

```
In [9]: # Import python libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt # visualizing data
%matplotlib inline #tells the IPython environment to draw the plots immediately after the current cell
import seaborn as sns
```

```
In [11]: # Import csv file
df = pd.read_csv('D:\ml\1 Sales Data.csv', encoding='unicode_escape')
#to avoid encoding error use 'unicode_escape'
```

```
In [12]: df.shape # tells us how many rows and columns are in the file
(11251, 15)
```

```
Out[12]:
```

```
In [14]: df.head(10) # shows top 5 rows of the file by default but in bracket you can customize the number of rows to be displayed
```

Out[14]:

	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	Karik	P00110942	F	26-35	35	1		Andhra Pradesh	Southern	Govt	Auto	3	23934.00	NaN	NaN
2	1001990	Bindu	P00118542	F	26-35	35	1		Uttar Pradesh	Central	Automobile	Auto	3	23924.00	NaN	NaN
3	1001425	Sudevi	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	Central	Lawyer	Auto	4	23841.00	NaN	NaN
7	1002932	Shivangi	P00273442	F	55+	61	0		Maharashtra	Western	IT Sector	Auto	1	NaN	NaN	NaN
8	1003244	Kunal	P00205642	M	26-35	35	0		Uttar Pradesh	Central	Govt	Auto	2	23809.00	NaN	NaN
9	1003550	Gimny	P00031142	F	26-35	26	1		Andhra Pradesh	Southern	Media	Auto	4	23799.99	NaN	NaN

```
In [15]: df.info() # df means dataframe and this function gives information about the file
```

```
<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: 1.3+ MB
```

```
In [16]: #drop unrelated/blank columns
df.drop(['Status', 'unnamed1'], axis=1, inplace=True)
# axis=1 means complete row
# inplace=True means implement the required task
```

```
In [19]: #check for null values
pd.isnull(df).sum()
```

Out[19]:

```
User_ID      0
Cust_name    0
Product_ID   0
Gender        0
Age Group    0
Age           0
Marital_Status 0
Zone          0
Occupation    0
Product_Category 0
Orders        0
Amount       12
dtype: int64
```

```
In [30]: # drop null values
df.dropna(inplace=True)
```

```
In [32]: # change data type
df['Amount'] = df['Amount'].astype('int')
```

```
In [33]: df['Amount'].dtypes # use to check datatype of a column
```

```
Out[33]: dtype('int32')
```

```
In [11]: df.columns # shows name of all columns in a file
```

```
Out[11]: Index(['User_ID', 'Cust_name', 'Product_ID', 'Gender', 'Age Group', 'Age', 'Marital_Status', 'State', 'Zone', 'Occupation', 'Product_Category', 'Orders', 'Amount'], dtype='object')
```

```
In [12]: #rename column
df.rename(columns={'Marital_Status':'Shaadi'})
```

Out[12]:

	User_ID	Cust_name	Product_ID	Gender	Age Group	Age	Shaadi		State	Zone	Occupation	Product_Category	Orders	Amount
0	1002903	Sanskriti	P00125942	F	26-35	28	0		Maharashtra	Western	Healthcare	Auto	1	23952
1	1000732	Karik	P00110942	F	26-35	35	1		Andhra Pradesh	Southern	Govt	Auto	3	23934
2	1001990	Bindu	P00118542	F	26-35	35	1		Uttar Pradesh	Central	Automobile	Auto	3	23924
3	1001425	Sudevi	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	1000695	Manning	P00296942	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	1001109	Oshin	P00201342	F	36-45	40	0		Madhya Pradesh	Central	Textile	Office	4	213
11249	1004023	Noonan	P00059442	M	36-45	37	0		Karnataka	Southern	Agriculture	Office	3	206
11250	1002744	Brumley	P00281742	F	18-25	19	0		Maharashtra	Western	Healthcare	Office	3	188

