

# Importing libraries

**Definition** - Importing Pandas, Numpy, Matplotlib, Seaborn, Glob and OS, for accessing, assessing, cleaning, Engineering, Analysing and Visualising Given Data.

In [ ]:

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sb
import glob
import os
from datetime import datetime, date
%matplotlib inline
```

# Gathering Data

**Definition** - Reading file from cloud stored csv file using file drive path . **Function Used** : pandas.read\_csv()

In [ ]:

```
df = pd.read_csv('/content/drive/MyDrive/Utkarsh doc/201902-fordgobike-tripdata.csv')
```

# Assessing Data

**Definition** - Assessing given Raw database using functions such as:

- dataframe.head()
- dataframe.tail()
- dataframe.shape - It is an attribute not a function.
- dataframe.info()
- dataframe.dtypes - It is an attribute not a function.
- dataframe.describe() - It is an attribute not a function.

In [ ]:

```
df.head()
```

Out[ ]:

	duration_sec	start_time	end_time	start_station_id	start_station_name	start_station
0	52185	2019-02-28 17:32:10.1450	2019-03-01 08:01:55.9750	21.0	Montgomery St BART Station (Market St at 2nd St)	37
1	42521	2019-02-28 18:53:21.7890	2019-03-01 06:42:03.0560	23.0	The Embarcadero at Steuart St	37
2	61854	2019-02-28 12:13:13.2180	2019-03-01 05:24:08.1460	86.0	Market St at Dolores St	37
3	36490	2019-02-28 17:54:26.0100	2019-03-01 04:02:36.8420	375.0	Grove St at Masonic Ave	37
4	1585	2019-02-28 23:54:18.5490	2019-03-01 00:20:44.0740	7.0	Frank H Ogawa Plaza	37

In [ ]:

```
df.tail()
```

Out[ ]:

	duration_sec	start_time	end_time	start_station_id	start_station_name	start_s
183407	480	2019-02-01 00:04:49.7240	2019-02-01 00:12:50.0340	27.0	Beale St at Harrison St	
183408	313	2019-02-01 00:05:34.7440	2019-02-01 00:10:48.5020	21.0	Montgomery St BART Station (Market St at 2nd St)	
183409	141	2019-02-01 00:06:05.5490	2019-02-01 00:08:27.2200	278.0	The Alameda at Bush St	
183410	139	2019-02-01 00:05:34.3600	2019-02-01 00:07:54.2870	220.0	San Pablo Ave at MLK Jr Way	
183411	271	2019-02-01 00:00:20.6360	2019-02-01 00:04:52.0580	24.0	Spear St at Folsom St	

In [ ]:

```
df.shape
```

Out[ ]:

(183412, 16)

In [ ]:

df.dtypes

Out[ ]:

```

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

```

In [ ]:

df.describe()

Out[ ]:

	duration_sec	start_station_id	start_station_latitude	start_station_longitude	end_station_id
<b>count</b>	183412.000000	183215.000000	183412.000000	183412.000000	183215.000000
<b>mean</b>	726.078435	138.590427	37.771223	-122.352664	136.240427
<b>std</b>	1794.389780	111.778864	0.099581	0.117097	111.515151
<b>min</b>	61.000000	3.000000	37.317298	-122.453704	3.000000
<b>25%</b>	325.000000	47.000000	37.770083	-122.412408	44.000000
<b>50%</b>	514.000000	104.000000	37.780760	-122.398285	100.000000
<b>75%</b>	796.000000	239.000000	37.797280	-122.286533	235.000000
<b>max</b>	85444.000000	398.000000	37.880222	-121.874119	398.000000

### What is the structure of your dataset?

The dataset includes 183,412 trips 'rows' with 16 features 'columns'. Out of 16 features, seven are float64, two are int64, and seven objects. Also, start and end time have the wrong datatype( string instead of date-time object).

### What is/are the main feature(s) of interest in your dataset?

According to me, the main features are the duration of the trip, origination and conclusion of trips, the user-type and most used stations.

**What features in the dataset do you think will help support your investigation into your feature(s) of interest?**

- duration\_sec
- start\_time
- end\_time
- user\_type
- start\_station\_name
- end\_station\_name

## Cleaning Data

**1. Definition - Making a copy of provided dataframe to retain original values for future comparison.**

### Code

In [ ]:

```
df_copy = df.copy()
```

**2. Definition - Start and End times of trips are given in String format and thus need to be changed into Date-time object as it would further facilitate the calculations and visualisation process.**

### Code

In [ ]:

```
#changing datatype of 'start_time' & 'end_time' into datetime object  
df.start_time = pd.to_datetime(df.start_time)  
df.end_time = pd.to_datetime(df.end_time)
```

**testing acquired Changes**

In [ ]:

```
df.dtypes
```

Out[ ]:

```
duration_sec          int64
start_time            datetime64[ns]
end_time              datetime64[ns]
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
```

### 3.Definition - Ahead, Checking for data redundancy and duplicacy in particular.

#### Code

In [ ]:

```
#check if any rows are duplicated
sum(df.duplicated())
```

Out[ ]:

0

#### Test Result -

No Duplicacy Found

### 4.Definition - Extracting Hours and minutes from Date-Time object("start\_time" and "end\_time") into new added respectable columns for further extraction of data into subpart for reducing complexity and better visualization.

#### Code

In [ ]:

```
#extract the hours from start time
df['start_time_hours']=df['start_time'].dt.hour
#extract the minute from start time
df['start_time_minutes']=df['start_time'].dt.minute + df['start_time_hours']*60
```

In [ ]:

