

```
In [54]: #import python libraries

import numpy as np
import pandas as pd
import matplotlib.pyplot as plt # visualizing data
%matplotlib inline
import seaborn as sns

In [56]: #import csv file
df = pd.read_csv('G:\\swati\\Sales Data.csv', encoding='unicode_escape')

In [53]: df.shape
Out[53]: (11251, 15)

In [32]: df.head(10)
Out[32]:
   User_ID  Cust_Name  Product_ID  Gender  Age Group  Age  Marital_Status  State  Zone  Occupation  Product_Category  Orders  Amount  Status  unnamed1
0  1002903  Sanskriti  P00125942  F      26-35  28      0  Maharashtra  Western  Healthcare  Auto  1  23952.00  NaN  NaN
1  1000732  Karkit  P00110942  F      26-35  35      1  Andhra Pradesh  Southern  Govt  Auto  3  23934.00  NaN  NaN
2  1001990  Bindu  P00118942  F      26-35  35      1  Uttar Pradesh  Central  Automobile  Auto  3  23924.00  NaN  NaN
3  1001425  Sudovi  P00237842  M      0-17  16      0  Karnataka  Southern  Construction  Auto  2  23912.00  NaN  NaN
4  1000588  Joni  P00057942  M      26-35  28      1  Gujarat  Western  Food Processing  Auto  2  23877.00  NaN  NaN
5  1000588  Joni  P00057942  M      26-35  28      1  Himachal Pradesh  Northern  Food Processing  Auto  1  23877.00  NaN  NaN
6  1001132  Balk  P00018042  F      18-25  25      1  Uttar Pradesh  Western  Lawyer  Auto  4  23841.00  NaN  NaN
7  1002092  Shivangi  P00273442  F      55+  61      0  Maharashtra  Central  IT Sector  Auto  1  NaN  NaN  NaN
8  1000224  Kushal  P0005642  M      26-35  35      0  Uttar Pradesh  Central  Govt  Auto  2  23809.00  NaN  NaN
9  1002650  Gmny  P0031142  F      26-35  26      1  Andhra Pradesh  Southern  Media  Auto  4  23799.00  NaN  NaN

In [5]: df.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 11251 entries, 0 to 11250
Data columns (total 15 columns):
#   Column  Non-Null Count  Dtype
---  ---
0  User_ID  11251 non-null  int64
1  Cust_Name  11251 non-null  object
2  Product_ID  11251 non-null  object
3  Gender  11251 non-null  object
4  Age Group  11251 non-null  object
5  Age  11251 non-null  int64
6  Marital_Status  11251 non-null  int64
7  State  11251 non-null  object
8  Zone  11251 non-null  object
9  Occupation  11251 non-null  object
10 Product_Category  11251 non-null  object
11 Orders  11251 non-null  int64
12 Amount  11239 non-null  float64
13 Status  0 non-null  float64
14 unnamed1  0 non-null  float64
dtypes: float64(3), int64(4), object(8)
memory usage: 1.1+ MB

In [33]: #drop unrelated/blank columns
df.drop(['Status', 'unnamed1'], axis=1, inplace=True)

In [21]: df.info()
<class 'pandas.core.frame.DataFrame'>
Int64Index: 0 entries
