

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from scipy.stats import binom
6 from scipy.stats import poisson
7 from scipy.stats import norm
8 from matplotlib.ticker import ScalarFormatter
9 sns.set_palette('bright')

```

```

1 df = pd.read_csv('walmart_data.csv')
2 df.head(2)

```

	User_ID	Product_ID	Gender	Age	Occupation	City_Category	Stay_In_Current_City_Years
0	1000001	P00069042	F	0-17	10	A	

```
1 df.shape
```

```
(550068, 10)
```

```
1 df.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 550068 entries, 0 to 550067
Data columns (total 10 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   User_ID          550068 non-null   int64  
 1   Product_ID       550068 non-null   object  
 2   Gender           550068 non-null   object  
 3   Age              550068 non-null   object  
 4   Occupation       550068 non-null   int64  
 5   City_Category    550068 non-null   object  
 6   Stay_In_Current_City_Years  550068 non-null   object  
 7   Marital_Status   550068 non-null   int64  
 8   Product_Category 550068 non-null   int64  
 9   Purchase         550068 non-null   int64  
dtypes: int64(5), object(5)
memory usage: 42.0+ MB

```

```
1 df.nunique()
```

```

User_ID            5891
Product_ID        3631
Gender             2
Age                7
Occupation         21
City_Category      3
Stay_In_Current_City_Years  5
Marital_Status     2
Product_Category   20
Purchase           18105
dtype: int64

```

There are 5891 distinct customers, 3631 distinct products, 21 different occupations, 3 city category, 5 distinct Stay In Current City in Years, 20 product categories.

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

```

User_ID            0
Product_ID        0
Gender             0
Age                0
Occupation         0
City_Category      0
Stay_In_Current_City_Years  0
Marital_Status     0
Product_Category   0
Purchase           0
dtype: int64

```

- There are NO NULL values in the data.

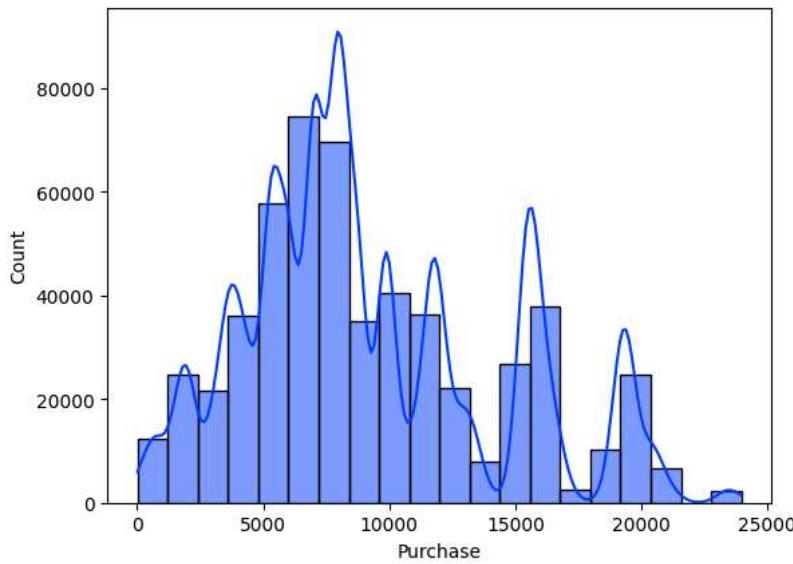
```
1 #Total Sales till date is
```

```
2 df['Purchase'].sum()
```

```
5095812742
```

```
1 sns.histplot(df, x= 'Purchase' , kde= True , bins = 20)
```

```
2 plt.show()
```

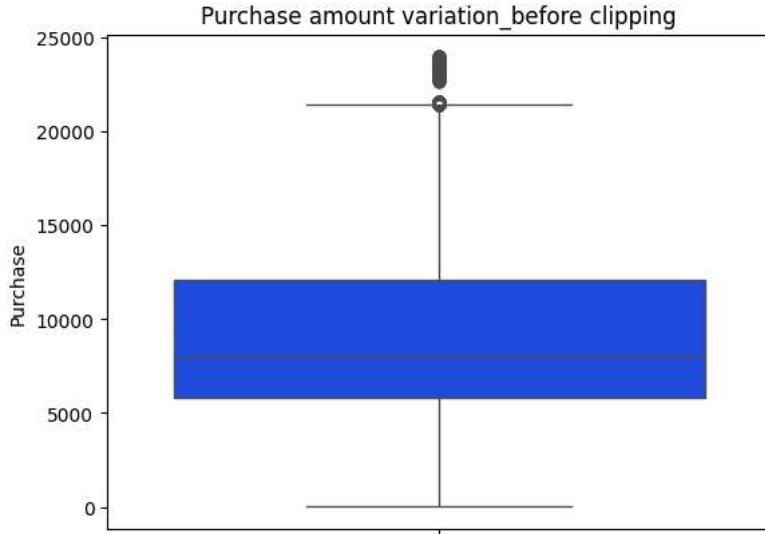


- Total Sales till date is 5,095,812,742 \$

```
1 sns.boxplot(df['Purchase'])
```

```
2 plt.title('Purchase amount variation_before clipping')
```

```
3 plt.show()
```



```
1 a,b = np.percentile(df['Purchase'],5), np.percentile(df['Purchase'],95)
```

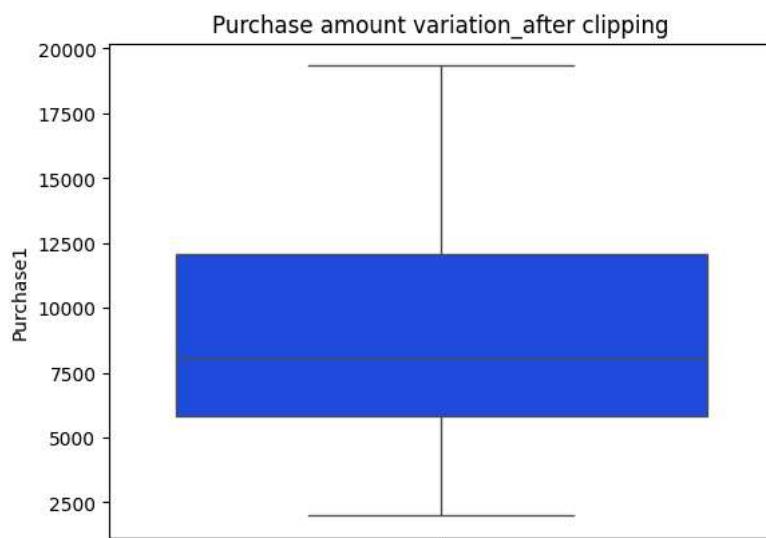
```
2 # df1 = df
```

```
3 df['Purchase1'] = np.clip(df['Purchase'], a_min = a, a_max= b)
```

```
4 sns.boxplot(df['Purchase1'])
```

```
5 plt.title('Purchase amount variation_after clipping')
```

```
6 plt.show()
```



```
1 #Total Sales till date is (after clipping)
2 df['Purchase1'].sum()
```

```
5091820225
```

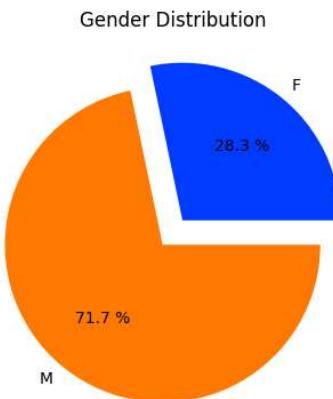
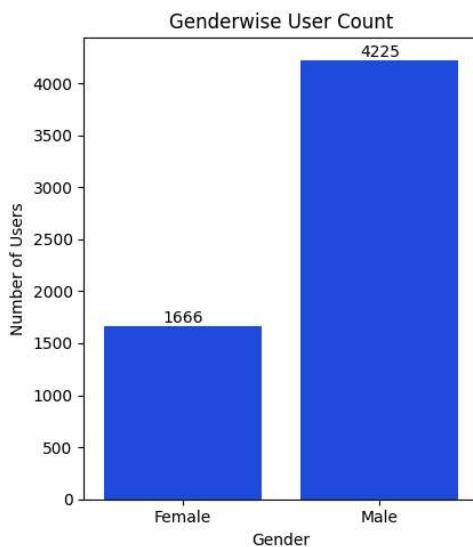
The difference between total purchase before and after clipping is 3992517(5095812742 - 5091820225\$)

```
1 5095812742 - 5091820225
```

```
3992517
```

## ▼ Graphical Analysis

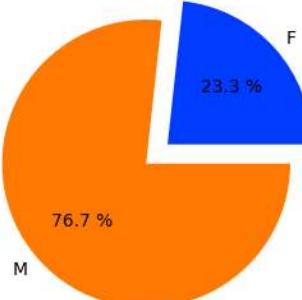
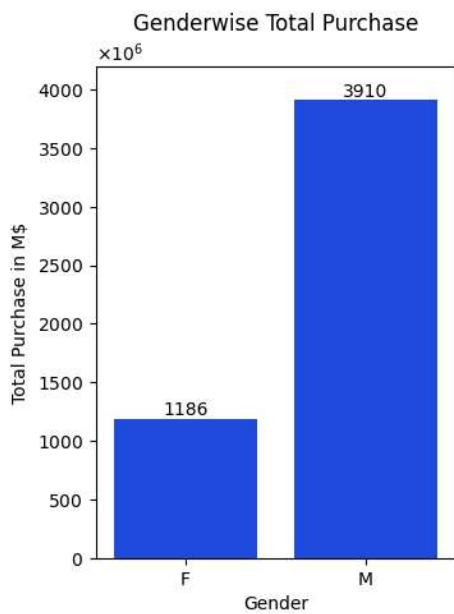
```
1 grp1 = pd.DataFrame(df.groupby(['Gender']).agg({'User_ID':'nunique'})).reset_index()
2
3 plt.figure(figsize=[8,5], num= 2)
4
5 plt.subplot(121)
6
7 ax = sns.barplot(grp1, x = 'Gender' , y = 'User_ID')
8 plt.title('Genderwise User Count', loc= 'center')
9 plt.ylabel('Number of Users')
10 plt.xticks([0, 1], ['Female', 'Male'])
11
12 for p in ax.patches:
13     ax.annotate(format(p.get_height(), '.0f'),
14                 (p.get_x() + p.get_width() / 2., p.get_height()),
15                 ha = 'center', va = 'center',
16                 xytext = (0, 5),
17                 textcoords = 'offset points')
18
19 plt.subplot(122, aspect = 'equal')
20 plt.pie(grp1['User_ID'], labels= grp1['Gender'], autopct= '%1.1f %%', explode= [0.1,0.1])
21 plt.title('Gender Distribution')
22 plt.tight_layout()
23 plt.show()
```



```

1 grp2 = pd.DataFrame(df.groupby(['Gender']).agg({'Purchase': 'sum'}).reset_index()
2 plt.figure(figsize=[8,5] , num = 2)
3
4 plt.subplot(121)
5 ax =sns.barplot(grp2, x = 'Gender' , y = 'Purchase')
6 plt.title('Genderwise Total Purchase', loc= 'center')
7 plt.ylabel('Total Purchase in M$')
8 ax.yaxis.set_major_formatter(ScalarFormatter(useMathText=True, useOffset=False))
9 ax.ticklabel_format(axis='y', style='sci', scilimits=(6,6))
10 for p in ax.patches:
11     ax.annotate(format(p.get_height()/10**6, '.0f'),
12                 (p.get_x() + p.get_width() / 2., p.get_height(),
13                  ha = 'center', va = 'center',
14                  xytext = (0, 5),
15                  textcoords = 'offset points',
16                  )
17 plt.ylim(0,4200*10**6)
18
19 plt.subplot(122, aspect = 'equal')
20
21 plt.pie(grp2['Purchase'] , labels= grp2['Gender'], autopct= '%1.1f %%', explode= [0.1,0.1])
22 plt.show()

