

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [92]: df= pd.read_excel ("D:\assignment\facebook_user_data.xlsx")
```

```
In [93]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 99093 entries, 0 to 99092
Data columns (total 15 columns):
#   Column                Non-Null Count  Dtype  ---  ---  ---  ---  ---  ---  ---  ---  ---  ---  ---  ---  ---  ---
0   user_id                99093 non-null  int64  
1   age                   99093 non-null  int64  
2   dob_day               99093 non-null  int64  
3   dob_year              99093 non-null  int64  
4   dob_month             99093 non-null  int64  
5   gender                98826 non-null  object 
6   tenure                99091 non-null  float64
7   friend_count          99093 non-null  int64  
8   friendships_initiated 99093 non-null  int64  
9   likes                 99093 non-null  int64  
10  likes_received        99093 non-null  int64  
11  mobile_likes          99093 non-null  int64  
12  mobile_likes_received 99093 non-null  int64  
13  www_likes             99093 non-null  int64  
14  www_likes_received    99093 non-null  int64  
dtypes: float64(1), int64(13), object(1)
memory usage: 11.3+ MB
```

```
In [4]: df.head(5)
```

	user_id	age	dob_day	dob_year	dob_month	gender	tenure	friend_count	friendships_initiated	likes	likes_received	mobile_likes	mobile_likes_received	www_likes	www_likes_received
	0	2094382	14	19	1999	11	male	266.0	0	0	0	0	0	0	0
1	1192001	14	2	1999	11	female	6.0	0	0	0	0	0	0	0	0
2	2083894	14	16	1999	11	male	13.0	0	0	0	0	0	0	0	0
3	1203168	14	25	1999	12	female	93.0	0	0	0	0	0	0	0	0
4	1733186	14	4	1999	12	male	82.0	0	0	0	0	0	0	0	0

Imputation of missing values:

```
In [5]: df.columns
```

```
Out[5]: Index(['user_id', 'age', 'dob_day', 'dob_year', 'dob_month', 'gender', 'tenure', 'friend_count', 'friendships_initiated', 'likes', 'likes_received', 'mobile_likes', 'mobile_likes_received', 'www_likes', 'www_likes_received', 'dtype: object'])
```

```
In [6]: df.isna().sum()
```

```
Out[6]: user_id      0
age           0
dob_day       0
dob_year      0
dob_month     0
gender        179
tenure        2
friend_count   0
friendships_initiated  0
likes         0
likes_received 0
mobile_likes   0
mobile_likes_received 0
www_likes      0
www_likes_received 0
dtype: int64
```

Replace the null values (NA) of gender column with its mode or median and explain why mode/median used to replace NA values

```
In [19]: df.duplicated().sum()
```

```
Out[19]: 0
```

```
In [8]: df['gender'].fillna(df['gender'].mode,inplace= True)
```

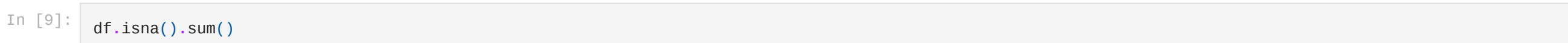
ANS- Its a categorical variable so i used mode.

```
In [9]: df.isna().sum()
```

```
Out[9]: user_id      0
age           0
dob_day       0
dob_year      0
dob_month     0
gender        0
tenure        2
friend_count   0
friendships_initiated  0
likes         0
likes_received 0
mobile_likes   0
mobile_likes_received 0
www_likes      0
www_likes_received 0
dtype: int64
```

Replace the null values (NA) of tenure column (numerical variable) with its median, and explain why mode/median used to replace NA values

```
In [10]: sns.displot(df.tenure)
```



```
In [11]: df['tenure'].fillna(df['tenure'].median,inplace= True)
```

When the data is skewed, it is good to consider using the median value for replacing the missing values.