Data columns (total 13 columns):
#   Column  Non-Null Count  Dtype
---  ---
0  User_ID  0 non-null  int64
1  Cust_Name  0 non-null  object
2  Product_ID  0 non-null  object
3  Gender  0 non-null  object
4  Age Group  0 non-null  object
5  Age  0 non-null  int64
6  Marital_Status  0 non-null  int64
7  State  0 non-null  object
8  Zone  0 non-null  object
9  Occupation  0 non-null  object
10 Product_Category  0 non-null  object
11 Orders  0 non-null  int64
12 Amount  0 non-null  float64
dtypes: float64(1), int64(4), object(8)
memory usage: 0.0+ bytes

In [24]: pd.isnull(df)

In [42]:
   User_ID  Cust_Name  Product_ID  Gender  Age Group  Age  Marital_Status  State  Zone  Occupation  Product_Category  Orders  Amount
0  1002903  Sanskriti  P00125942  F      26-35  28      0  Maharashtra  Western  Healthcare  Auto  1  23952
1  1000732  Karkit  P00110942  F      26-35  35      1  Andhra Pradesh  Southern  Govt  Auto  3  23934
2  1001990  Bindu  P00118942  F      26-35  35      1  Uttar Pradesh  Central  Automobile  Auto  3  23924
3  1001425  Sudovi  P00237842  M      0-17  16      0  Karnataka  Southern  Construction  Auto  2  23912
4  1000588  Joni  P00057942  M      26-35  28      1  Gujarat  Western  Food Processing  Auto  2  23877
...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...
11246  1000595  Manning  P00236942  M      18-25  19      1  Maharashtra  Western  Chemical  Office  4  370
11247  1004089  Reichenbach  P00171342  M      26-35  33      0  Haryana  Northern  Healthcare  Veterinary  3  367
11248  1001209  Oshin  P00201342  F      36-45  40      0  Madhya Pradesh  Central  Textile  Office  4  213
11249  1004023  Noorjan  P00059442  M      36-45  37      0  Karnataka  Southern  Agriculture  Office  3  206
11250  1002714  Brumley  P00281742  F      18-25  19      0  Maharashtra  Western  Healthcare  Office  3  188
11239 rows x 13 columns

In [43]: # describe() method returns description of the data in the DataFrame (i.e. count, mean, std, etc)
df.describe()
Out[43]:
   User_ID  Age  Marital_Status  Orders  Amount
count  1.123900e+04  11239.000000  11239.000000  11239.000000  11239.000000
mean  1.003004e+06  35.410357  0.420055  2.489634  9453.610553
std  1.716039e+03  12.753866  0.493589  1.114967  5222.355168
min  1.000001e+06  12.000000  0.000000  1.000000  188.000000
25%  1.001492e+06  27.000000  0.000000  2.000000  5443.000000
50%  1.003064e+06  33.000000  0.000000  2.000000  8109.000000
75%  1.004426e+06  43.000000  1.000000  3.000000  12675.000000
max  1.006040e+06  92.000000  1.000000  4.000000  23952.000000
```