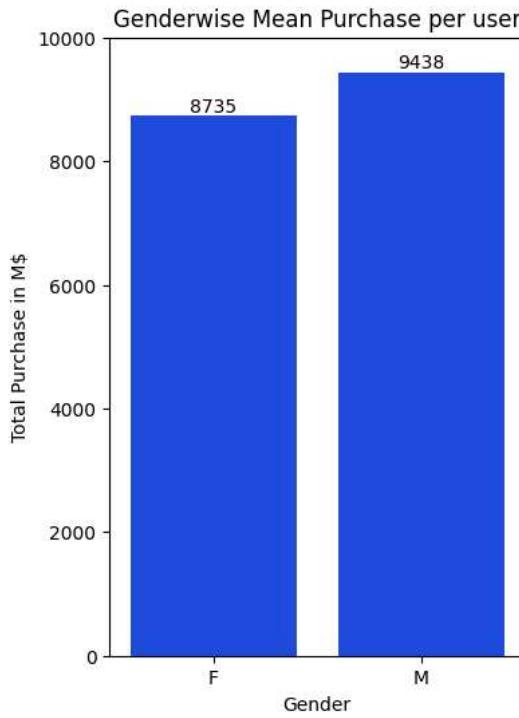
```



```

1 grp2_1 = pd.DataFrame(df.groupby(['Gender']).agg({'Purchase':'mean'})).reset_index()
2 plt.figure(figsize=[4,6])
3 ax =sns.barplot(grp2_1, x = 'Gender' , y = 'Purchase')
4 plt.title('Genderwise Mean Purchase per user', loc= 'center')
5 plt.ylabel('Total Purchase in M$')
6 for p in ax.patches:
7     ax.annotate(format(p.get_height(), '.0f'),
8                 (p.get_x() + p.get_width() / 2., p.get_height()),
9                 ha = 'center', va = 'center',
10                xytext = (0, 5),
11                textcoords = 'offset points',
12                )
13 plt.ylim(0,10000)
14 plt.show()

```

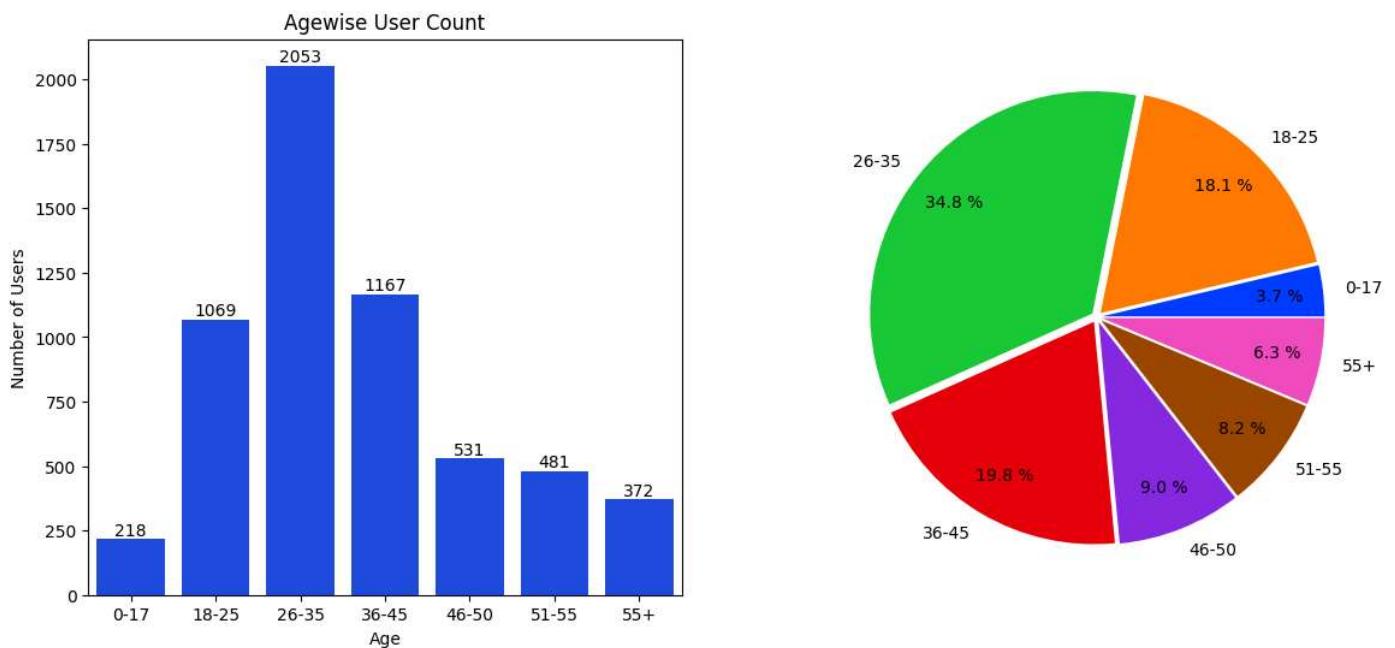


- ✓ The total purchase amount made by males significantly exceeds that of females. Additionally, the average purchase per user among males is higher than that among females.

```

1 grp3 = pd.DataFrame(df.groupby(['Age']).agg({'User_ID':'nunique'})).reset_index()
2 plt.figure(figsize=[14,6], num = 2)
3
4 plt.subplot(121)
5 ax = sns.barplot(grp3, x = 'Age' , y = 'User_ID')
6 plt.title('Agewise User Count', loc= 'center')
7 plt.ylabel('Number of Users')
8 for p in ax.patches:
9     ax.annotate(format(p.get_height(), '.0f'),
10                 (p.get_x() + p.get_width() / 2., p.get_height()),
11                 ha = 'center', va = 'center',
12                 xytext = (0, 5),
13                 textcoords = 'offset points')
14
15 plt.subplot(122, aspect = 'equal')
16 plt.pie(grp3['User_ID'] , labels= grp3['Age'] , explode = [0.025] * len(grp3['Age']), autopct= '%1.1f %%' , pctdistance= 0.8 )
17 plt.show()

```

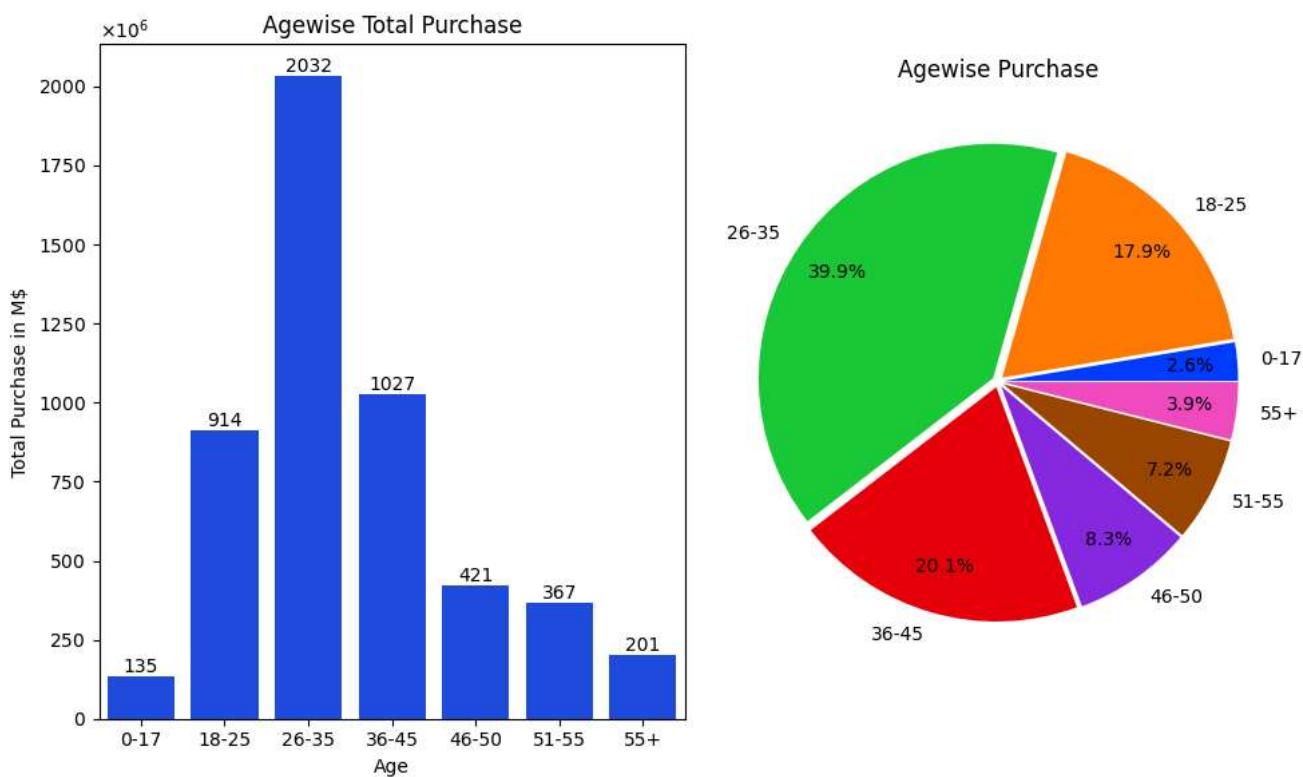


- Maximum users are from Age range 26-35 followed by 36-45 with 2053 and 1167 users respectively.

```

1 grp4 = pd.DataFrame(df.groupby(['Age']).agg({'Purchase':'sum'}).reset_index())
2 plt.figure(figsize=[10,6], num = 2)
3
4 plt.subplot(121)
5 ax = sns.barplot(grp4, x = 'Age' , y = 'Purchase')
6 plt.title('Agewise Total Purchase', loc= 'center')
7 plt.ylabel('Total Purchase in M$')
8 ax.yaxis.set_major_formatter(ScalarFormatter(useMathText=True, useOffset=False))
9 ax.ticklabel_format(axis='y', style='sci', scilimits=(6,6))
10 for p in ax.patches:
11     ax.annotate(format(p.get_height()/10**6, '.0f'),
12                 (p.get_x() + p.get_width() / 2., p.get_height()),
13                 ha = 'center', va = 'center',
14                 xytext = (0, 5),
15                 textcoords = 'offset points')
16
17
18 plt.subplot(122, aspect = 'equal')
19 plt.pie(grp4['Purchase'], labels=grp4['Age'], explode = [0.025] * len(grp4), autopct= '%1.1f%%', pctdistance= 0.8)
20 plt.title('Agewise Purchase', loc= 'center')
21 plt.tight_layout()
22 plt.show()

```

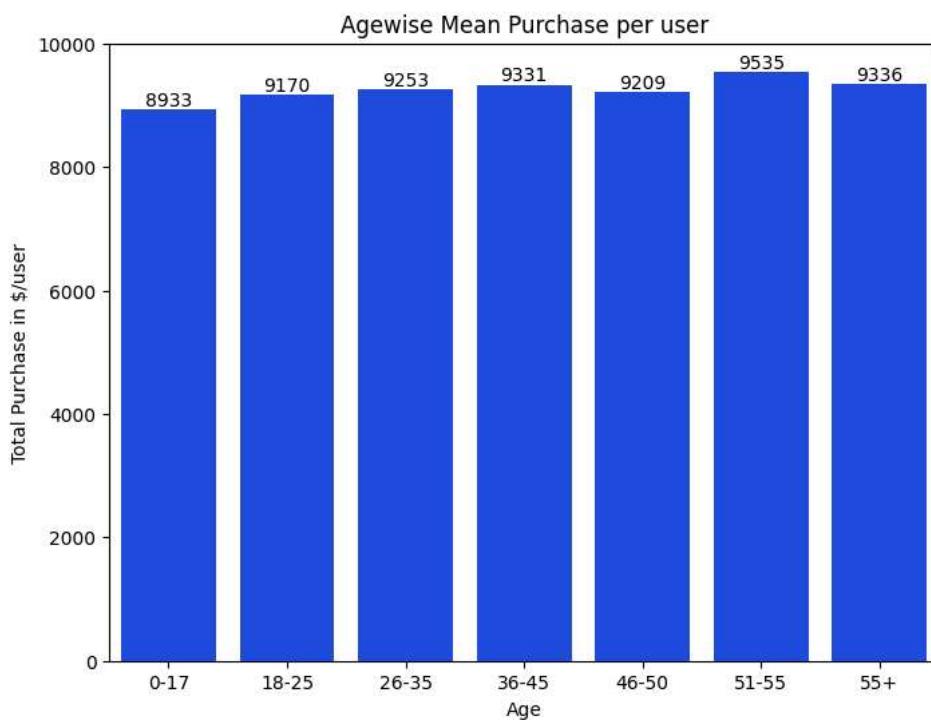


- ✓ The maximum purchase has been done by the Agegroup of 26-35, leading by 2032 M (39.9 (20.1%). Lowest is of 0-17 and 55+ with 135 Mand 201M

```

1 grp4_1 = pd.DataFrame(df.groupby(['Age']).agg({'Purchase':'mean'})).reset_index()
2 plt.figure(figsize=[8,6])
3 ax = sns.barplot(grp4_1, x = 'Age' , y = 'Purchase')
4 plt.title('Agewise Mean Purchase per user', loc= 'center')
5 plt.ylabel('Total Purchase in $/user')
6 for p in ax.patches:
7     ax.annotate(format(p.get_height(), '.0f'),
8                 (p.get_x() + p.get_width() / 2., p.get_height()),
9                 ha = 'center', va = 'center',
10                xytext = (0, 5),
11                textcoords = 'offset points')
12 plt.ylim(0,10000)
13 plt.show()

```

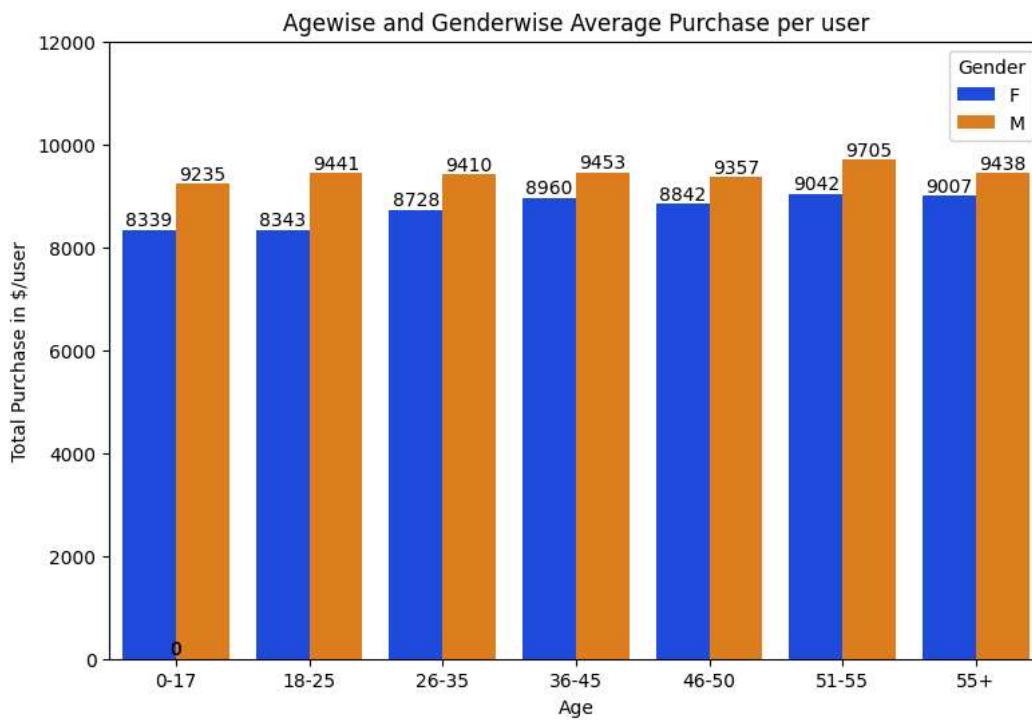


- ✓ The most spending user is age group 51-55 , followed by 36-45 with 9535 and 9331 respectively. The lowest spending age group is 0-17 and 18-25 with 8933 and 9170.