```
In [12]: df.isna().sum()
```

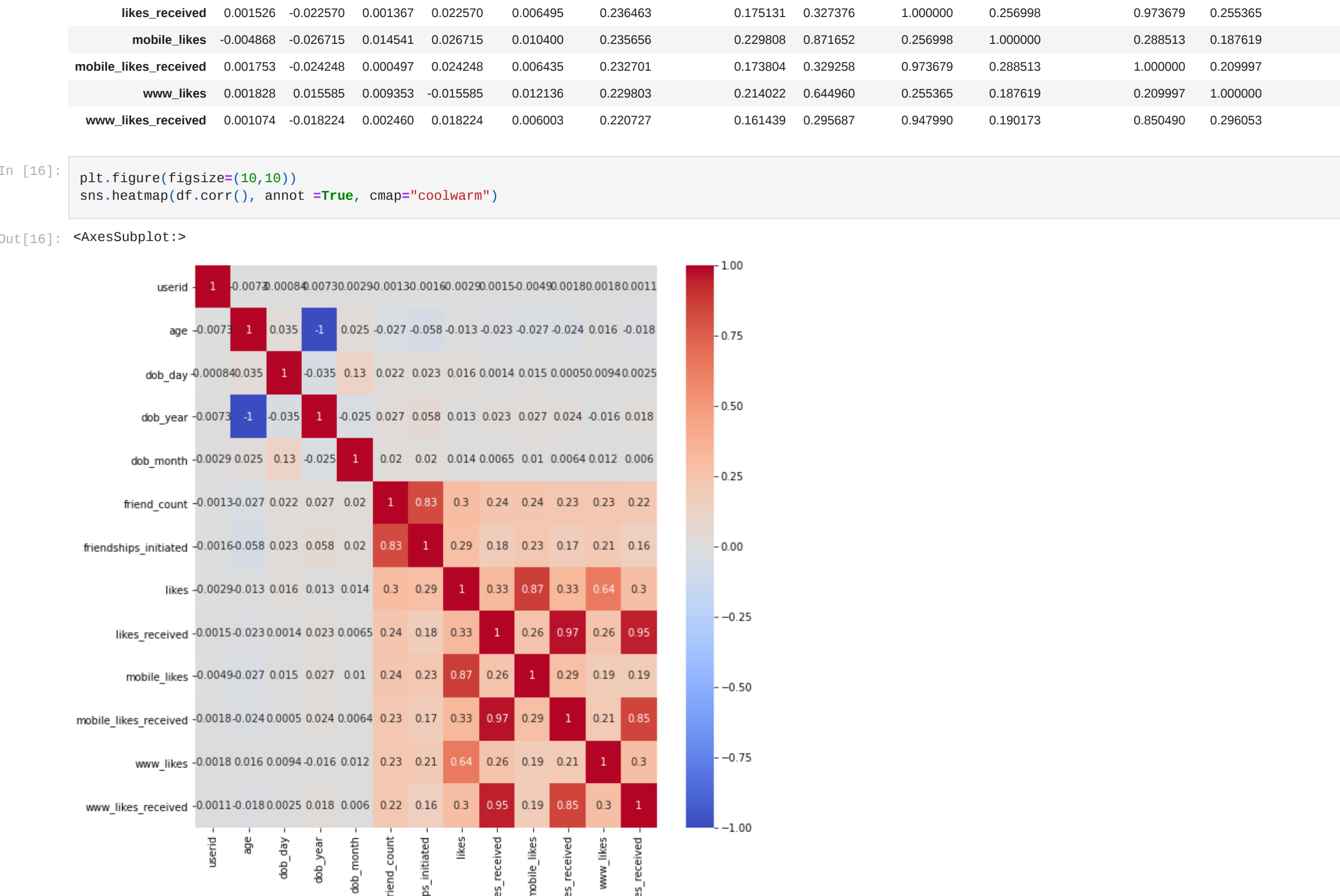
```
Out[12]: user_id      0
age           0
dob_day       0
dob_year      0
dob_month     0
gender        0
tenure        0
friend_count   0
friendships_initiated  0
likes         0
likes_received 0
mobile_likes   0
mobile_likes_received 0
www_likes      0
www_likes_received 0
dtype: int64
```

2) Plot heatmap / correlation matrix on all the columns.