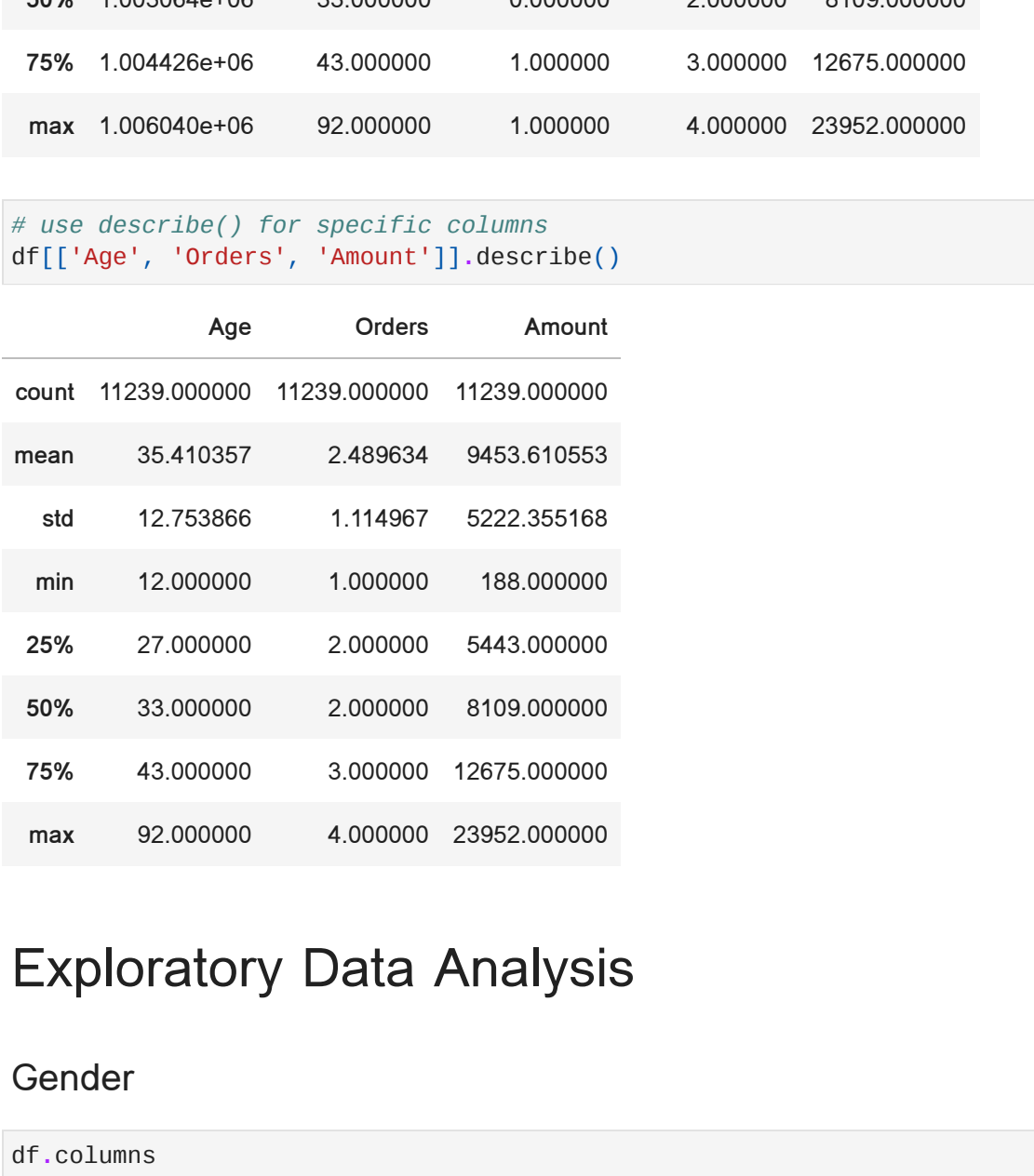
```
In [45]: df.columns
Out[45]: Index(['User_ID', 'Cust_name', 'Product_ID', 'Gender', 'Age Group', 'Age', 'Marital_Status', 'State', 'Zone', 'Occupation', 'Product_Category', 'Orders', 'Amount'],
              dtype='object')

In [47]: #plotting a bar chart for Gender and it's count
ax = sns.countplot(x = 'Gender', data = df)
for bars in ax.containers:
    ax.bar_label(bars)

Out[47]:
   Gender
F      7832
M      3407

In [46]: #plotting a bar chart for gender vs total amount
sales_gen = df.groupby(['Gender'], as_index=False)['Amount'].sum().sort_values(by='Amount', ascending=False)
sns.barplot(x = 'Gender', y = 'Amount', data = sales_gen)

Out[46]:
<Axes: xlabel='Gender', ylabel='Amount'>
```