```

1 grp4_2 = pd.DataFrame(df.groupby(['Age', 'Gender']).agg({'Purchase':'mean'})).reset_index()
2 plt.figure(figsize=[9,6])
3 ax = sns.barplot(grp4_2, x = 'Age' , y = 'Purchase', hue= 'Gender')
4 plt.title('Agewise and Genderwise Average Purchase per user', loc= 'center')
5 plt.ylabel('Total Purchase in $/user')
6 for p in ax.patches:
7     ax.annotate(format(p.get_height(), '.0f'),
8                 (p.get_x() + p.get_width() / 2., p.get_height()),
9                 ha = 'center', va = 'center',
10                xytext = (0, 5),
11                textcoords = 'offset points')
12 plt.ylim(0,12000)
13 plt.show()

```

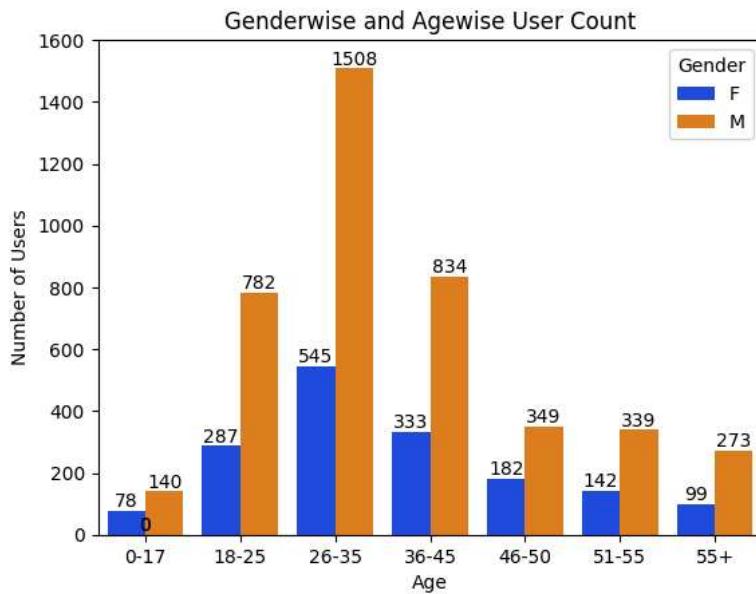


- ✓ The Average spendings of men is higher than women in all age ranges.

```

1 grp5 = pd.DataFrame(df.groupby(['Gender', 'Age']).agg({'User_ID':'nunique'}).reset_index())
2 ax = sns.barplot(grp5, x = 'Age' , y = 'User_ID', hue='Gender')
3 plt.title('Genderwise and Agewise User Count', loc= 'center')
4 plt.ylabel('Number of Users')
5 for p in ax.patches:
6     ax.annotate(format(p.get_height(), '.0f'),
7                 (p.get_x() + p.get_width() / 2., p.get_height()),
8                 ha = 'center', va = 'center',
9                 xytext = (0, 5),
10                textcoords = 'offset points')
11 plt.ylim(0,1600)
12 plt.show()

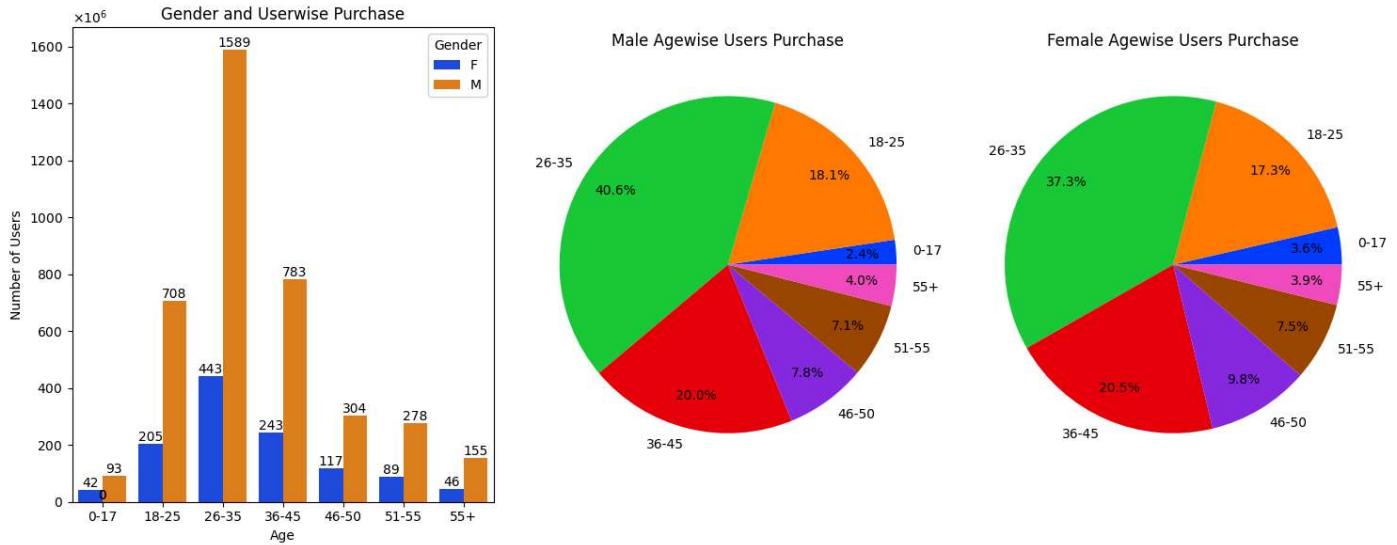
```



```

1 grp6 = pd.DataFrame(df.groupby(['Gender', 'Age']).agg({'Purchase':'sum'})).reset_index()
2 plt.figure(num= 3, figsize= (15,6))
3
4 plt.subplot(131)
5 ax = sns.barplot(grp6 , x = 'Age' , y = 'Purchase', hue='Gender')
6 plt.title('Gender and Userwise Purchase', loc= 'center')
7 plt.ylabel('Number of Users')
8 ax.yaxis.set_major_formatter(ScalarFormatter(useMathText=True, useOffset=False))
9 ax.ticklabel_format(axis='y', style='sci', scilimits=(6,6))
10 for p in ax.patches:
11     ax.annotate(format(p.get_height()/10**6, '.0f'),
12                 (p.get_x() + p.get_width() / 2., p.get_height() + 25),
13                 ha = 'center', va = 'center',
14                 xytext = (0, 5),
15                 textcoords = 'offset points',
16                 rotation = 0)
17
18 plt.subplot(132, aspect = 'equal')
19 plt.pie(grp6[grp6['Gender'] == 'M']['Purchase'], labels= (grp6[grp6['Gender']=='M']['Age']), autopct= '%1.1f%%', pctdistance= 0.8)
20 plt.title('Male Agewise Users Purchase')
21
22 plt.subplot(133, aspect = 'equal')
23 plt.pie(grp6[grp6['Gender'] == 'F']['Purchase'], labels= (grp6[grp6['Gender']=='F']['Age']), autopct= '%1.1f%%', pctdistance= 0.8)
24 plt.title('Female Agewise Users Purchase')
25
26 plt.tight_layout()
27 plt.show()

```

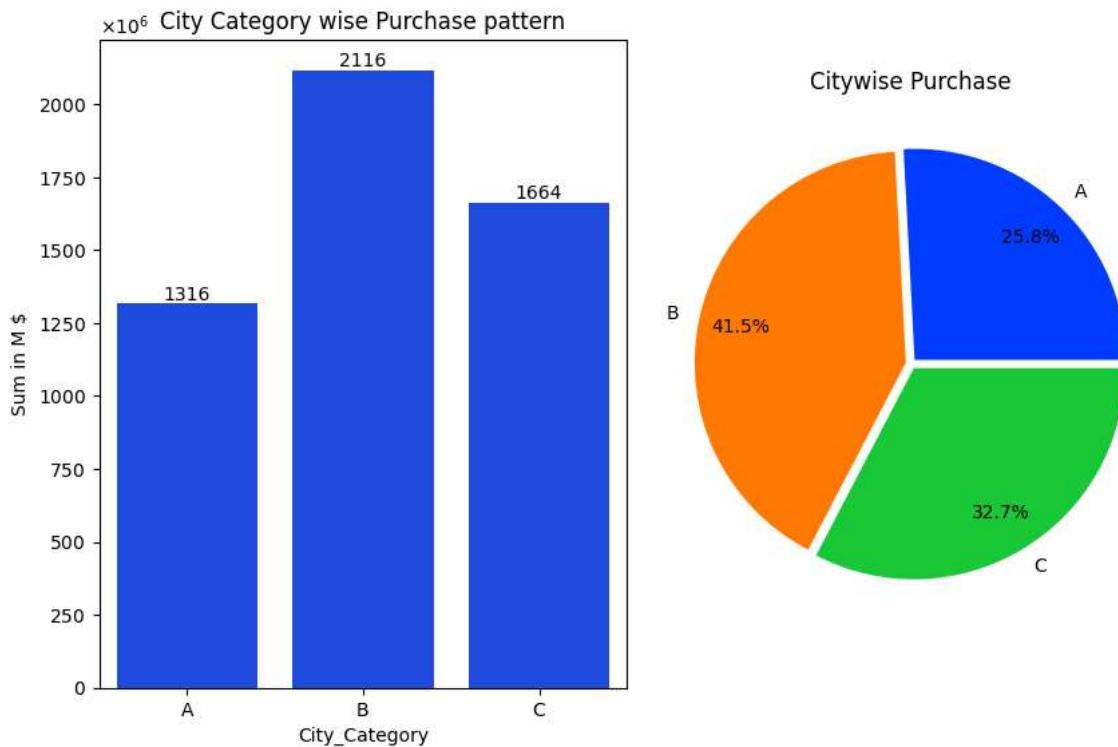


- There are 4225 Males and 1666 Female customers with purchase amounting 3910 Million and 1186 Million respectively.
- Most of the customers are in the age range 18-45.

```

1 grp7 = pd.DataFrame(df.groupby(['City_Category']).agg({'Purchase':'sum'})).reset_index()
2 plt.figure(num= 2, figsize= (9,6))
3
4 plt.subplot(121)
5 ax = sns.barplot(grp7 , x = 'City_Category' , y = 'Purchase')
6 plt.title('City Category wise Purchase pattern', loc= 'center')
7 plt.ylabel('Sum in M $')
8 ax.yaxis.set_major_formatter(ScalarFormatter(useMathText=True, useOffset=False))
9 ax.ticklabel_format(axis='y', style='sci', scilimits=(6,6))
10 for p in ax.patches:
11     ax.annotate(format(p.get_height()/10**6, '.0f'),
12                 (p.get_x() + p.get_width() / 2., p.get_height()+25),
13                 ha = 'center', va = 'center',
14                 xytext = (0, 5),
15                 textcoords = 'offset points',
16                 rotation = 0)
17
18 plt.subplot(122, aspect = 'equal')
19 plt.pie(grp7['Purchase'], labels=grp7['City_Category'], explode = [0.025] * len(grp7), autopct= '%1.1f%%', pctdistance= 0.8)
20 plt.title('Citywise Purchase', loc= 'center')
21 plt.tight_layout()
22
23 plt.show()