```
#extract the hours from end time
df['end_time_hours']=df['end_time'].dt.hour
#extract the minute from end time
df['end_time_minutes']=df['end_time'].dt.minute + df['start_time_hours']*60
```

## Testing above made changes

In [ ]:

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 183412 entries, 0 to 183411
Data columns (total 20 columns):
 #   Column                                Non-Null Count  Dtype
---  -
 0   duration_sec                          183412 non-null  int64
 1   start_time                            183412 non-null  datetime64[ns]
 2   end_time                              183412 non-null  datetime64[ns]
 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
16   start_time_hours                       183412 non-null  int64
17   start_time_minutes                     183412 non-null  int64
18   end_time_hours                         183412 non-null  int64
19   end_time_minutes                       183412 non-null  int64
dtypes: datetime64[ns](2), float64(7), int64(6), object(5)
memory usage: 28.0+ MB
```

## 5.Definition - Finally Extracting Month Number from given timestamps.

(We do not need to extract Year as this data is collected from 2019 only.)

## Code

In [ ]:

```
df['month']=df.start_time.dt.month
```

## Testing above made changes

In [ ]:

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 183412 entries, 0 to 183411
Data columns (total 21 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   duration_sec                          183412 non-null  int64
1   start_time                            183412 non-null  datetime64[ns]
2   end_time                              183412 non-null  datetime64[ns]
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
16  start_time_hours                       183412 non-null  int64
17  start_time_minutes                     183412 non-null  int64
18  end_time_hours                         183412 non-null  int64
19  end_time_minutes                       183412 non-null  int64
20  month                                  183412 non-null  int64
dtypes: datetime64[ns](2), float64(7), int64(7), object(5)
memory usage: 29.4+ MB
```

## 6.Definition -

- Dropping irrelevant columns from our Data Base.
- Dropping Rows with Null values for removing all NAN values from our database.

## Code

In [ ]:

```
#dropping irrelevant columns
df.drop(['start_station_latitude', 'start_station_longitude', 'end_station_latitude', 'end_station_longitude'], axis =1 , inplace = True)
```

In [ ]:

```
df=df.dropna()
```

## Testing above made alteration.

In [ ]:

```
df.isna().sum()
```

Out[ ]:

```
duration_sec          0
start_time            0
end_time              0
start_station_id      0
start_station_name    0
end_station_id        0
end_station_name      0
bike_id              0
user_type             0
member_birth_year     0
member_gender         0
bike_share_for_all_trip 0
start_time_hours      0
start_time_minutes    0
end_time_hours        0
end_time_minutes      0
month                0
dtype: int64
```

## Storing Data

### Saving Engineered Database in a csv file

In [ ]:

```
df.to_csv("fordbike_engineered.csv")
```

## Analyzing Data

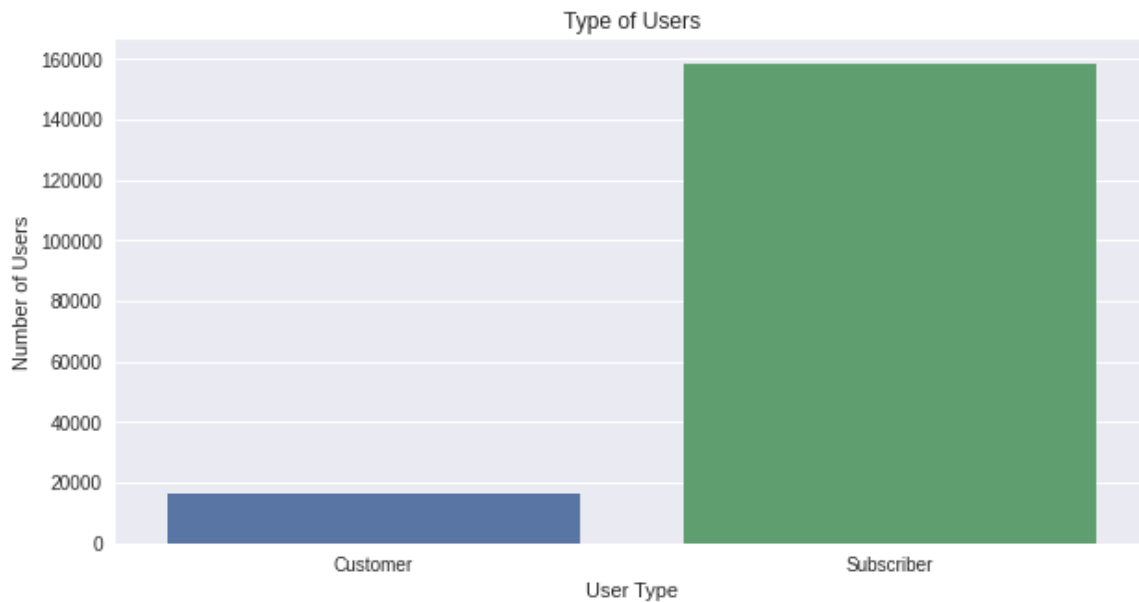
### Univariate Exploration

#### Analyse the ratio of User Type : Customer/Subscriber



In [79]:

```
#BoxPlot representing User Type : Customer/Subscriber  
plt.figure(figsize = (10, 5))  
sb.countplot(data=df,x='user_type')  
plt.title("Type of Users")  
plt.xlabel('User Type')  
plt.ylabel('Number of Users')  
plt.show()
```

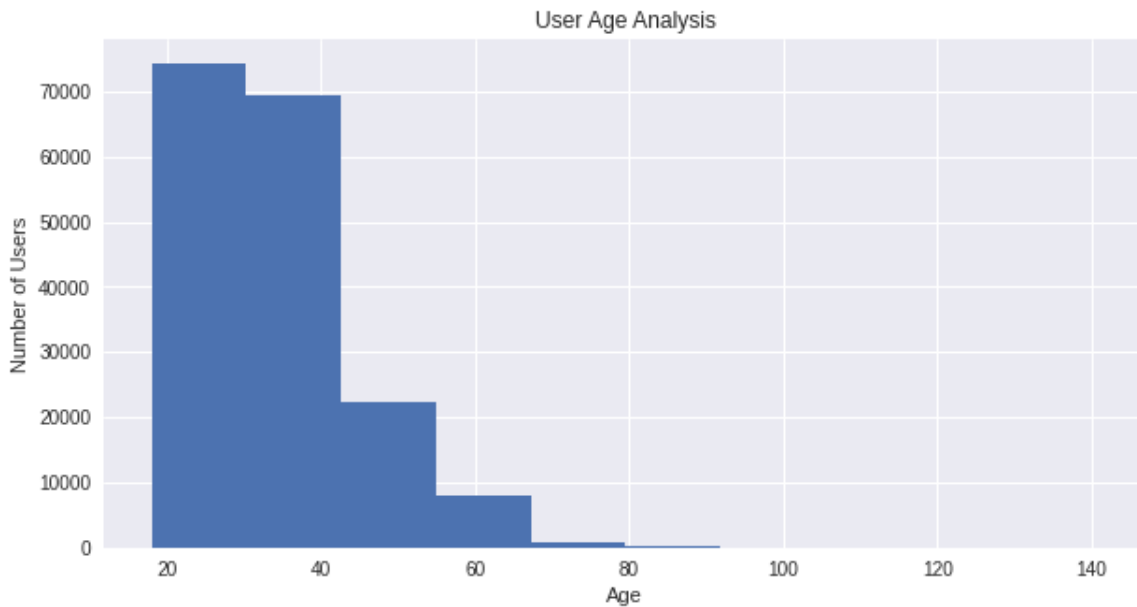


Above chart suggests that majority of users are Subscribers while, a lesser number of users fall under 'Customer' category

**Which Age group comprises most of our User Database?**

In [ ]:

```
#Histogram Plot depicting member's spread on the basis of their Age
df['Age'] = 2019- df['member_birth_year']
plt.figure(figsize=(10,5))
plt.hist(data=df,x='Age')
plt.title('User Age Analysis')
plt.xlabel('Age')
plt.ylabel('Number of Users')
plt.show()
```



Above histogram clearly shows,

- People of age '18' to '25', comprise most of our Users.
- We see clear fall in user count as their Age progress.

**Function to plot Univariate Barplots with specified X-Label.**

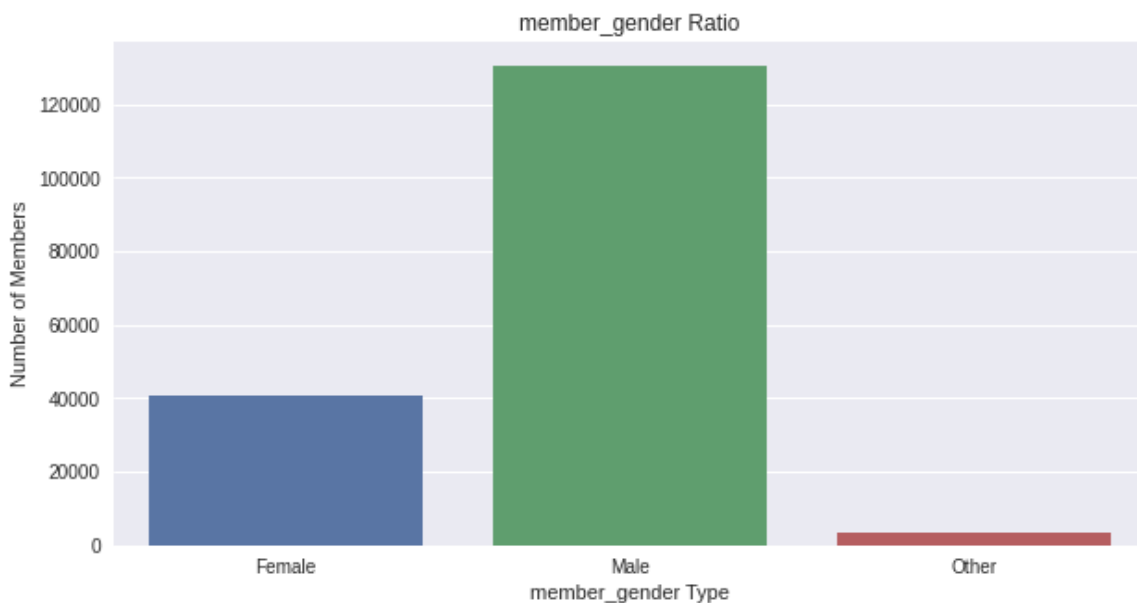
In [96]:

```
#Function to plot Univariate Barplots with specified X-Label.  
def funcbarplot(df, x):  
    mg=df.groupby(x).size()  
    plt.figure(figsize = (10, 5))  
    sb.barplot(x = mg.index, y = mg.values)  
    plt.title(f"{x} Ratio")  
    plt.xlabel(f'{x} Type')  
    plt.ylabel('Number of Members')  
    plt.show()
```

## Which Gender type comprises most of our Users?

In [97]:

```
#BarPlot representing member gender ratio  
funcbarplot(df, 'member_gender')
```



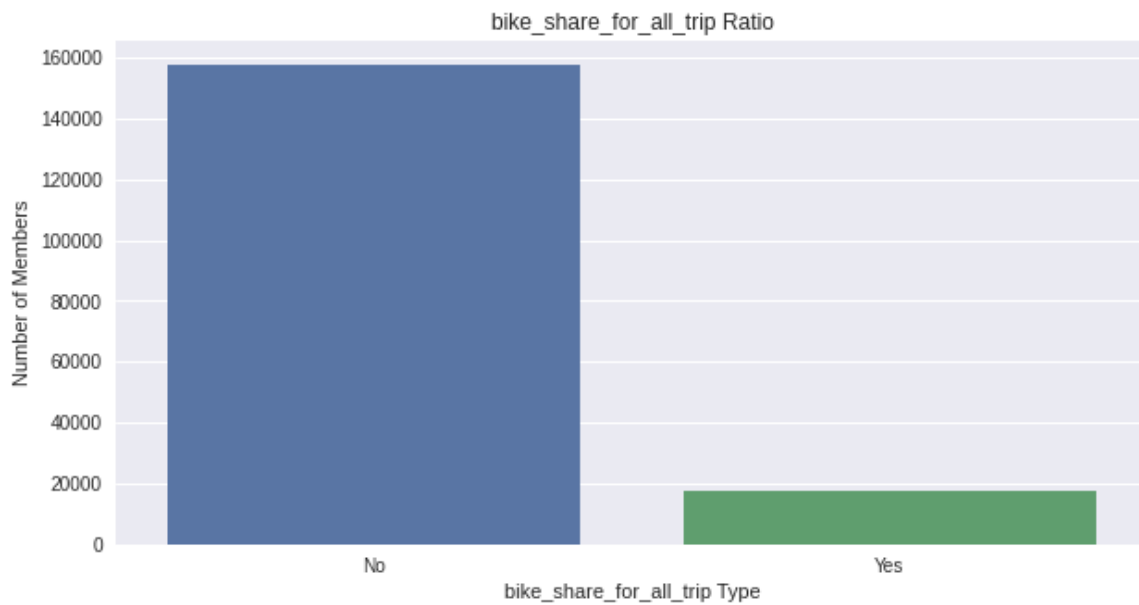
Above BarPlot shows,

- "Male" users dominate in number, with a total of 135000(approx) Users.
- "Female" users fall 2nd in count ratio, with a total number close to 40000.
- Users of "Other" gender type are least in number. Their total number lies close to 3000 users.

## Do users prefer to share their ride with other users?

In [58]:

```
#BarPlot representing member gender ratio
funcbarplot(df, 'bike_share_for_all_trip')
```



No, We can clearly see, majority of our users do not prefer sharing their rides with other users.

## Function to plot barchart of Station on basis of trips onboarded and deboarded.

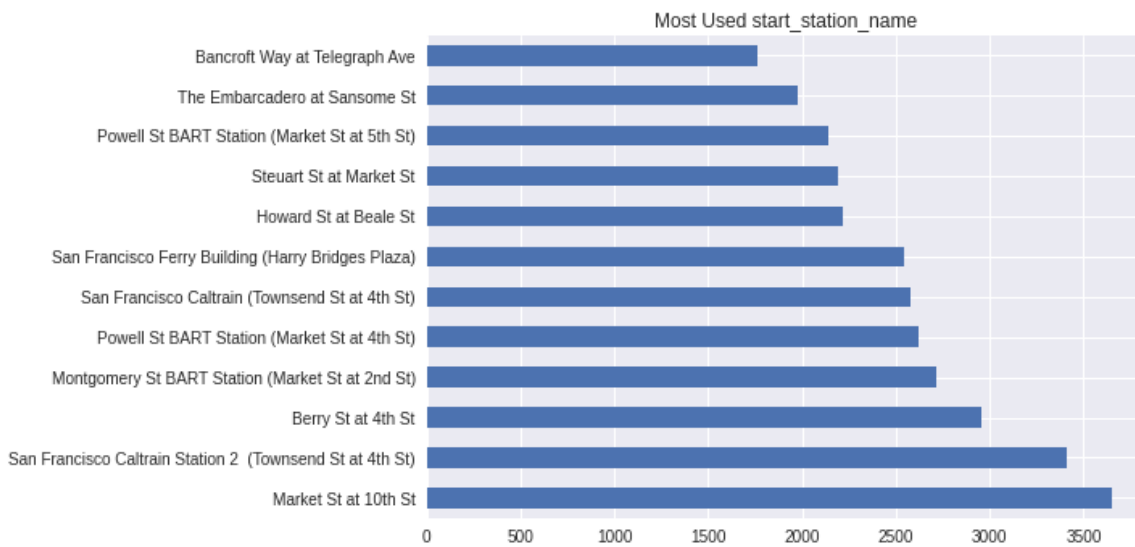
In [71]:

```
def funcbarh(df, x):
    df[x].value_counts()[ :12 ].plot(kind='barh')
    plt.title(f'Most Used {x}')
    plt.show()
```

## A Barchart for depicting journey origination frequency of all stations.

In [74]:

```
#A Barchart for depicting journey origination frquency of all stations.
funcbarh(df, 'start_station_name')
```

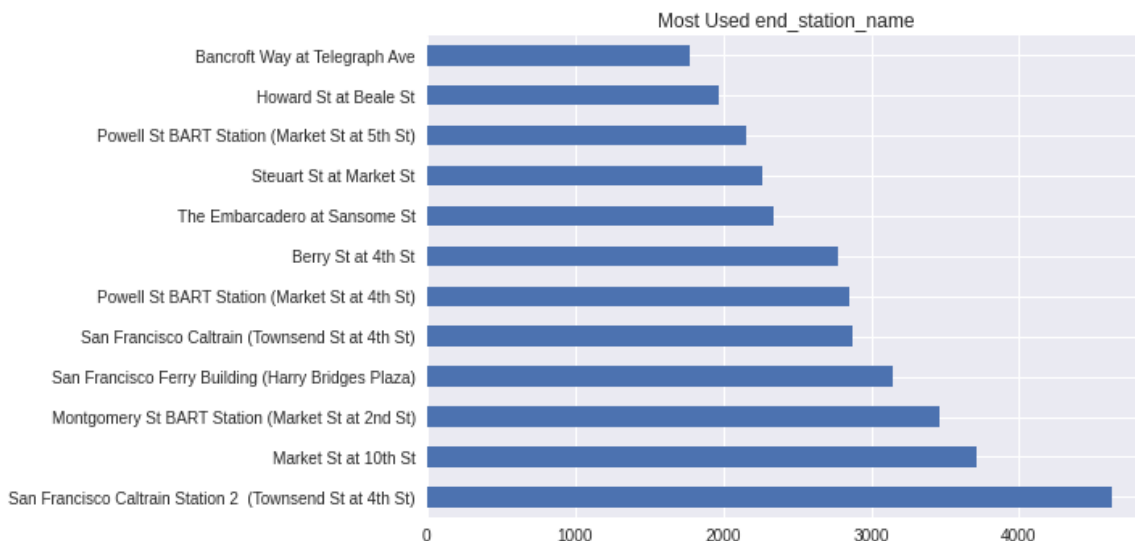


Market st at 10th station is the most used as trip origination station.

## A Barchart for depicting journey conclusion frquency of all stations.

In [75]:

```
#A Barchart for depicting journey conclusion frquency of all stations.
funcbarh(df, 'end_station_name')
```

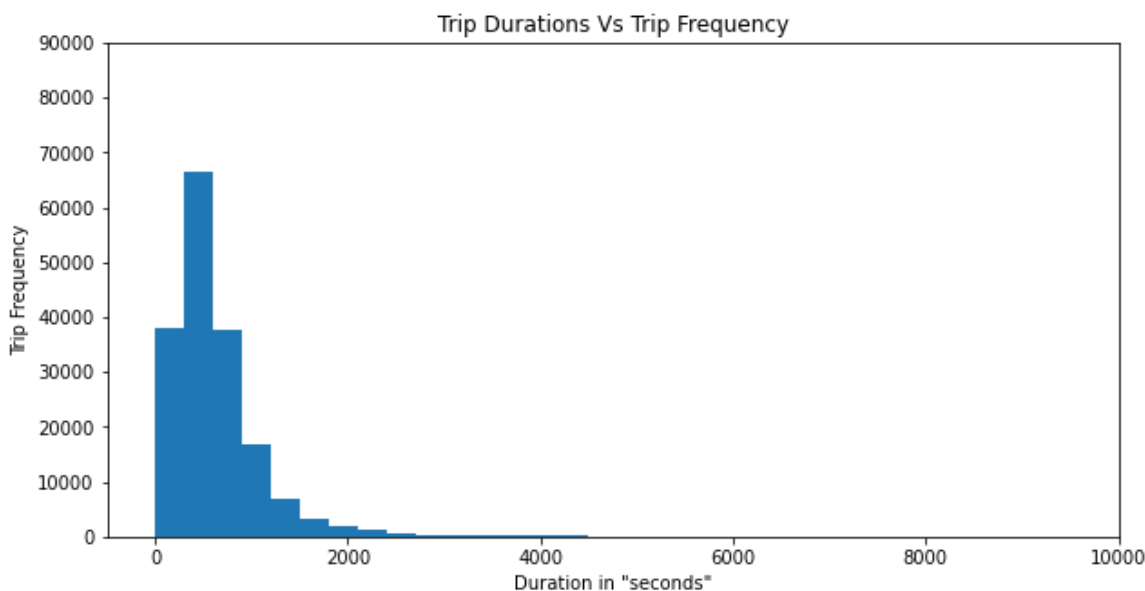


San Francisco Caltrain Station 2 is most used as trip conclusion station.

## Histogram plot for plotting Frequency Vs time duration in seconds

In [ ]:

```
#Histogram plot for plotting Frequency Vs time duration in seconds
plt.figure(figsize=(10,5))
duration_bins = np.arange(1,df.duration_sec.max()+300,300)
plt.hist(data=df,x='duration_sec',bins=duration_bins)
plt.title('Trip Durations Vs Trip Frequency')
plt.xlabel('Duration in "seconds"')
plt.ylabel('Trip Frequency')
plt.axis([-500, 10000, 0, 90000])
plt.show()
```



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

The graph shows the most common trip durationz take 0-2000 seconds. The number of subscribers are significantly higher than customers in User-types. also, Market st at 10th street is the most used as trip origination station. And, San Francisco Caltrain Station 2 is most used as trip conclusion station.

**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?**

**YES!** There were some data redundancy and distribution issues in the used dataset and needed some cleaning. The problems addressed are as follows:

- Start and End times of trips are given in String format and thus need to be changed into Date-time object as it would further facilitate the calculations and visualisation process.
- Extracting Hours and minutes from Date-Time object("start\_time" and "end\_time") into new added respectable columns for further extraction of data into subpart for reducing complexity and better visualization.
- Extracting Month Number from given timestamps. (We do not need to extract Year as this data is collected from 2019 only.)
- Dropping irrelevant columns from our Data Base.
- Dropping Rows with Null values for removing all NAN values from our database.

## Bivariate Exploration

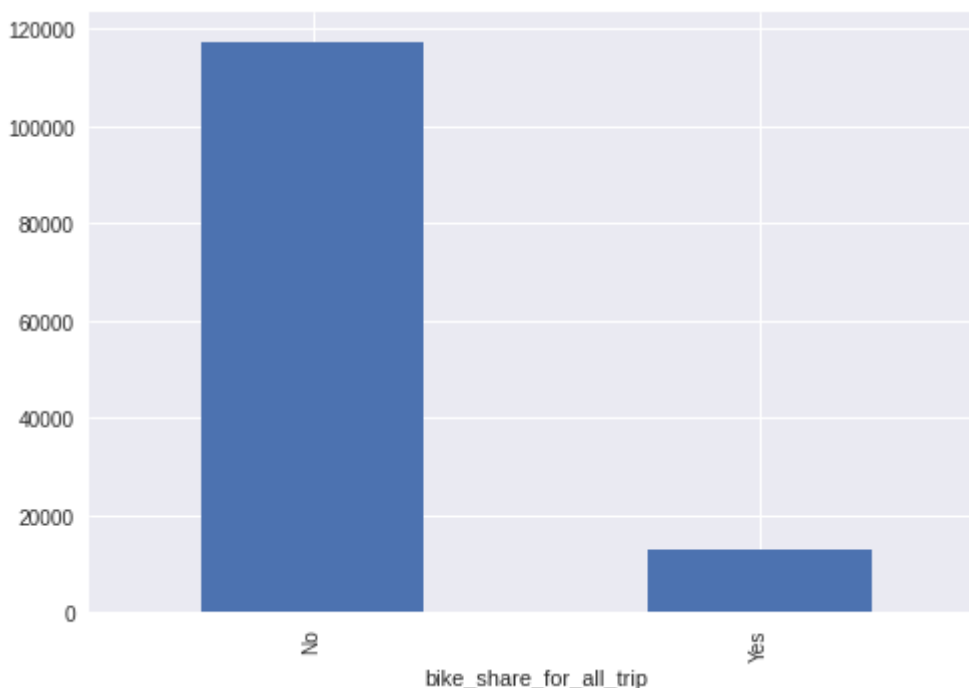
**Discuss trip sharing tendency of Users of different gender types.**

In [122]:

```
#Function to plot Bike Sharing Tendency of users
def funcMemberG(df,x):
    df['Memberg'] = df['member_gender'].mask(df['member_gender'].ne(x))
    df.groupby(['bike_share_for_all_trip'])['Memberg'].count().plot.bar()
```

In [123]:

```
#Function to Plot Bike Sharing Tendency of MALE users
funcMemberG(df,"Male")
```

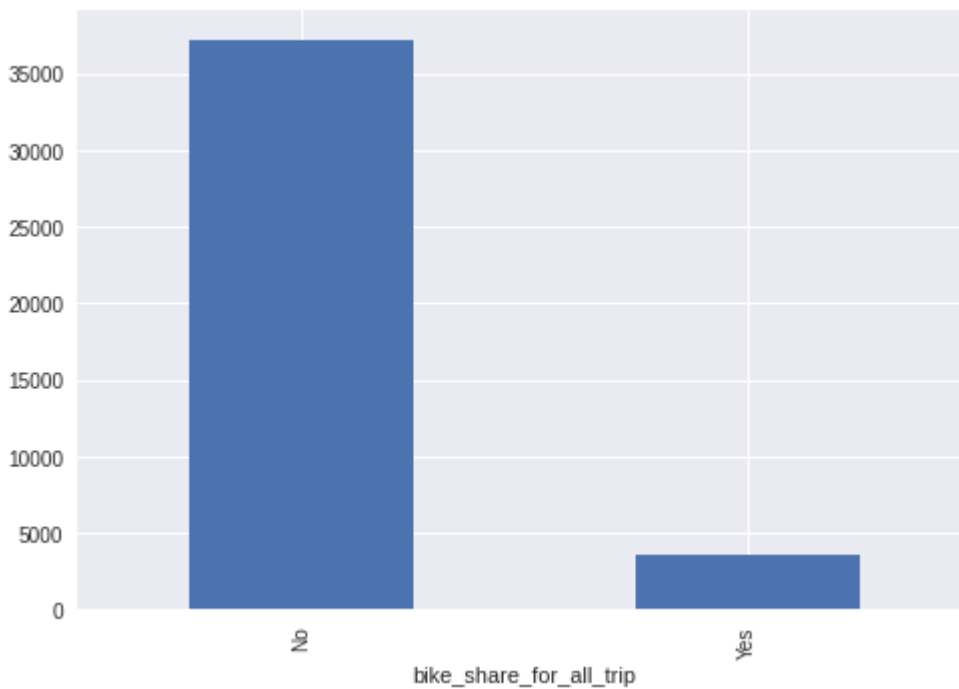


We can clearly infer,

- Users with "Male" gender type are dominant in number.
- Their total number is close to 120000 users.
- Approximately 17500 Users of "Male" gender prefer to share their ride.
- Above given point indicates that such users have medium tendency of all three genders to share their ride. Since approximately 1/6th of them share their ride.

In [124]:

```
#Function to Plot Bike Sharing Tendency of FEMALE users  
funcMemberG(df, "Female")
```



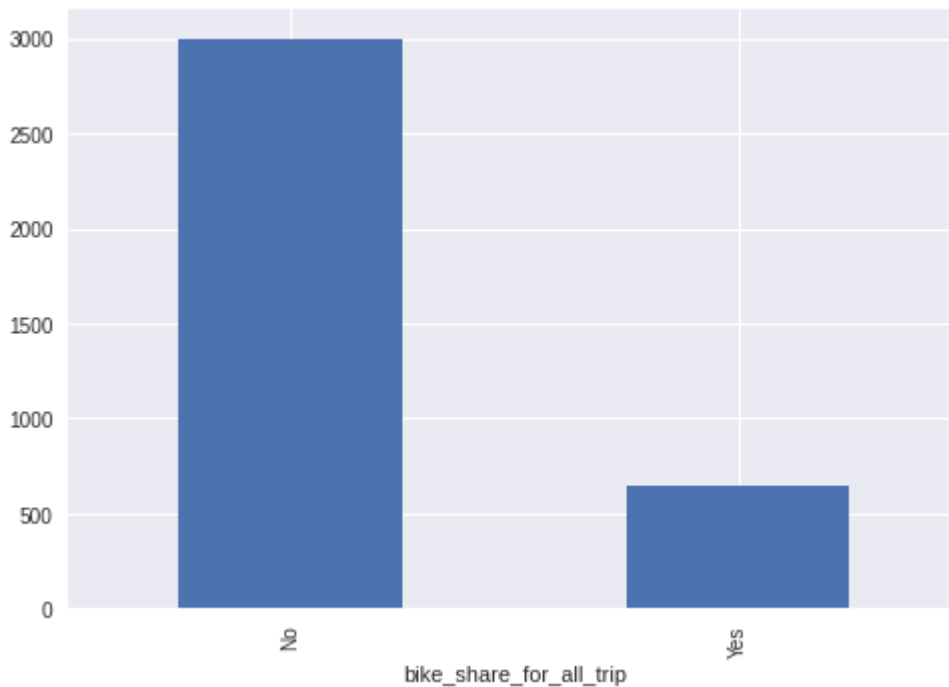
We can clearly infer,

- Users with "Female" gender type are less in number than "Male" users.
- Their total number is close to 37500 users.
- Approximately 4000 Users of "Female" gender prefer to share their ride.
- Above given point indicates that such users have least tendency to share their ride. Since approximately only 1/9th of them share their ride.



In [125]:

```
#Function to Plot Bike Sharing Tendency of OTHER users
funcMemberG(df, "Other")
```



We can clearly infer,

- Users with "Other" gender type are least in number.
- Their total number is close to 3000 users.
- Approximately 600 Users of "Other" gender prefer to share their ride.
- Above given point indicates that such users have maximum tendency to share their ride. Since approximately 1/5th of them share their ride.

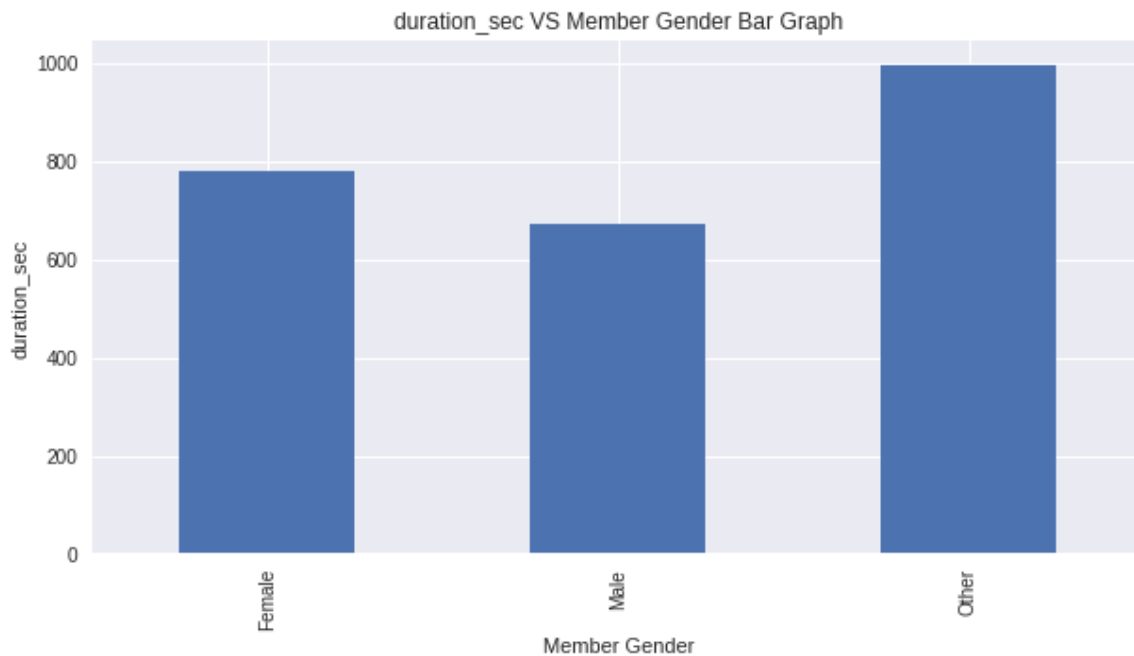
## Discuss trip duration tendency of Users of different gender types.

In [91]:

```
#Function to Plot Trip Duration Mean for Each Gender
def funcplot(df, x):
    plt.figure(figsize = (10, 5))
    df.groupby(['member_gender'])[x].mean().plot.bar()
    plt.title(f'{x} VS Member Gender Bar Graph')
    plt.xlabel('Member Gender')
    plt.ylabel(f'{x}')
    plt.show()
```

In [92]:

```
funcplot(df, 'duration_sec')
```



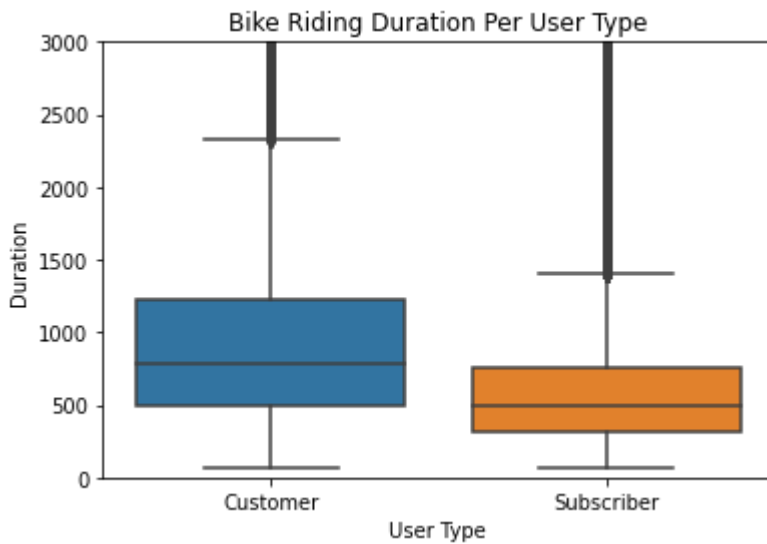
We can clearly infer,

- Users of "Other" gender type tend to have longer trip duration.
- Female users travel more than Male users.
- Men tend to have travel for least duration. (We are all the same and lazy !)

**Discuss trip duration tendency of Users of different User types.**

In [ ]:

```
# Boxplot showing trip duration per user type, ie; Customer/Subscriber
sb.boxplot(data=df, x='user_type', y='duration_sec')
plt.ylim(0,3000)
plt.title("Bike Riding Duration Per User Type")
plt.xlabel('User Type')
plt.ylabel('Duration')
plt.show()
```



We can infer here,

- Customers tend to travel for longer durations.
- Subscribers travel for comparatively less duration.

**Discuss trip duration tendency of Users from top 3 most Onboarded stations.**

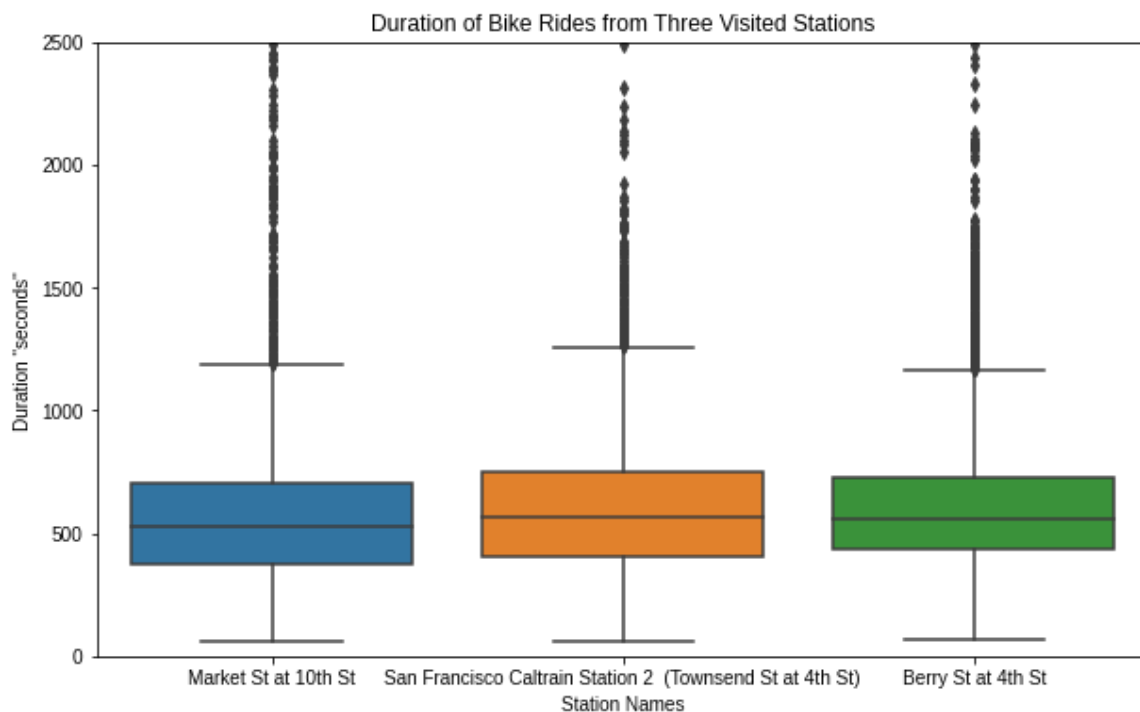
In [ ]:

```
#boxplot for trip duration analysis of journeys from top three most used station
s.
start=df['start_station_name'].value_counts().index[:3]
s_stations = df.loc[df['start_station_name'].isin(start)]

plt.figure(figsize=(10,6))
sb.boxplot(data = s_stations, x='start_station_name', y='duration_sec')
plt.ylim(0, 2500)
plt.style.use('seaborn')
plt.title('Duration of Bike Rides from Three Onboarded Stations')
plt.xlabel('Station Names')
plt.ylabel('Duration "seconds"')
```

Out[ ]:

```
Text(0, 0.5, 'Duration "seconds"')
```



We can infer that,

- All stations share close to common trip durations.
- Trips started from Station 'San Francisco Caltrain Station 2', tend to have longest duration of the three.

**Discuss traffic tendency of Users over 24hr period.**

In [ ]:

```
start=df['start_station_name'].value_counts().index[:5]
s_stations = df.loc[df['start_station_name'].isin(start)]
plt.figure(figsize = (15,6))
sb.countplot(data=s_stations, x='start_time_hours', hue='start_station_name')
x_tick= np.arange(0,24,2)
x_label= [str(x)+":00" for x in x_tick]
plt.xticks(x_tick, x_label)
plt.xlabel('Time in Hours')
plt.ylabel('Users')
plt.show()
```