From above graphs we can see that most of the buyers are females and even the purchasing power of females are greater than men

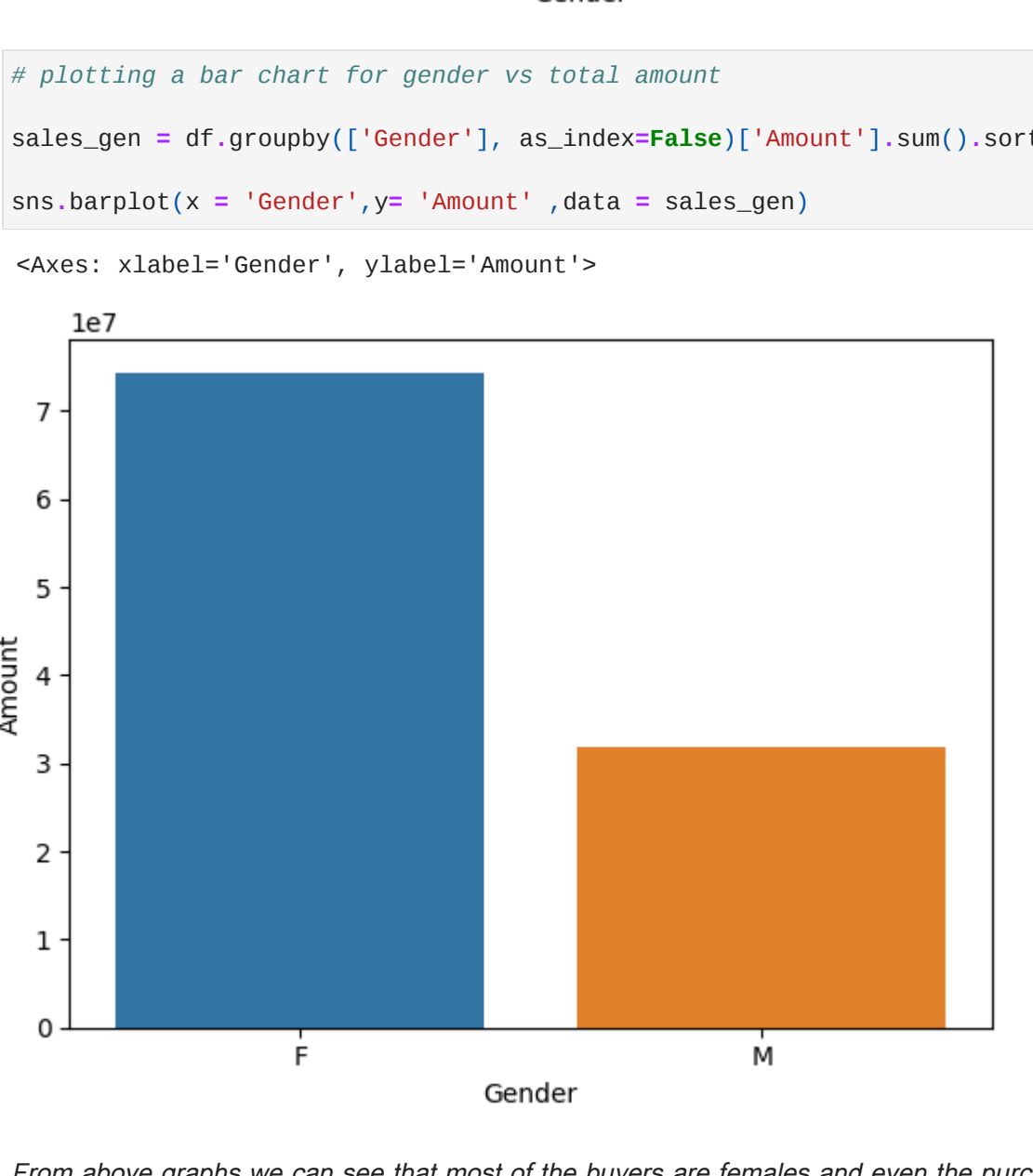
Age

```
In [17]: ax = sns.countplot(data = df, x = 'Age Group', hue = 'Gender')
for bars in ax.containers:
    ax.bar_label(bars)

Out[17]:
   Age Group  Gender
26-35      F      3269
26-35      M      1272
0-17       F      162
0-17       M      134
18-25     F      1305
18-25     M       574
51-55     F       553
51-55     M       272
46-50     F       693
46-50     M       290
55+       F       272
55+       M       155
36-45     F      1578
36-45     M       705

In [18]: # Total Amount vs Age Group
sales_age = df.groupby(['Age Group'], as_index=False)['Amount'].sum().sort_values(by='Amount', ascending=False)
sns.barplot(x = 'Age Group', y = 'Amount', data = sales_age)

Out[18]:
<Axes: xlabel='Age Group', ylabel='Amount'>
```