```
In [13]: df.corr()
```

	user_id	age	dob_day	dob_year	dob_month	friend_count	friendships_initiated	likes	likes_received	mobile_likes	mobile_likes_received	www_likes	www_likes_received
user_id	1.000000	-0.007265	-0.000839	0.007265	0.002924	-0.001314	0.002924	-0.001591	-0.002875	0.001526	-0.004868	0.001753	0.001828
age	-0.007265	1.000000	0.035035	-1.000000	0.025167	-0.027407	-0.027407	-0.058059	-0.013009	-0.022570	-0.026715	-0.024248	0.015585
dob_day	-0.000839	0.035035	1.000000	-0.030035	0.129443	0.021961	0.022999	0.015980	0.001367	0.014541	0.000497	0.009353	0
dob_year	0.007265	-1.000000	-0.035035	1.000000	-0.025167	0.027407	0.027407	0.058059	0.013009	0.022570	0.026715	-0.024248	-0.015585
dob_month	0.002924	0.025167	0.129443	-0.025167	1.000000	0.019804	0.019804	0.020075	0.014147	0.006495	0.010400	0.006435	0.012136
friend_count	-0.001314	-0.027407	0.021961	0.027407	0.019804	1.000000	0.825850	0.298017	0.236463	0.235656	0.232701	0.229803	0
friendships_initiated	-0.001591	-0.058059	0.022999	0.058059	0.020075	0.019804	1.000000	0.285592	0.175131	0.229808	0.173904	0.214022	0
likes	-0.002875	-0.013009	0.021590	0.013009	0.014147	0.298017	0.285592	1.000000	0.327376	0.871652	0.329258	0.644960	0
likes_received	0.001526	-0.022570	0.001367	0.022570	0.006495	0.236463	0.229808	0.327376	1.000000	0.256998	0.973679	0.255365	0
mobile_likes	-0.004868	-0.026715	0.014541	0.026715	0.010400	0.235656	0.229808	0.871652	0.256998	1.000000	0.288513	0.187619	0
mobile_likes_received	0.001753	-0.024248	0.000497	0.024248	0.006435	0.232701	0.173904	0.329258	0.973679	0.288513	1.000000	0.209997	0
www_likes	0.001828	0.015585	0.009353	-0.015585	0.012136	0.229803	0.214022	0.644960	0.255365	0.187619	0.209997	1.000000	0
www_likes_received	0.001074	-0.018224	0.002460	-0.018224	0.006003	0.220277	0.161439	0.295687	0.947990	0.190173	0.850490	0.296053	1

```
In [16]: plt.figure(figsize=(10,10))
sns.heatmap(df.corr(), annot= True,cmap="coolwarm")
```



### 3) Analysis based on gender of the users

• What is composition of male and female users?

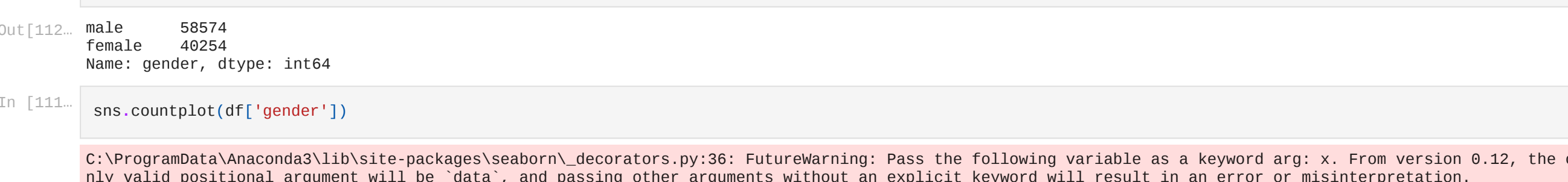
```
In [112]: df['gender'].value_counts()
```

```
Out[112]: male      58574
female    40254
Name: gender, dtype: int64
```

```
In [111]: sns.countplot(df['gender'])
```

C:\ProgramData\Anaconda3\lib\site-packages\seaborn\decorators.py:96: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.

```
Out[111]: <AxesSubplot: xlabel='gender', ylabel='count'>
```



Which category of gender has more friends?

```
In [113]: df.groupby('gender').agg({'friend_count': ['sum']})
```

gender	friend_count
female	9740259
male	9666787

```
In [114]: plt.figure(figsize=(5,5))
sns.barplot(x=df.gender,y=df.friend_count)
```

```
Out[114]: <AxesSubplot: xlabel='gender', ylabel='friend_count'>
```



ANS- categories of Female has more friends.

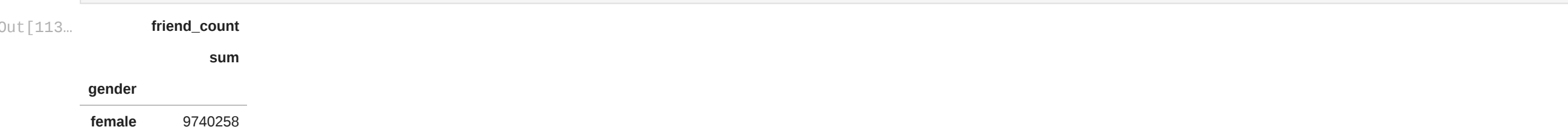
### Which category of gender initiated more friendships?