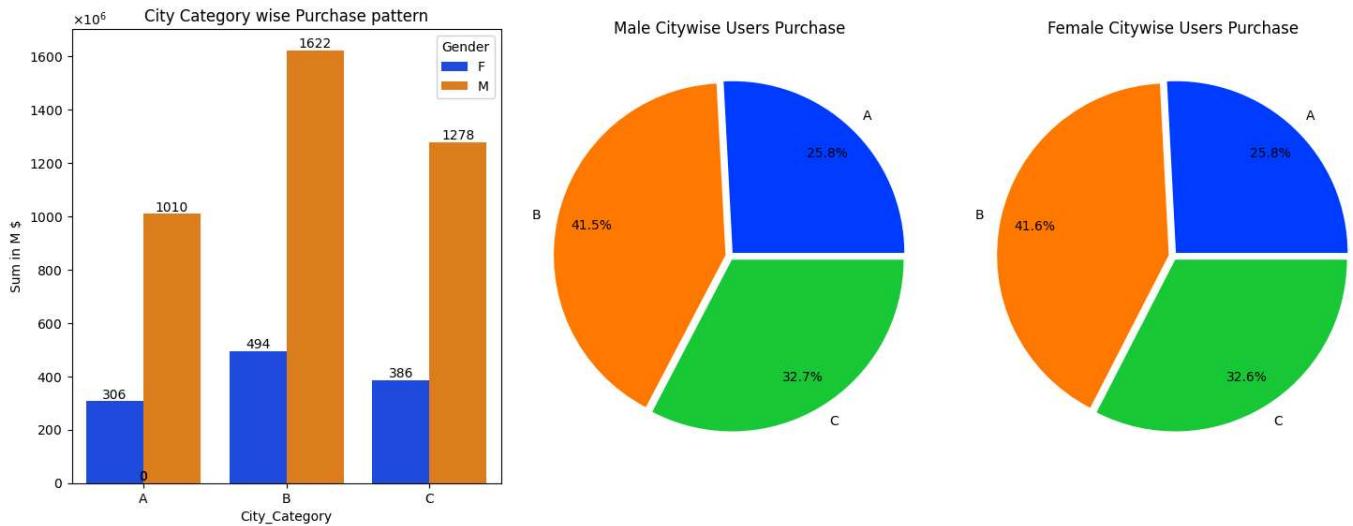
```



```

1 grp7_1 = pd.DataFrame(df.groupby(['City_Category', 'Gender']).agg({'Purchase':'sum'}).reset_index())
2 plt.figure(num = 3 , figsize = (15,6))
3
4 plt.subplot(131)
5 ax = sns.barplot(grp7_1 , x = 'City_Category' , y = 'Purchase', hue = 'Gender')
6 plt.title('City Category wise Purchase pattern', loc= 'center')
7 plt.ylabel('Sum in M $')
8 ax.yaxis.set_major_formatter(ScalarFormatter(useMathText=True, useOffset=False))
9 ax.ticklabel_format(axis='y', style='sci', scilimits=(6,6))
10 for p in ax.patches:
11     ax.annotate(format(p.get_height()/10**6, '.0f'),
12                 (p.get_x() + p.get_width() / 2., p.get_height()+25),
13                 ha = 'center', va = 'center',
14                 xytext = (0, 5),
15                 textcoords = 'offset points',
16                 rotation = 0)
17
18 plt.subplot(132, aspect = 'equal')
19 plt.pie(grp7_1[grp7_1['Gender'] == 'M']['Purchase'], labels= (grp7_1[grp7_1['Gender']=='M']['City_Category']),
20          autopct= '%.1f%%', pctdistance= 0.8,explode= [0.025]*(len(grp7_1[grp7_1['Gender'] == 'F'])))
21 plt.title('Male Citywise Users Purchase')
22
23 plt.subplot(133, aspect = 'equal')
24 plt.pie(grp7_1[grp7_1['Gender'] == 'F']['Purchase'], labels= (grp7_1[grp7_1['Gender']=='F']['City_Category']),
25          autopct= '%.1f%%', pctdistance= 0.8, explode= [0.025]*(len(grp7_1[grp7_1['Gender'] == 'F'])))
26 plt.title('Female Citywise Users Purchase')
27
28
29 plt.tight_layout()
30

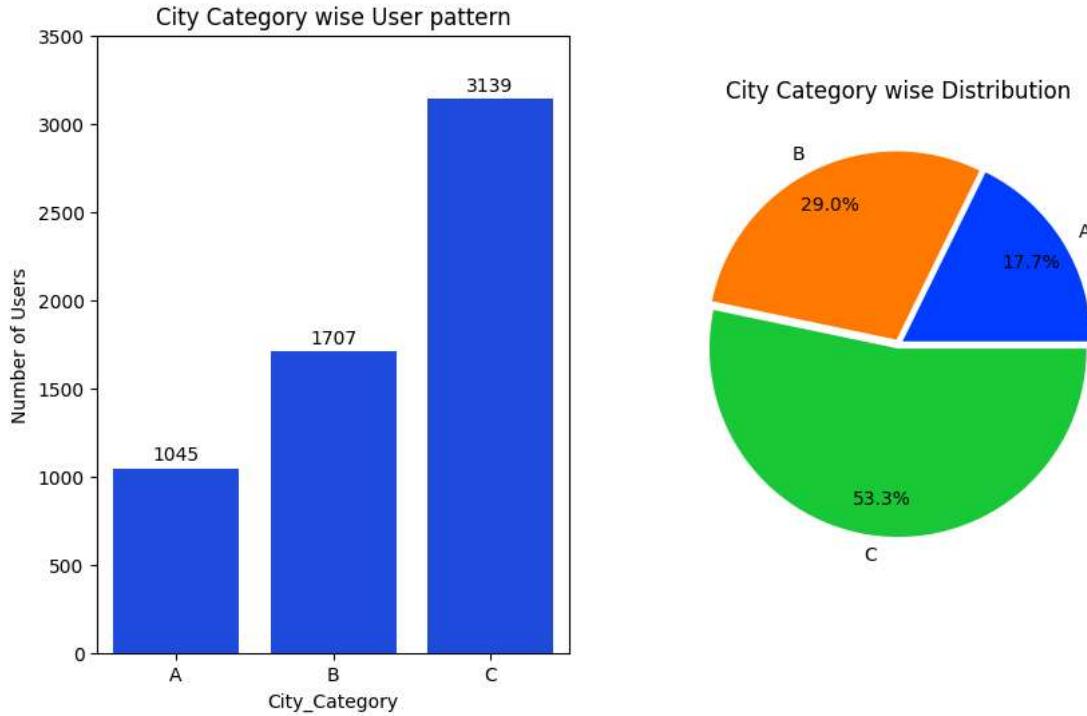
```



```

1 grp8 = pd.DataFrame(df.groupby(['City_Category']).agg({'User_ID':'nunique'})).reset_index()
2 plt.figure(figsize = (10,6), num = 2)
3
4 plt.subplot(121)
5 ax = sns.barplot(grp8 , x = 'City_Category' , y = 'User_ID')
6 plt.title('City Category wise User pattern', loc= 'center')
7 plt.ylabel('Number of Users')
8 # ax.yaxis.set_major_formatter(ScalarFormatter(useMathText=True, useOffset=False))
9 # ax.ticklabel_format(axis='y', style='sci', scilimits=(9,9))
10 for p in ax.patches:
11     ax.annotate(format(p.get_height(), '.0f'),
12                 (p.get_x() + p.get_width() / 2., p.get_height() + 25),
13                 ha = 'center', va = 'center',
14                 xytext = (0, 5),
15                 textcoords = 'offset points',
16                 rotation = 0)
17 plt.ylim(0,3500)
18
19 plt.subplot(122, aspect = 'equal')
20
21 plt.pie(grp8['User_ID'], labels= grp8['City_Category'], autopct= '%1.1f%%', pctdistance= 0.8,
22          explode= [0.025]*len(grp8['User_ID']))
23 plt.title('City Category wise Distribution')
24
25 plt.show()

```



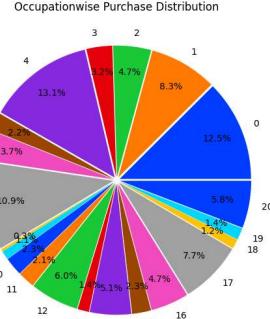
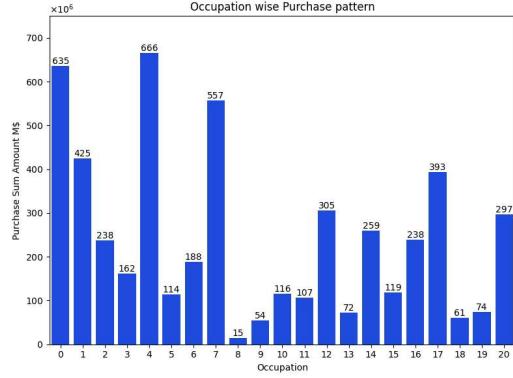
- Majority of users are in the city\_category C followed by B.

The Category B city is leading in sales and Genderwise as well.

```

1 grp9 = pd.DataFrame(df.groupby(['Occupation']).agg({'Purchase' : 'sum'})).reset_index()
2 plt.figure(figsize=(15,6),num= 2)
3
4 plt.subplot(121)
5 ax = sns.barplot(grp9, x = 'Occupation', y = 'Purchase')
6 plt.title('Occupation wise Purchase pattern', loc= 'center')
7 plt.ylabel('Purchase Sum Amount M$')
8 ax.yaxis.set_major_formatter(ScalarFormatter(useMathText=True, useOffset=False))
9 ax.ticklabel_format(axis='y', style='sci', scilimits=(6,6))
10 for p in ax.patches:
11     ax.annotate(format(p.get_height()/10**6, '.0f'),
12                 (p.get_x() + p.get_width() / 2, p.get_height()),
13                 ha = 'center', va = 'center',
14                 xytext = (0, 5),
15                 textcoords = 'offset points',
16                 rotation = 0)
17 plt.ylim(0,750*10**6)
18
19 plt.subplot(122)
20 plt.pie(grp9['Purchase'], labels = grp9['Occupation'],
21         autopct = "%1.1f%%", pctdistance= 0.8,
22         explode= [0.025]*len(grp9['Occupation']))
23 plt.title('Occupationwise Purchase Distribution')
24 plt.tight_layout()
25 plt.show()

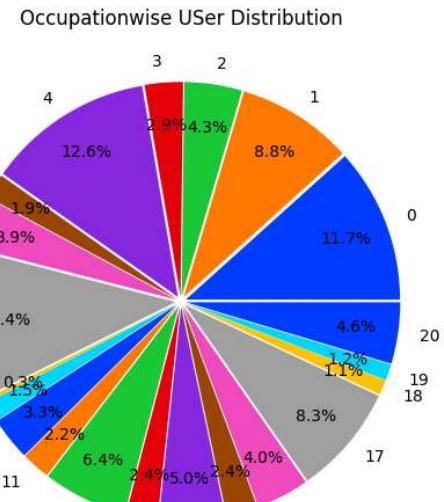
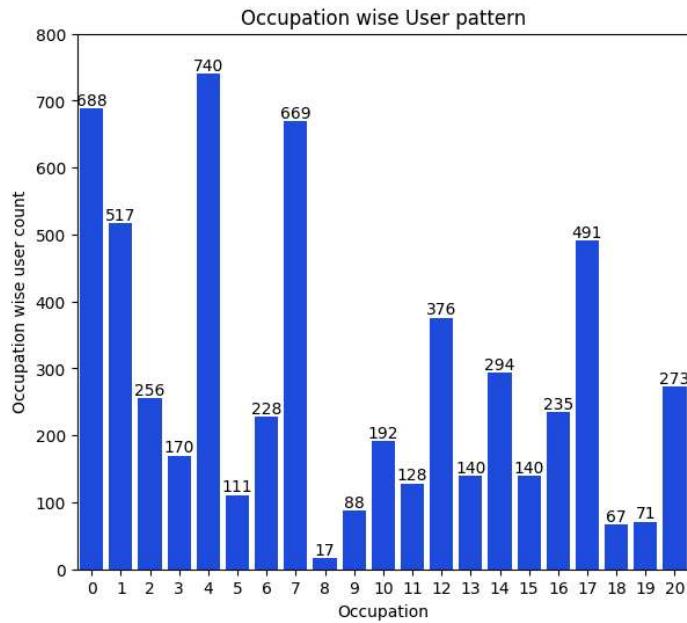
```



```

1 grp10 = pd.DataFrame(df.groupby(['Occupation']).agg({'User_ID' : 'nunique'})).reset_index()
2 plt.figure(figsize=(15,6), num =2)
3
4 plt.subplot(121)
5 ax = sns.barplot(grp10, x = 'Occupation', y = 'User_ID')
6 plt.title('Occupation wise User pattern', loc= 'center')
7 plt.ylabel('Occupation wise user count')
8 for p in ax.patches:
9     ax.annotate(format(p.get_height(), '.0f'),
10                 (p.get_x() + p.get_width() / 2, p.get_height()),
11                 ha = 'center', va = 'center',
12                 xytext = (0, 5),
13                 textcoords = 'offset points',
14                 rotation = 0)
15 plt.ylim(0,800)
16
17 plt.subplot(122)
18
19 plt.pie(grp10['User_ID'], labels= grp10['Occupation'],
20         autopct = '%1.1f%%', pctdistance= .8,
21         explode=[0.025]*len(grp10['Occupation']))
22 plt.title('Occupationwise USeR Distribution')
23 plt.show()