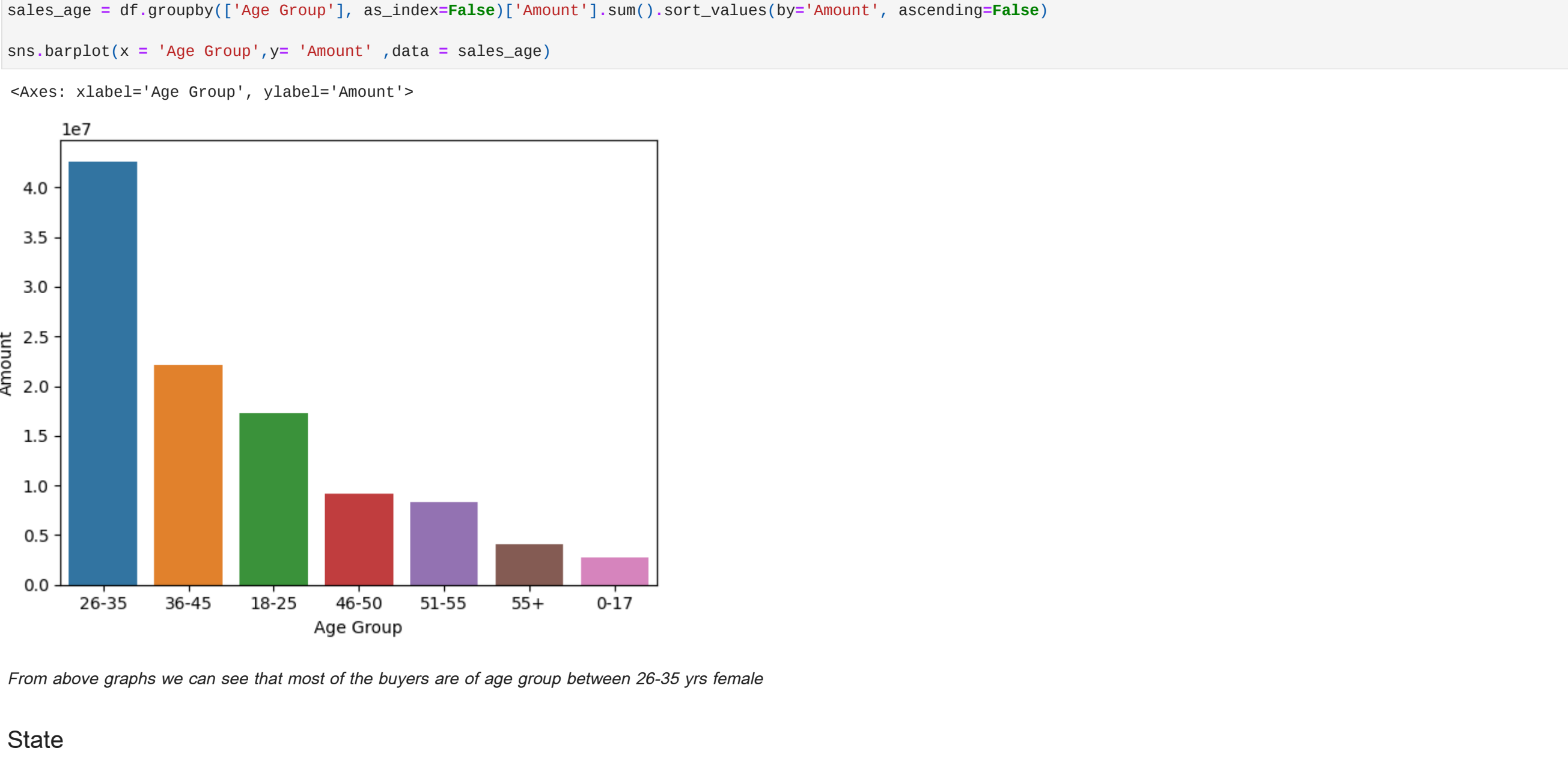


From above graphs we can see that most of the buyers are of age group between 26-35 yrs female