```
In [116]: df.groupby('gender').agg({'friendships_initiated': ['sum']})
```

gender	friendships_initiated
female	4584694
male	6037023

```
In [124]: plt.figure(figsize=(5,5))
sns.barplot(x=df.gender,y=df.friendships_initiated)
```

```
Out[124]: <AxesSubplot: xlabel='gender', ylabel='friendships_initiated'>
```



ANS- Female calagory has more friends initiated than male.

### What is the distribution of tenure across different categories of gender?

```
In [117]: df.groupby('gender').agg({'tenure': ['sum']})
```

gender	tenure
female	23637151.0
male	29238972.0

```
In [126]: plt.figure(figsize=(5,5))
sns.barplot(y=df.tenure,x=df.gender)
```

```
Out[126]: <AxesSubplot: xlabel='gender', ylabel='tenure'>
```



ANS- Females are spend more Number of days active on Facebook than males.

### 4) Analysis based on the least active users on Facebook

• How many users have no friends?

```
In [96]: df.friend_count.value_counts()
```

0	1962
1	1816
2	1117
3	189
4	789
...	...
3299	1
4576	1
1359	1
4384	1
4687	1

Name: friend\_count, Length: 2562, dtype: int64

ANS-1962 users dont have friends

### How many users did not like any posts?

```
In [104]: df.likes.value_counts()
```

0	22308
1	6929
2	4434
3	3240
4	2597
...	...
2887	1
3467	1
6119	1
4868	1
2947	1

Name: likes, Length: 2924, dtype: int64

ANS-22308 users did not like any post

### How many users did not receive any likes?

```
In [102]: df.likes_received.value_counts()
```

0	24428
1	7395
2	4541
3	3347
4	2669
...	...
1610	1
3859	1
3723	1
2660	1
2947	1

Name: likes\_received, Length: 2681, dtype: int64

ANS-24428 no of users didnt receive any likes

### Analysis based on the user accessibility (Mobile Devices vs. Web Devices)

What is the average number of posts liked by users (based on gender) through web vs mobile devices?

```
In [153]: a=df.groupby('gender').agg({'mobile_likes': ['mean']})
a
```

gender	mobile_likes
female	177.912938
male	60.261328

```
In [154]: b=df.groupby('gender').agg({'www_likes': ['mean']})
b
```

gender	www_likes
female	87.138297
male	24.616550

```
In [157]: plt.figure(figsize=(5,5))
x=["mobile_like_by_female","web_likes_by_female"]
y=list(a.loc['female'])[0].list(b.loc['female'])[0]
plt.bar(x,y)
```

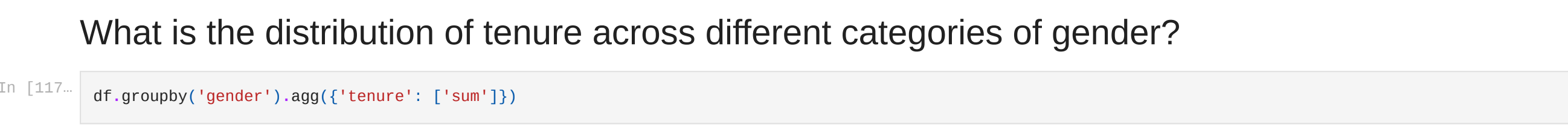
C:\ProgramData\Anaconda3\lib\site-packages\numpy\lib\stride\_tricks.py:536: VisibleDeprecationWarning: Creating an ndarray from ragged nested sequences (which is a list-or-tuple of lists-or-tuples-or ndarrays with different lengths or shapes) is deprecated. If you meant to do this, you must specify 'dtype=object' when creating the ndarray.

args = (np.array([], copy=False, subok=subok) for \_ in args)

C:\ProgramData\Anaconda3\lib\site-packages\numpy\core\asarray.py:102: VisibleDeprecationWarning: Creating an ndarray from ragged nested sequences (which is a list-or-tuple of lists-or-tuples-or ndarrays with different lengths or shapes) is deprecated. If you meant to do this, you must specify 'dtype=object' when or eating the ndarray.

return array(a, dtype, copy=False, order=order)

<BarContainer object of 2 artists>



In [161]: plt.figure(figsize=(5,5))
x=["mobile\_like\_by\_male","web\_likes\_by\_male"]
y=list(a.loc['male'])[0].list(b.loc['male'])[0]
plt.bar(x,y)

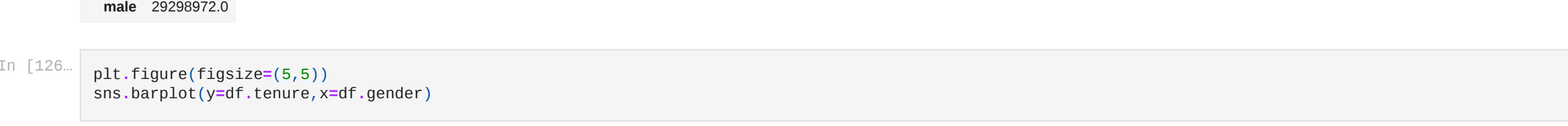
C:\ProgramData\Anaconda3\lib\site-packages\numpy\lib\stride\_tricks.py:536: VisibleDeprecationWarning: Creating an ndarray from ragged nested sequences (which is a list-or-tuple of lists-or-tuples-or ndarrays with different lengths or shapes) is deprecated. If you meant to do this, you must specify 'dtype=object' when creating the ndarray.

args = (np.array([], copy=False, subok=subok) for \_ in args)

C:\ProgramData\Anaconda3\lib\site-packages\numpy\core\asarray.py:102: VisibleDeprecationWarning: Creating an ndarray from ragged nested sequences (which is a list-or-tuple of lists-or-tuples-or ndarrays with different lengths or shapes) is deprecated. If you meant to do this, you must specify 'dtype=object' when or eating the ndarray.

return array(a, dtype, copy=False, order=order)

<BarContainer object of 2 artists>



What is the average number of likes received by users (based on gender) through web vs. mobile devices?