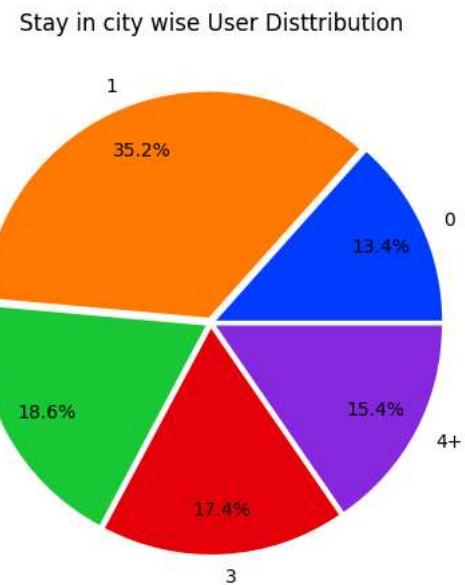
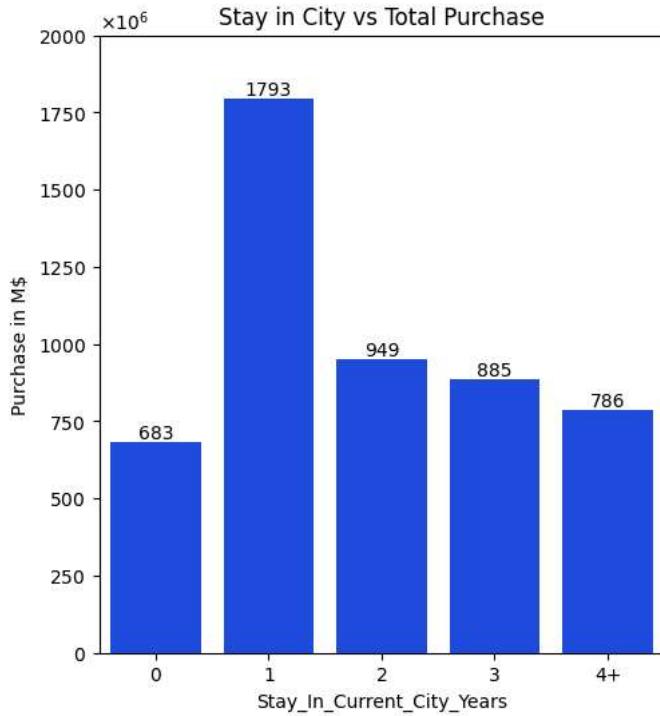
```



```

1 grp11 = pd.DataFrame(df.groupby('Stay_In_Current_City_Years').agg({'Purchase': 'sum'})).reset_index()
2 plt.figure(num = 2, figsize=(12,6))
3
4 plt.subplot(121)
5 ax = sns.barplot(grp11, x = 'Stay_In_Current_City_Years', y= 'Purchase')
6 ax.yaxis.set_major_formatter(ScalarFormatter(useMathText=True, useOffset=False))
7 ax.ticklabel_format(axis='y', style='sci', scilimits=(6,6))
8 for p in ax.patches:
9     ax.annotate(format(p.get_height()/10**6, '.0f'),
10                 (p.get_x() + p.get_width() / 2, p.get_height()),
11                 ha = 'center', va = 'center',
12                 xytext = (0, 5),
13                 textcoords = 'offset points',
14                 rotation = 0)
15 plt.ylim(0,2*10**9)
16 plt.ylabel('Purchase in M$')
17 plt.title('Stay in City vs Total Purchase')
18
19 plt.subplot(122)
20 plt.pie(grp11['Purchase'], labels = grp11['Stay_In_Current_City_Years'],
21         autopct = '%1.1f%%', pctdistance= 0.8,
22         explode=[0.025]*len(grp11['Stay_In_Current_City_Years']))
23 plt.title('Stay in city wise User Distribution')
24
25 plt.show()

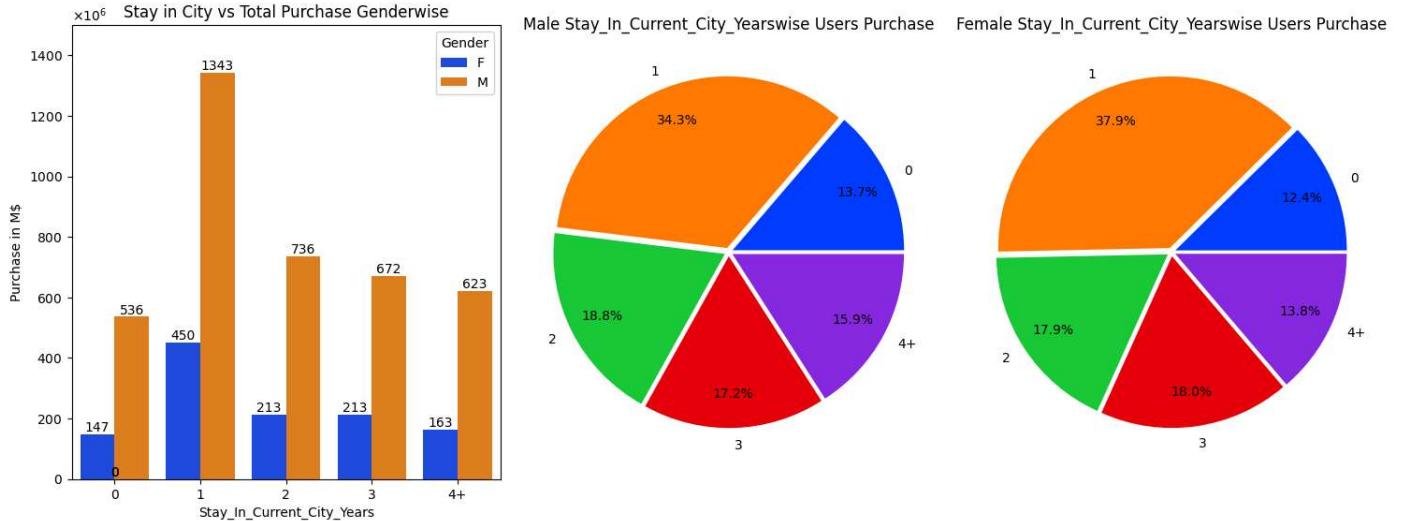
```



```

1 grp11_1 = pd.DataFrame(df.groupby(['Stay_In_Current_City_Years', 'Gender']).agg({'Purchase': 'sum'}).reset_index())
2 plt.figure(num=3, figsize=(15,6))
3
4 plt.subplot(131)
5 ax = sns.barplot(grp11_1, x = 'Stay_In_Current_City_Years', y= 'Purchase', hue = 'Gender')
6 ax.yaxis.set_major_formatter(ScalarFormatter(useMathText=True, useOffset=False))
7 ax.ticklabel_format(axis='y', style='sci', scilimits=(6,6))
8 for p in ax.patches:
9     ax.annotate(format(p.get_height()/10**6, '.0f'),
10                 (p.get_x() + p.get_width() / 2, p.get_height()),
11                 ha = 'center', va = 'center',
12                 xytext = (0, 5),
13                 textcoords = 'offset points',
14                 rotation = 0)
15 plt.ylim(0,1.5*10**9)
16 plt.ylabel('Purchase in M$')
17 plt.title('Stay in City vs Total Purchase Genderwise')
18
19 plt.subplot(132, aspect = 'equal')
20 plt.pie(grp11_1[grp11_1['Gender'] == 'M']['Purchase'], labels= (grp11_1[grp11_1['Gender']=='M']['Stay_In_Current_City_Years']),
21         autopct= '%1.1f%%', pctdistance= 0.8,explode= [0.025]*(len(grp11_1[grp11_1['Gender']] == 'F')))
22 plt.title('Male Stay_In_Current_City_Yearswise Users Purchase')
23
24 plt.subplot(133, aspect = 'equal')
25 plt.pie(grp11_1[grp11_1['Gender'] == 'F']['Purchase'], labels= (grp11_1[grp11_1['Gender']=='F']['Stay_In_Current_City_Years']),
26         autopct= '%1.1f%%', pctdistance= 0.8, explode= [0.025]*(len(grp11_1[grp11_1['Gender']] == 'F')))
27 plt.title('Female Stay_In_Current_City_Yearswise Users Purchase')
28
29
30
31 plt.tight_layout()
32
33
34 plt.show()

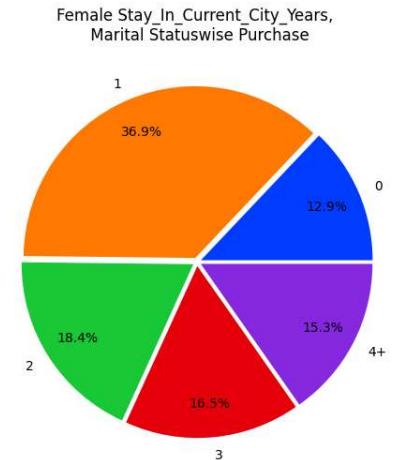
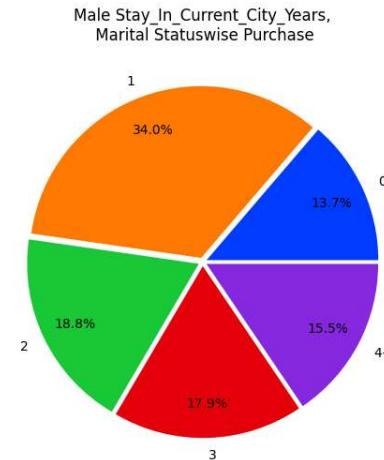
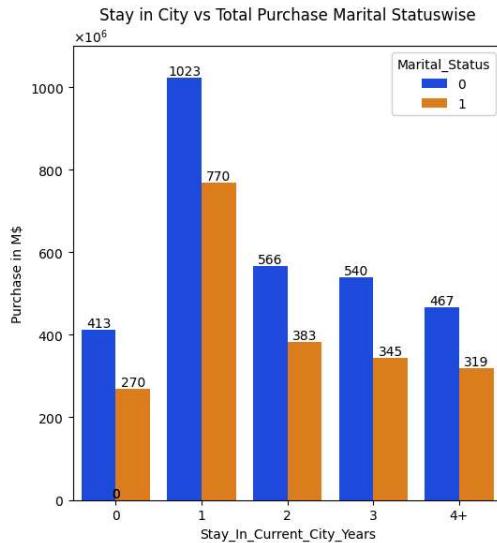
```



```

1 grp11_2 = pd.DataFrame(df.groupby(['Stay_In_Current_City_Years','Marital_Status']).agg({'Purchase': 'sum'}).reset_index())
2 plt.figure(num=3, figsize=(15,6))
3
4 plt.subplot(131)
5 ax = sns.barplot(grp11_2, x = 'Stay_In_Current_City_Years', y= 'Purchase', hue = 'Marital_Status')
6 ax.yaxis.set_major_formatter(ScalarFormatter(useMathText=True, useOffset=False))
7 ax.ticklabel_format(axis='y', style='sci', scilimits=(6,6))
8 for p in ax.patches:
9     ax.annotate(format(p.get_height()/10**6, '.0f'),
10                 (p.get_x() + p.get_width() / 2, p.get_height()),
11                 ha = 'center', va = 'center',
12                 xytext = (0, 5),
13                 textcoords = 'offset points',
14                 rotation = 0)
15 plt.ylim(0,1.1*10**9)
16 plt.ylabel('Purchase in M$')
17 plt.title('Stay in City vs Total Purchase Marital Statuswise')
18
19 plt.subplot(132, aspect = 'equal')
20 plt.pie(grp11_2[grp11_2['Marital_Status'] == 0]['Purchase'], labels= (grp11_2[grp11_2['Marital_Status']== 0]['Stay_In_Current_City_Years']
21             autopct= '%1.1f%%', pctdistance= 0.8,
22             explode= [0.025]*(len(grp11_2[grp11_2['Marital_Status'] == 0])))
23 plt.title('Male Stay_In_Current_City_Years, \n Marital Statuswise Purchase')
24
25 plt.subplot(133, aspect = 'equal')
26 plt.pie(grp11_2[grp11_2['Marital_Status'] == 1]['Purchase'], labels= (grp11_2[grp11_2['Marital_Status']==1]['Stay_In_Current_City_Years']
27             autopct= '%1.1f%%', pctdistance= 0.8, explode= [0.025]*(len(grp11_2[grp11_2['Marital_Status'] == 1])))
28 plt.title('Female Stay_In_Current_City_Years, \n Marital Statuswise Purchase')
29
30
31
32 plt.tight_layout()
33 plt.show()