State

```
In [19]: # total number of orders from top 10 states
sales_state = df.groupby(['State'], as_index=False)['Orders'].sum().sort_values(by='Orders', ascending=False).head(10)
sns.set(rc={'figure.figsize':(15,6)})
sns.barplot(data = sales_state, x = 'State', y = 'Orders')

Out[19]:
<Axes: xlabel='State', ylabel='Orders'>
```



From above graphs we can see that most of the orders & total sales/amount are from Uttar Pradesh, Maharashtra and Karnataka respectively

Marital Status

```
In [21]: ax = sns.countplot(data = df, x = 'Marital_Status')
for bars in ax.containers:
    ax.bar_label(bars)

Out[21]:
   Marital_Status
0      6518
1      4721

In [22]: sales_state = df.groupby(['Marital_Status', 'Gender'], as_index=False)['Amount'].sum().sort_values(by='Amount', ascending=False)
sns.set(rc={'figure.figsize':(6,5)})
sns.barplot(data = sales_state, x = 'Marital_Status', y = 'Amount', hue='Gender')

Out[22]:
<Axes: xlabel='Marital_Status', ylabel='Amount'>
```



From above graphs we can see that most of the buyers are married (women) and they have high purchasing power

Occupation

```
In [23]: sns.set(rc={'figure.figsize':(25,5)})
ax = sns.countplot(data = df, x = 'Occupation')
for bars in ax.containers:
    ax.bar_label(bars)

Out[23]:
   Occupation
IT Sector  1583
Healthcare 1408
Govt      854
Automobile 565
Construction 414
Food Processing 423
Lawyer     531
Media      637
Banking Occupation 501
Retail     1137
IT Sector  1583
Aviation   1310
Hospitality 703
Agriculture 283
Textile    349
Chemical   541

In [24]: sales_state = df.groupby(['Occupation'], as_index=False)['Amount'].sum().sort_values(by='Amount', ascending=False)
sns.set(rc={'figure.figsize':(20,5)})
sns.barplot(data = sales_state, x = 'Occupation', y = 'Amount')

Out[24]:
<Axes: xlabel='Occupation', ylabel='Amount'>
```



From above graphs we can see that most of the buyers are working in IT, Healthcare and Aviation sector