```
In [162]: d=df.groupby('gender').agg({'mobile_likes_received': ['mean']})
d
```

gender	mobile_likes_received
female	147.80084
male	40.833015

```
In [163]: e=df.groupby('gender').agg({'www_likes_received': ['mean']})
e
```

gender	www_likes_received
female	104.334451
male	27.078533

```
In [164]: plt.figure(figsize=(5,5))
x=["mobile_like_recive_by_male","web_likes_recive_by_male"]
y=list(d.loc['male'])[0].list(e.loc['male'])[0]
plt.bar(x,y)
```

C:\ProgramData\Anaconda3\lib\site-packages\numpy\lib\stride\_tricks.py:536: VisibleDeprecationWarning: Creating an ndarray from ragged nested sequences (which is a list-or-tuple of lists-or-tuples-or ndarrays with different lengths or shapes) is deprecated. If you meant to do this, you must specify 'dtype=object' when or eating the ndarray.

args = (np.array([], copy=False, subok=subok) for \_ in args)

C:\ProgramData\Anaconda3\lib\site-packages\numpy\core\asarray.py:102: VisibleDeprecationWarning: Creating an ndarray from ragged nested sequences (which is a list-or-tuple of lists-or-tuples-or ndarrays with different lengths or shapes) is deprecated. If you meant to do this, you must specify 'dtype=object' when or eating the ndarray.

return array(a, dtype, copy=False, order=order)

<BarContainer object of 2 artists>



```
In [165]: plt.figure(figsize=(5,5))
x=["mobile_like_recive_by_female","web_likes_recive_by_female"]
y=list(d.loc['female'])[0].list(e.loc['female'])[0]
plt.bar(x,y)
```

C:\ProgramData\Anaconda3\lib\site-packages\numpy\lib\stride\_tricks.py:536: VisibleDeprecationWarning: Creating an ndarray from ragged nested sequences (which is a list-or-tuple of lists-or-tuples-or ndarrays with different lengths or shapes) is deprecated. If you meant to do this, you must specify 'dtype=object' when or eating the ndarray.

args = (np.array([], copy=False, subok=subok) for \_ in args)

C:\ProgramData\Anaconda3\lib\site-packages\numpy\core\asarray.py:102: VisibleDeprecationWarning: Creating an ndarray from ragged nested sequences (which is a list-or-tuple of lists-or-tuples-or ndarrays with different lengths or shapes) is deprecated. If you meant to do this, you must specify 'dtype=object' when or eating the ndarray.

return array(a, dtype, copy=False, order=order)

<BarContainer object of 2 artists>