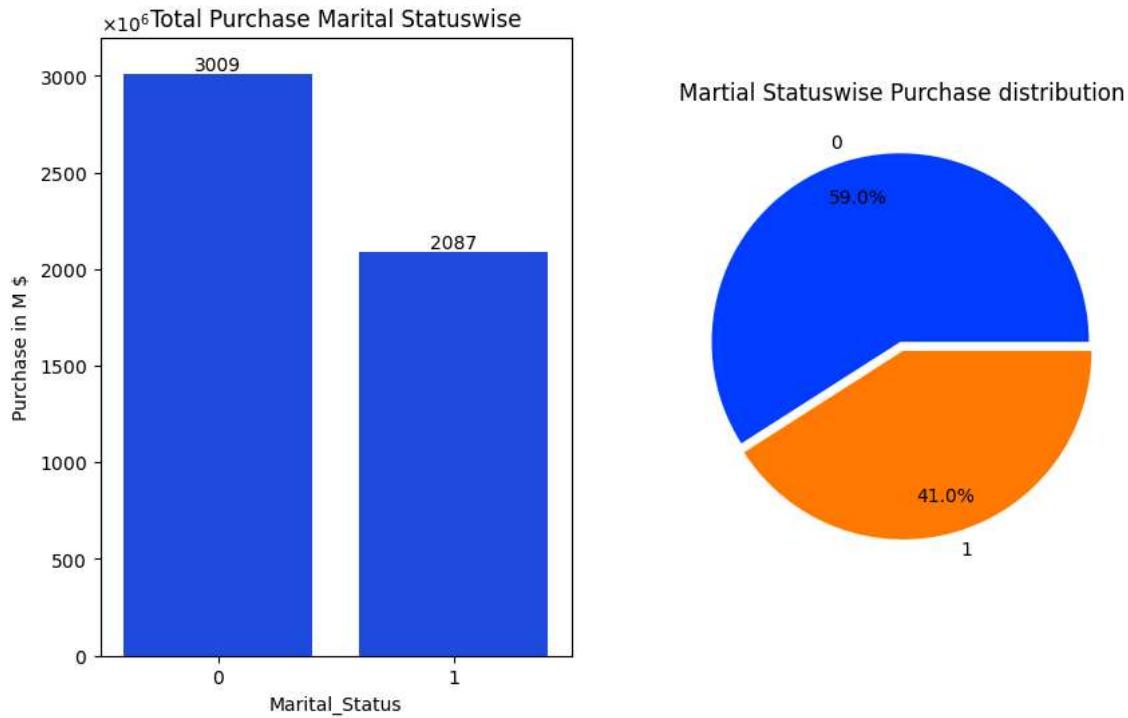
```



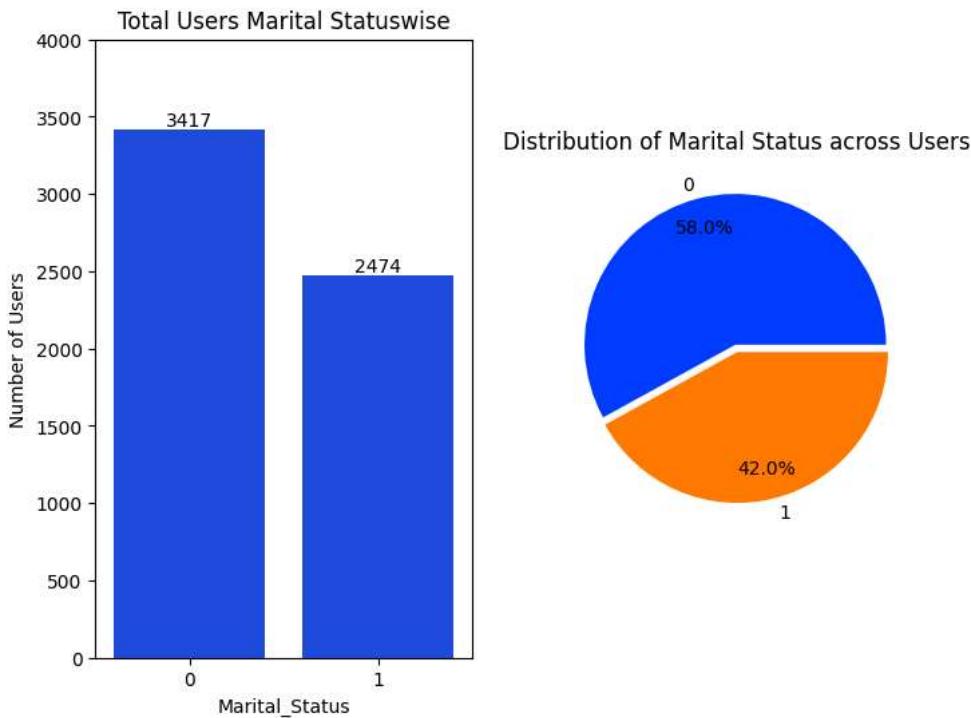
```

1 grp12 = pd.DataFrame(df.groupby(['Marital_Status']).agg({'Purchase': 'sum'}).reset_index())
2
3 plt.figure(figsize = (10,6), num= 2)
4 plt.subplot(121)
5 ax = sns.barplot(grp12, x = 'Marital_Status', y = 'Purchase')
6 ax.yaxis.set_major_formatter(ScalarFormatter(useMathText=True, useOffset=False))
7 ax.ticklabel_format(axis='y', style='sci', scilimits=(6,6))
8 for p in ax.patches:
9     ax.annotate(format(p.get_height()/10**6, '.0f'),
10                 (p.get_x() + p.get_width() / 2, p.get_height()),
11                 ha = 'center', va = 'center',
12                 xytext = (0, 5),
13                 textcoords = 'offset points',
14                 rotation = 0)
15 plt.ylim(0,3.2*10**9)
16 plt.ylabel('Purchase in M $')
17 plt.title('Total Purchase Marital Statuswise')
18
19 plt.subplot(122, aspect = 'equal')
20 plt.pie(grp12['Purchase'], labels=grp12['Marital_Status'],
21         autopct= '%1.1f%%', pctdistance= 0.8,
22         explode=[0.025]*len(grp12['Marital_Status']))
23 plt.title('Marital Statuswise Purchase distribution')
24 plt.show()

```



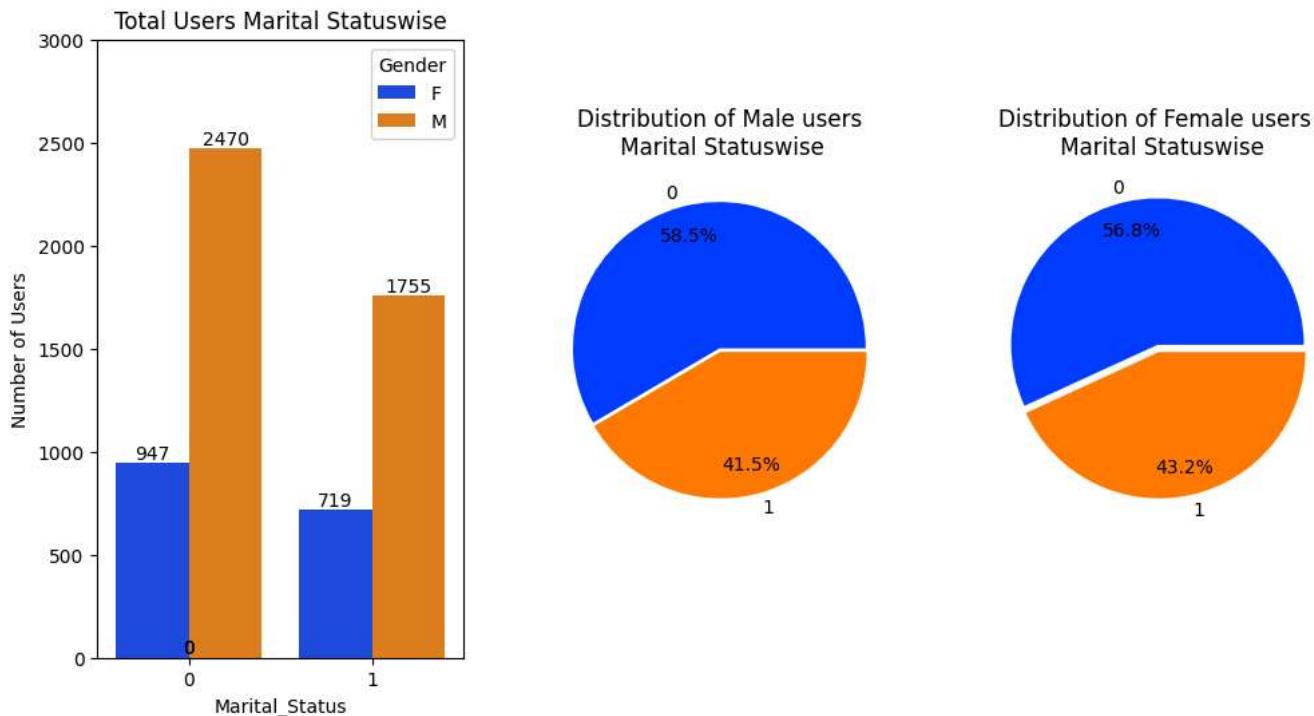
```
1 grp12_0 = pd.DataFrame(df.groupby(['Marital_Status']).agg({'User_ID': 'nunique'})).reset_index()
2 plt.figure(figsize = (8,6), num = 2)
3 plt.subplot(121)
4 ax = sns.barplot(grp12_0, x = 'Marital_Status', y = 'User_ID')
5 # ax.yaxis.set_major_formatter(ScalarFormatter(useMathText=True, useOffset=False))
6 # ax.ticklabel_format(axis='y', style='sci', scilimits=(6,6))
7 for p in ax.patches:
8     ax.annotate(format(p.get_height(), '.0f'),
9                 (p.get_x() + p.get_width() / 2, p.get_height()),
10                ha = 'center', va = 'center',
11                xytext = (0, 5),
12                textcoords = 'offset points',
13                rotation = 0)
14 plt.ylim(0,4000)
15 plt.ylabel('Number of Users')
16 plt.title('Total Users Marital Statuswise')
17
18 plt.subplot(122)
19 plt.pie(grp12_0['User_ID'], labels=grp12_0['Marital_Status'],
20         autopct='%.1f%%', pctdistance= 0.8,
21         explode= [0.025]*len(grp12_0['Marital_Status']))
22 plt.title('Distribution of Marital Status across Users', loc= 'center')
23 plt.show()
```



```

1 grp12_1 = pd.DataFrame(df.groupby(['Marital_Status','Gender']).agg({'User_ID': 'nunique'}).reset_index())
2 plt.figure(figsize = (12,6), num = 3)
3
4 plt.subplot(131)
5 ax = sns.barplot(grp12_1, x = 'Marital_Status', y = 'User_ID', hue = 'Gender')
6 # ax.yaxis.set_major_formatter(ScalarFormatter(useMathText=True, useOffset=False))
7 # ax.ticklabel_format(axis='y', style='sci', scilimits=(6,6))
8 for p in ax.patches:
9     ax.annotate(format(p.get_height(), '.0f'),
10                 (p.get_x() + p.get_width() / 2, p.get_height()),
11                 ha = 'center', va = 'center',
12                 xytext = (0, 5),
13                 textcoords = 'offset points',
14                 rotation = 0)
15 plt.ylim(0,3000)
16 plt.ylabel('Number of Users')
17 plt.title('Total Users Marital Statuswise')
18
19 plt.subplot(132, aspect = 'equal')
20 plt.pie(grp12_1[grp12_1['Gender']=='M']['User_ID'],labels= grp12_1[grp12_1['Gender']=='M']['Marital_Status'],
21         autopct = '%1.1f%%', pctdistance= 0.8,
22         explode = [0.025]*len(grp12_1[grp12_1['Gender']=='M']['Marital_Status']))
23 plt.title('Distribution of Male users \n Marital Statuswise')
24
25 plt.subplot(133)
26 plt.pie(grp12_1[grp12_1['Gender']=='F']['User_ID'],labels= grp12_1[grp12_1['Gender']=='F']['Marital_Status'],
27         autopct = '%1.1f%%', pctdistance= 0.8,
28         explode = [0.025]*len(grp12_1[grp12_1['Gender']=='F']['Marital_Status']))
29 plt.title('Distribution of Female users \n Marital Statuswise')
30 plt.show()

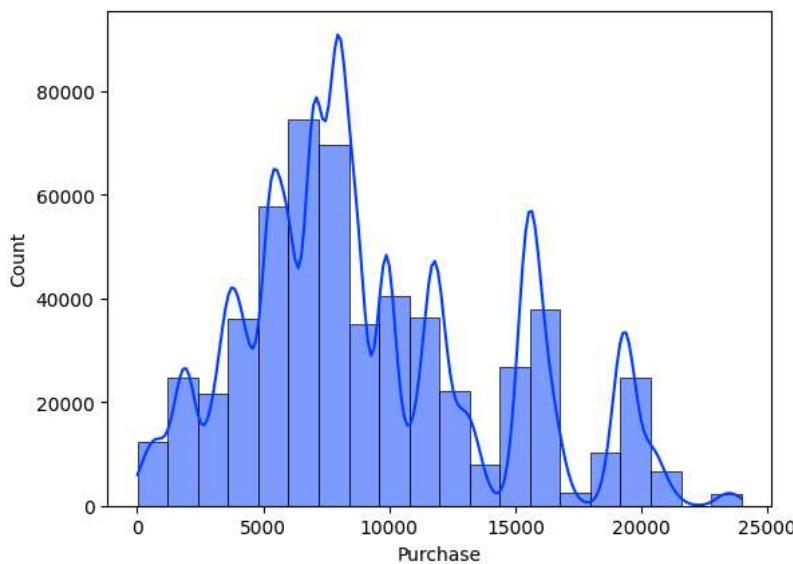
```



```

1 sns.histplot(df, x= 'Purchase', kde= True , bins = 20, linewidth = 0.5)
2 plt.show()

```

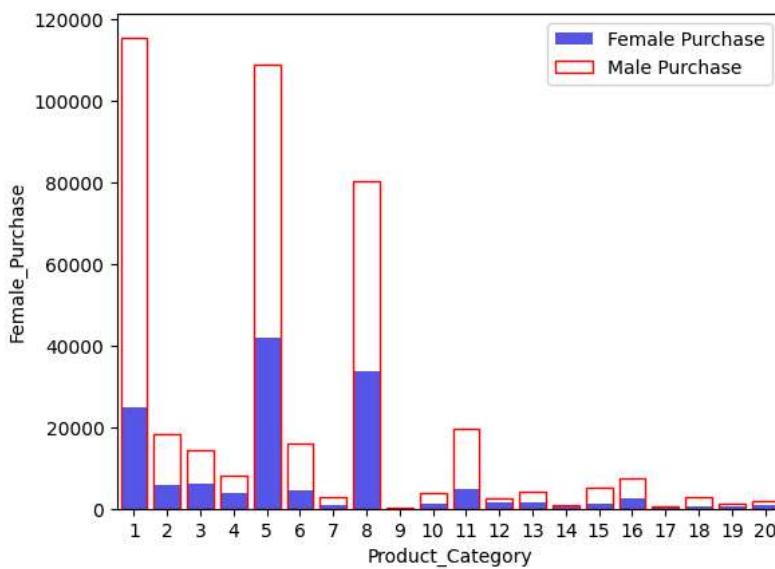