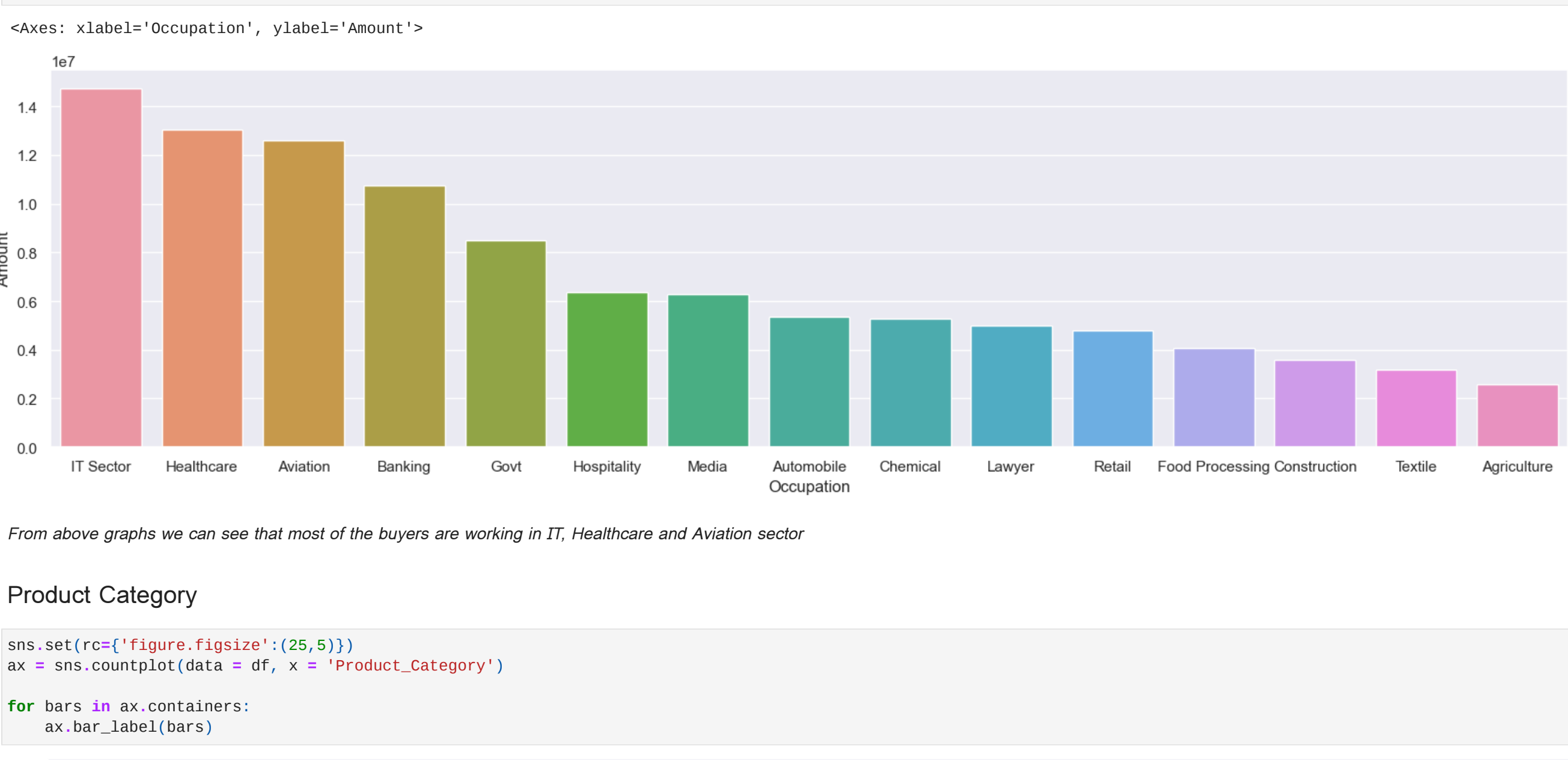
Product Category

```
In [52]: sns.set(rc={'figure.figsize':(25,5)})
ax = sns.countplot(data = df, x = 'Product_Category')
for bars in ax.containers:
    ax.bar_label(bars)

Out[52]:
   Product_Category
Food 2450
Clothing & Apparel 1029
Electronics & Gadgets 382
Footwear & Shoes 366
Furniture 386
Games & Toys 306
Sports Products 103
Decor 96
Beauty 422
Household Items 220
Pet Care 212
Veterinary 81
Office 113

In [26]: sales_state = df.groupby(['Product_Category'], as_index=False)['Amount'].sum().sort_values(by='Amount', ascending=False).head(10)
sns.set(rc={'figure.figsize':(20,5)})
sns.barplot(data = sales_state, x = 'Product_Category', y = 'Amount')

Out[26]:
<Axes: xlabel='Product_Category', ylabel='Amount'>
```



From above graphs we can see that most of the sold products are from Food, Clothing and Electronics category

```
In [28]: # Top 10 most sold products (same thing as above)
fig1, ax1 = plt.subplots(figsize=(12,7))
df.groupby('Product_ID')['Orders'].nlargest(10).sort_values(ascending=False).plot(kind='bar')

Out[28]:
<Axes: xlabel='Product_ID', ylabel='Orders'>
```



Conclusion:

Married women age group 26-35 yrs from UP, Maharashtra and Karnataka working in IT, Healthcare and Aviation are more likely to buy products from Food, Clothing and Electronics category

Thank you!