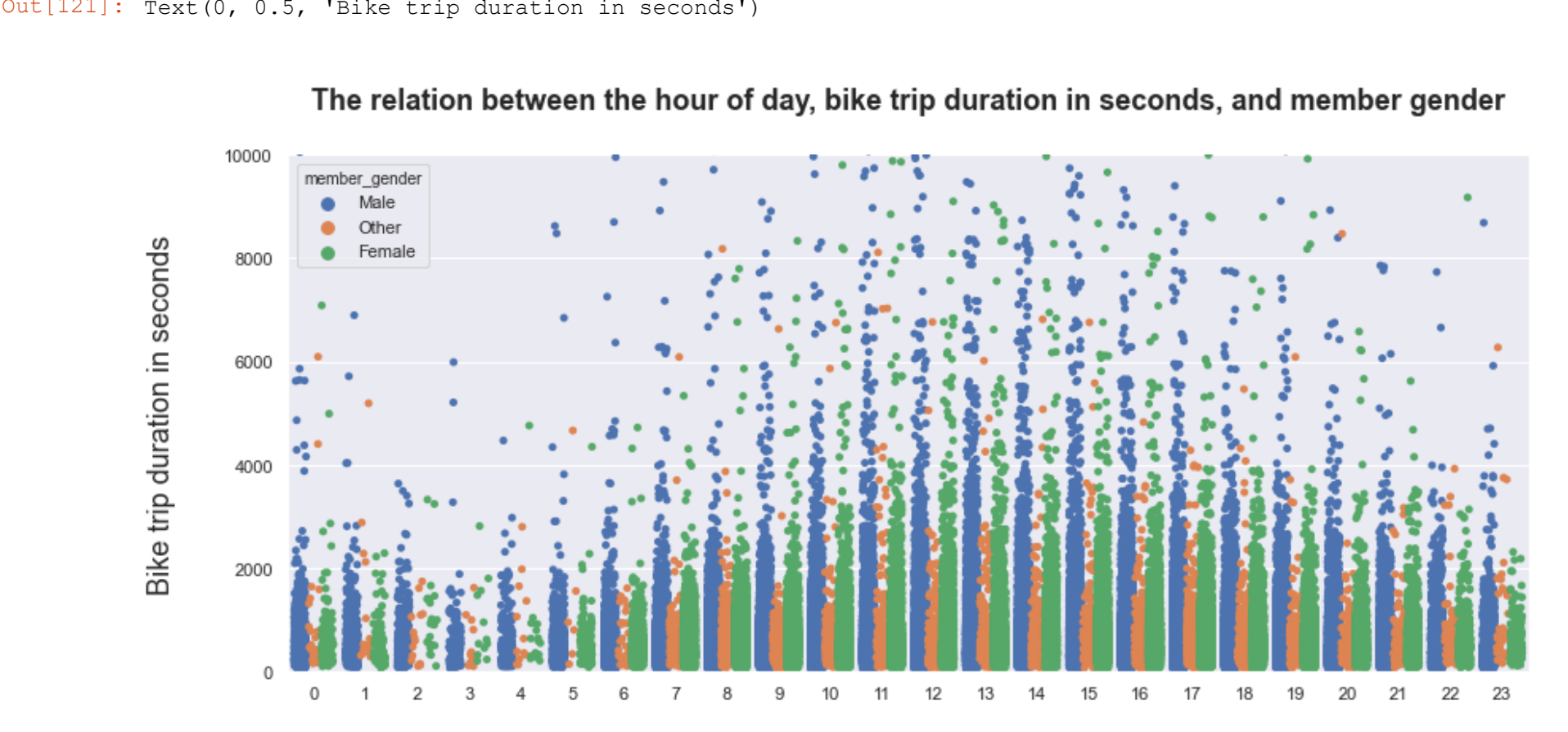
- I've observed that the number of bike trips for subscribers is less on weekends, as opposed to customers which stays consistent on all week days.

Multivariate Exploration

The relation between the hour of day, bike trip duration in seconds, and member gender



The relation between the hour of day, bike trip duration in seconds, and user type



Talk about some of the relationships you observed in this part of the investigation. Were there features that strengthened each other in terms of looking at your feature(s) of interest?

- I've observed that male customers take longer bike trips than female customers.
- I've observed that males take long bike trips from 11 AM to 4 PM.
- I've observed that subscribers take long bike trips from 11 AM to 5 PM.

Were there any interesting or surprising interactions between features?

- The interaction between the hour of the day, bike trip duration in seconds, and user type were surprising.
- The time between 11 AM to 5 PM is significantly different.

In [] :