```

1 PC_Genderwise = df[df['Gender'] == 'F'].groupby(['Product_Category']).agg({'Purchase' : 'count'}).reset_index()
2 PC_for_Male = df[df['Gender'] == 'M'].groupby(['Product_Category']).agg({'Purchase':'count'}).reset_index()
3 PC_Genderwise.columns = ['Product_Category', 'Female_Purchase']
4 PC_Genderwise['Male_Purchase'] = PC_for_Male['Purchase'].astype('int')
5 PC_Genderwise['Female_Purchase'] = PC_Genderwise['Female_Purchase'].astype('int')

1 sns.barplot(data=PC_Genderwise, x='Product_Category', y='Female_Purchase', color= 'blue', alpha = 0.75, label='Female Purchase')
2 sns.barplot(data=PC_Genderwise, x='Product_Category', y='Male_Purchase', color = 'red', linewidth=1, edgecolor='red', facecolor='none', 1
3 plt.legend()
4 plt.show()

```



- Most Favoured Categories of Male are 1,5,8,11 and for females are 5,8,1,3 in decending orders.

Q.4 Genderwise Spending habit on Black Friday

```

1 df_male = df[df['Gender'] == 'M'].groupby(['User_ID']).agg({'Purchase' : 'sum'}).reset_index()
2 df_female = df[df['Gender'] == 'F'].groupby(['User_ID']).agg({'Purchase' : 'sum'}).reset_index()
3 [df_male.shape, df_female.shape]

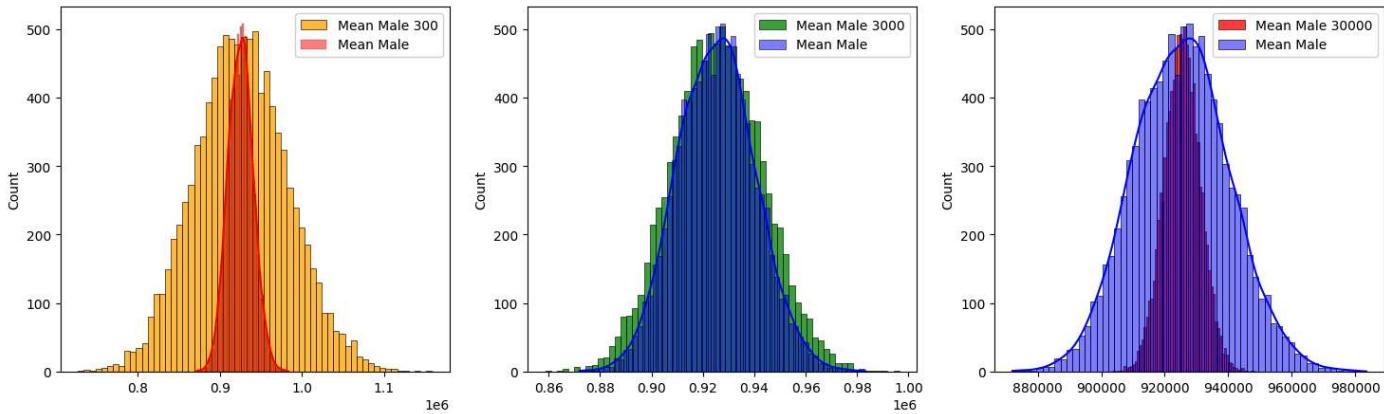
[(4225, 2), (1666, 2)]

```

```

1 mean_male = [np.mean(df_male['Purchase'].sample(len(df_male['User_ID']), replace = True)) for i in range (10000)]
2 mean_male_300 = [np.mean(df_male['Purchase'].sample(300, replace = True)) for i in range (10000)]
3 mean_male_3000 = [np.mean(df_male['Purchase'].sample(3000, replace = True)) for i in range (10000)]
4 mean_male_30000 = [np.mean(df_male['Purchase'].sample(30000, replace = True)) for i in range (10000)]
5 plt.figure(num=3, figsize=(18,5))
6
7 plt.subplot(131)
8 sns.histplot(data=mean_male_300, color='orange', label='Mean Male 300')
9 sns.histplot(data=mean_male,kde = True, color='red', label='Mean Male')
10 # plt.xlim(750000,1250000)
11 plt.legend()
12
13 plt.subplot(132)
14 sns.histplot(data=mean_male_3000, color='green', label='Mean Male 3000')
15 sns.histplot(data=mean_male,kde = True, color='blue', label='Mean Male')
16 # plt.xlim(750000,1250000)
17 plt.legend()
18
19 plt.subplot(133)
20 sns.histplot(data=mean_male_30000, color='red', label='Mean Male 30000')
21 sns.histplot(data=mean_male,kde = True, color='blue', label='Mean Male')
22 # plt.xlim(750000,1250000)
23
24 plt.legend()
25 plt.show()

```



```

1 mean_female = [np.mean(df_female['Purchase'].sample(len(df_female['User_ID']), replace= True)) for i in range (10000)]
2 mean_female_300 = [np.mean(df_female['Purchase'].sample(300, replace = True)) for i in range (10000)]
3 mean_female_3000 = [np.mean(df_female['Purchase'].sample(3000, replace = True)) for i in range (10000)]
4 mean_female_30000 = [np.mean(df_female['Purchase'].sample(30000, replace = True)) for i in range (10000)]

```

```

1 CI_male_purchase = np.percentile(mean_male, [5, 95])
2 CI_male_purchase = [int(np.round(value)) for value in CI_male_purchase]
3 CI_male_purchase_width = CI_male_purchase[1] - CI_male_purchase[0]
4 [CI_male_purchase , CI_male_purchase_width]

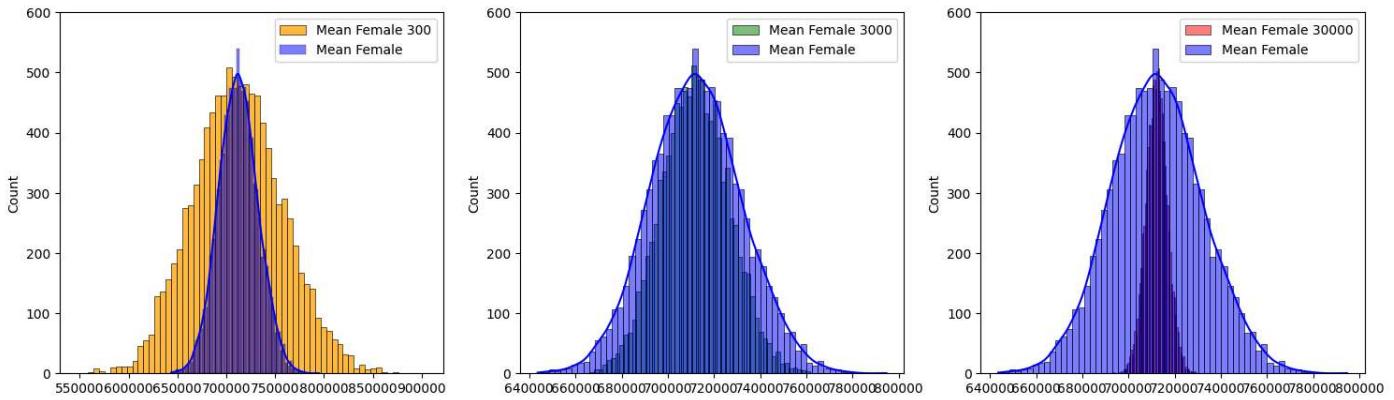
[[900868, 950448], 49580]

```

```

1 plt.figure(num=3, figsize=(18,5))
2
3 plt.subplot(131)
4 sns.histplot(data=mean_female_300, color='orange', label='Mean Female 300')
5 sns.histplot(data=mean_female,kde = True, color='blue', label='Mean Female')
6 # plt.xlim(550000,1000000)
7 plt.ylim(0,600)
8 plt.legend()
9
10 plt.subplot(132)
11 sns.histplot(data=mean_female_3000, color='green', label='Mean Female 3000', alpha=0.5)
12 sns.histplot(data=mean_female,kde = True, color='blue', label='Mean Female')
13 # plt.xlim(550000,1000000)
14 plt.ylim(0,600)
15 plt.legend()
16
17 plt.subplot(133)
18 sns.histplot(data=mean_female_30000, color='red', label='Mean Female 30000', alpha=0.5)
19 sns.histplot(data=mean_female,kde = True, color='blue', label='Mean Female')
20 # plt.xlim(550000,1000000)
21 plt.ylim(0,600)
22
23 plt.legend()
24 plt.show()
25

```



```

1 CI_female_purchase = np.percentile(mean_female, [5, 95])
2 CI_female_purchase = [int(np.round(value)) for value in CI_female_purchase]
3 CI_female_purchase_width = CI_female_purchase[1] - CI_female_purchase[0]
4 [CI_female_purchase , CI_female_purchase_width]

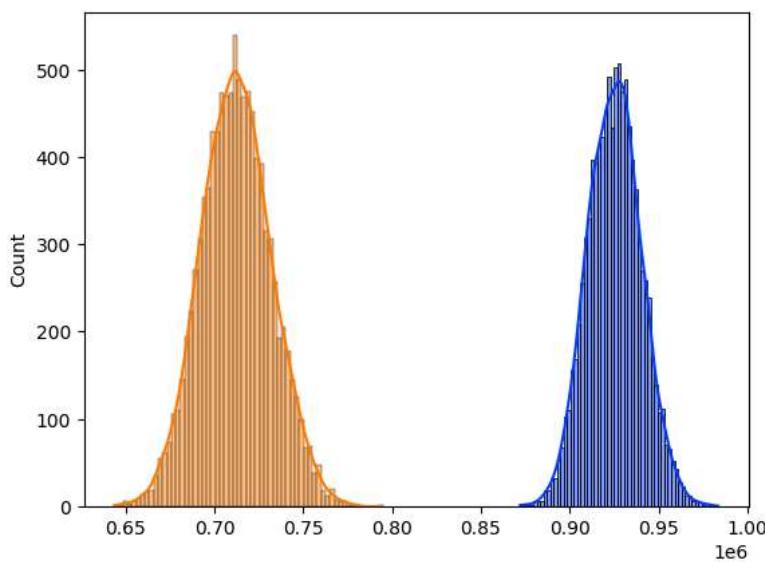
```

```
[[680226, 745747], 65521]
```

```

1 sns.histplot(data=mean_male,kde = True, label='Mean Male')
2 sns.histplot(data=mean_female,kde = True, label='Mean Female')
3 plt.show()

```



The confidence interval for Male and Female is of width 49770 and 64191 respectively. The variation is due to the Sample

- ✓ size and Variance of the Female dataset is lower than the Male dataset. The data is highly dependent on the gender for Black Friday Purchases.

As the sample size increases the CI becomes less. This is caused due to the reduction in the standard error by square root of the sample size, as evident from both the Histograms of male and female with variation in sample size.

As the sample size increases the width of the curve reduces with standard error. The widest dataset is with 300 samples and 30000 is lowest this is expected with increase in sample size. All samples overlap in the order 300 size, 3000 size, dataset size and 30000 size for Males. for Females the 3000 size and dataset size are exchanged in the order, this is due to population size being lower than sample size.

Q5. How does Marital\_Status affect the amount spent?

```

1 df_status0 = df[df['Marital_Status'] == 0].groupby(['User_ID']).agg({'Purchase' : 'sum'}).reset_index()
2 df_status1 = df[df['Marital_Status'] == 1].groupby(['User_ID']).agg({'Purchase' : 'sum'}).reset_index()
3 [df_status0.shape, df_status1.shape]

[(3417, 2), (2474, 2)]

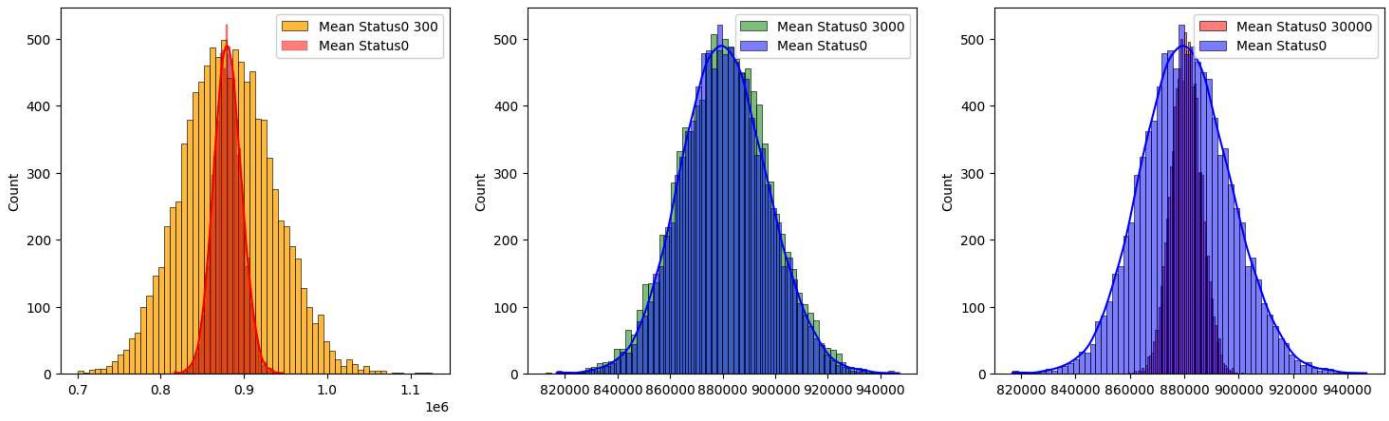
1 mean_status0 = [np.mean(df_status0['Purchase'].sample(len(df_status0['User_ID']), replace = True)) for i in range (10000)]
2 mean_status0_300 = [np.mean(df_status0['Purchase'].sample(300, replace = True)) for i in range (10000)]
3 mean_status0_3000 = [np.mean(df_status0['Purchase'].sample(3000, replace = True)) for i in range (10000)]
4 mean_status0_30000 = [np.mean(df_status0['Purchase'].sample(30000, replace = True)) for i in range (10000)]
5

```

```

1 plt.figure(num=3, figsize=(18,5))
2
3 plt.subplot(131)
4 sns.histplot(data=mean_status0_300, color='orange', label='Mean Status0 300')
5 sns.histplot(data=mean_status0,kde = True, color='red', label='Mean Status0')
6 # plt.xlim(650000,1100000)
7 plt.legend()
8
9 plt.subplot(132)
10 sns.histplot(data=mean_status0_3000, color='green', label='Mean Status0 3000', alpha=0.5)
11 sns.histplot(data=mean_status0,kde = True, color='blue', label='Mean Status0')
12 # plt.xlim(650000,1100000)
13 plt.legend()
14
15 plt.subplot(133)
16 sns.histplot(data=mean_status0_30000, color='red', label='Mean Status0 30000', alpha=0.5)
17 sns.histplot(data=mean_status0,kde = True, color='blue', label='Mean Status0')
18 # plt.xlim(650000,1100000)
19
20 plt.legend()
21 plt.show()
22

```



```

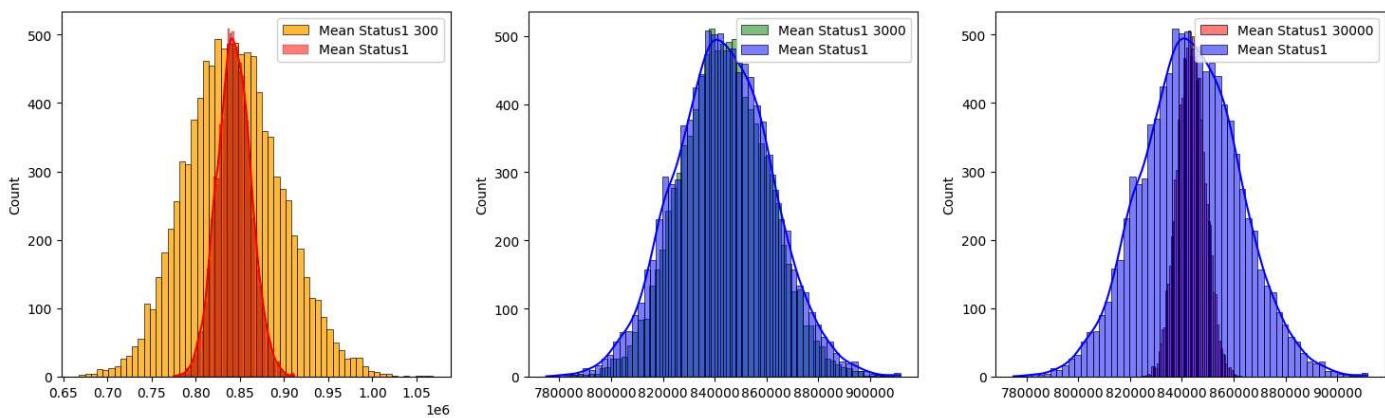
1 mean_status1 = [np.mean(df_status1['Purchase'].sample(len(df_status1['User_ID']), replace = True)) for i in range (10000)]
2 mean_status1_300 = [np.mean(df_status1['Purchase'].sample(300, replace = True)) for i in range (10000)]
3 mean_status1_3000 = [np.mean(df_status1['Purchase'].sample(3000, replace = True)) for i in range (10000)]
4 mean_status1_30000 = [np.mean(df_status1['Purchase'].sample(30000, replace = True)) for i in range (10000)]

```

```

1 plt.figure(num=3, figsize=(18,5))
2
3 plt.subplot(131)
4 sns.histplot(data=mean_status1_300, color='orange', label='Mean Status1 300')
5 sns.histplot(data=mean_status1,kde = True, color='red', label='Mean Status1')
6 # plt.xlim(650000,1100000)
7 plt.legend()
8
9 plt.subplot(132)
10 sns.histplot(data=mean_status1_3000, color='green', label='Mean Status1 3000', alpha=0.5)
11 sns.histplot(data=mean_status1,kde = True, color='blue', label='Mean Status1')
12 # plt.xlim(650000,1100000)
13 plt.legend()
14
15 plt.subplot(133)
16 sns.histplot(data=mean_status1_30000, color='red', label='Mean Status1 30000', alpha=0.5)
17 sns.histplot(data=mean_status1,kde = True, color='blue', label='Mean Status1')
18 # plt.xlim(650000,1100000)
19
20 plt.legend()
21 plt.show()
22

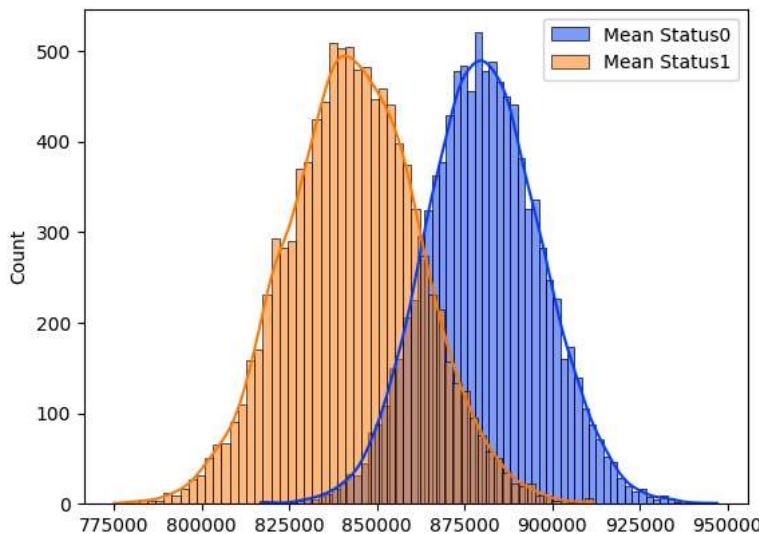
```



```

1 sns.histplot(data=mean_status0,kde = True, label='Mean Status0')
2 sns.histplot(data=mean_status1,kde = True, label='Mean Status1')
3 plt.legend()
4 plt.show()

```



```

1 CI_status0_purchase = np.percentile(mean_status0, [5, 95])
2 CI_status0_purchase = [int(np.round(value)) for value in CI_status0_purchase]
3 CI_status0_purchase_width = CI_status0_purchase[1] - CI_status0_purchase[0]
4 [CI_status0_purchase , CI_status0_purchase_width]

```

```
[[854054, 907992], 53938]
```

```

1 CI_status1_purchase = np.percentile(mean_status1, [5, 95])
2 CI_status1_purchase = [int(np.round(value)) for value in CI_status1_purchase]
3 CI_status1_purchase_width = CI_status1_purchase[1] - CI_status1_purchase[0]
4 [CI_status1_purchase , CI_status1_purchase_width]

```

```
[[812846, 874598], 61752]
```

The confidence interval for 0 and 1 is of width 52940 and 61966 respectively. The variation is due to the Variance of the 0

- ✓ dataset is lower than the 1 dataset and sample sizes are different. The curves are overlapping showing that 0 and 1 are mean are nearby. The data is similar and there is minimum effect of Marital status on the black Friday Purchases.

As the sample size increases the CI becomes less. This is caused due to the reduction in the standard error by square root of the sample size, as evident from both the Histograms of male and female with variation in sample size.

As the sample size increases the width of the curve reduces with standard error. The widest dataset is with 300 samples and 30000 is lowest this is expected with increase in sample size. All samples overlap in the order 300 size, 3000 size, dataset size

and 30000 size for Males. for Females the 3000 size and dataset size are exchanged in the order, this is due to population size being lower than sample size.

## 6. How does Age affect the amount spent?

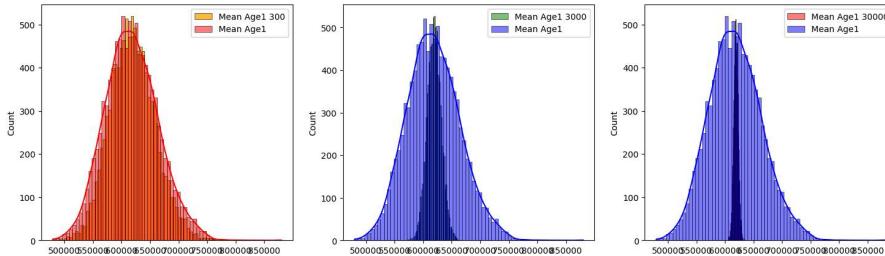
```

1 df_age1 = df[df['Age'] == '0-17'].groupby(['User_ID']).agg({'Purchase' : 'sum'}).reset_index()
2 df_age2 = df[df['Age'] == '18-25'].groupby(['User_ID']).agg({'Purchase' : 'sum'}).reset_index()
3 df_age3 = df[df['Age'] == '26-35'].groupby(['User_ID']).agg({'Purchase' : 'sum'}).reset_index()
4 df_age4 = df[df['Age'] == '36-45'].groupby(['User_ID']).agg({'Purchase' : 'sum'}).reset_index()
5 df_age5 = df[df['Age'] == '46-50'].groupby(['User_ID']).agg({'Purchase' : 'sum'}).reset_index()
6 df_age6 = df[df['Age'] == '51-55'].groupby(['User_ID']).agg({'Purchase' : 'sum'}).reset_index()
7 df_age7 = df[df['Age'] == '55+'].groupby(['User_ID']).agg({'Purchase' : 'sum'}).reset_index()
8 [df_age1.shape , df_age2.shape,df_age3.shape,df_age4.shape,df_age5.shape,df_age6.shape, df_age7.shape]

[(218, 2), (1069, 2), (2053, 2), (1167, 2), (531, 2), (481, 2), (372, 2)]
```

```

1 mean_age1 = [np.mean(df_age1['Purchase']).sample(len(df_age1['User_ID']), replace = True)) for i in range (10000)]
2 mean_age1_300 = [np.mean(df_age1['Purchase'].sample(300, replace = True)) for i in range (10000)]
3 mean_age1_3000 = [np.mean(df_age1['Purchase'].sample(3000, replace = True)) for i in range (10000)]
4 mean_age1_30000 = [np.mean(df_age1['Purchase'].sample(30000, replace = True)) for i in range (10000)]
5 plt.figure(num=3, figsize=(18,5))
6
7 plt.subplot(131)
8 sns.histplot(data=mean_age1_300, color='orange', label='Mean Age1 300')
9 sns.histplot(data=mean_age1,kde = True, color='red', label='Mean Age1')
10 plt.legend()
11
12 plt.subplot(132)
13 sns.histplot(data=mean_age1_3000, color='green', label='Mean Age1 3000', alpha=0.5)
14 sns.histplot(data=mean_age1,kde = True, color='blue', label='Mean Age1')
15
16 plt.legend()
17
18 plt.subplot(133)
19 sns.histplot(data=mean_age1_30000, color='red', label='Mean Age1 30000', alpha=0.5)
20 sns.histplot(data=mean_age1,kde = True, color='blue', label='Mean Age1')
21
22
23 plt.legend()
24 plt.show()
25
```



```
1 mean_age2 = [np.mean(df_age2['Purchase'].sample(len(df_age2['User_ID']), replace = True)) for i in range (10000)]
2 mean_age2_300 = [np.mean(df_age2['Purchase'].sample(300, replace = True)) for i in range (10000)]
3 mean_age2_3000 = [np.mean(df_age2['Purchase'].sample(3000, replace = True)) for i in range (10000)]
4 mean_age2_30000 = [np.mean(df_age2['Purchase'].sample(30000, replace = True)) for i in range (10000)]
5 plt.figure(num=3, figsize=(18,5))
6
7 plt.subplot(131)
8 sns.histplot(data=mean_age2_300, color='orange', label='Mean Age2 300')
9 sns.histplot(data=mean age2.kde = True, color='red', label='Mean Age2')
```