11239 rows × 13 columns

```
In [13]: # describe() method returns description of the data in the DataFrame (i.e. count, mean, std, etc) it works only with numeric columns
df.describe()
```

Out[13]:

	User_ID	Age	Marital_Status	Orders	Amount
count	11239.000000	11239.000000	11239.000000	11239.000000	11239.000000
mean	1.002094e+06	25.410357	0.420955	2.489634	9453.610553
std	1.716039e+03	12.753866	0.482699	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%	1.003056e+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 [14]: # use describe() for specific columns
df[['Age', 'Orders', 'Amount']].describe()
```

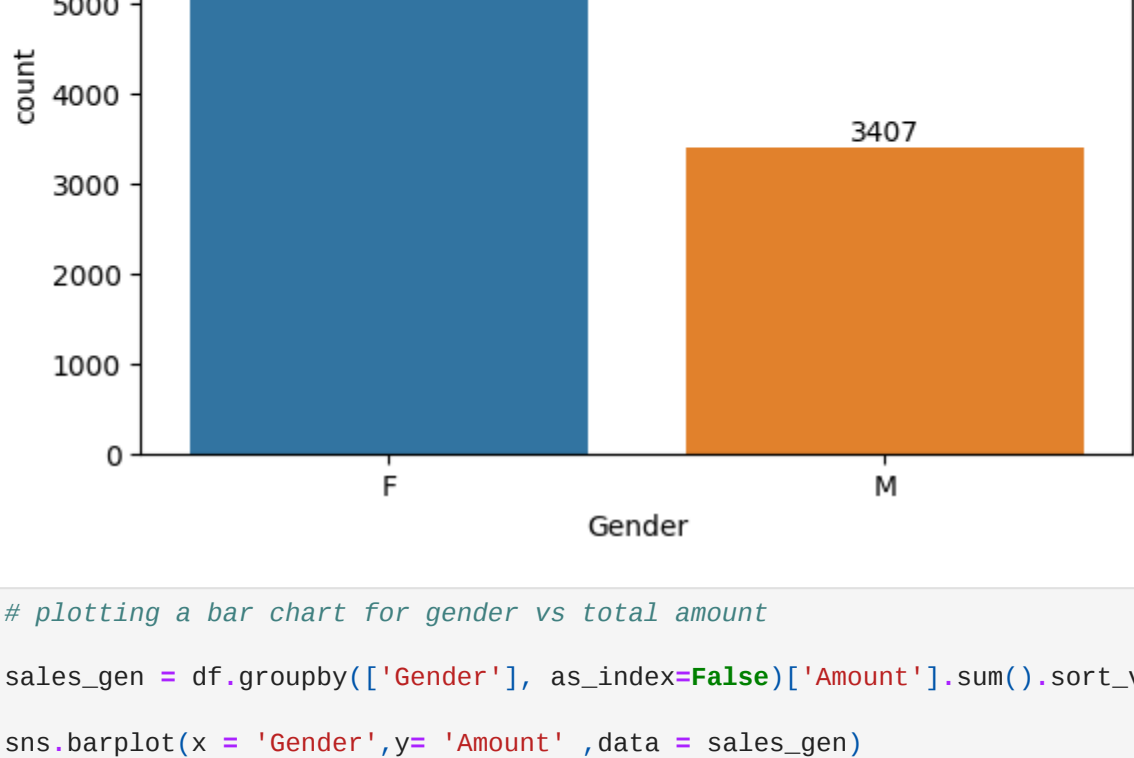
Out[14]:

	Age	Orders	Amount
count	11239.000000	11239.000000	11239.000000
mean	25.410357	2.489634	9453.610553
std	12.753866	1.114967	5222.355168
min	12.000000	1.000000	188.000000
25%	27.000000	2.000000	5443.000000
50%	33.000000	2.000000	8109.000000
75%	43.000000	3.000000	12675.000000
max	92.000000	4.000000	23952.000000

## Exploratory Data Analysis

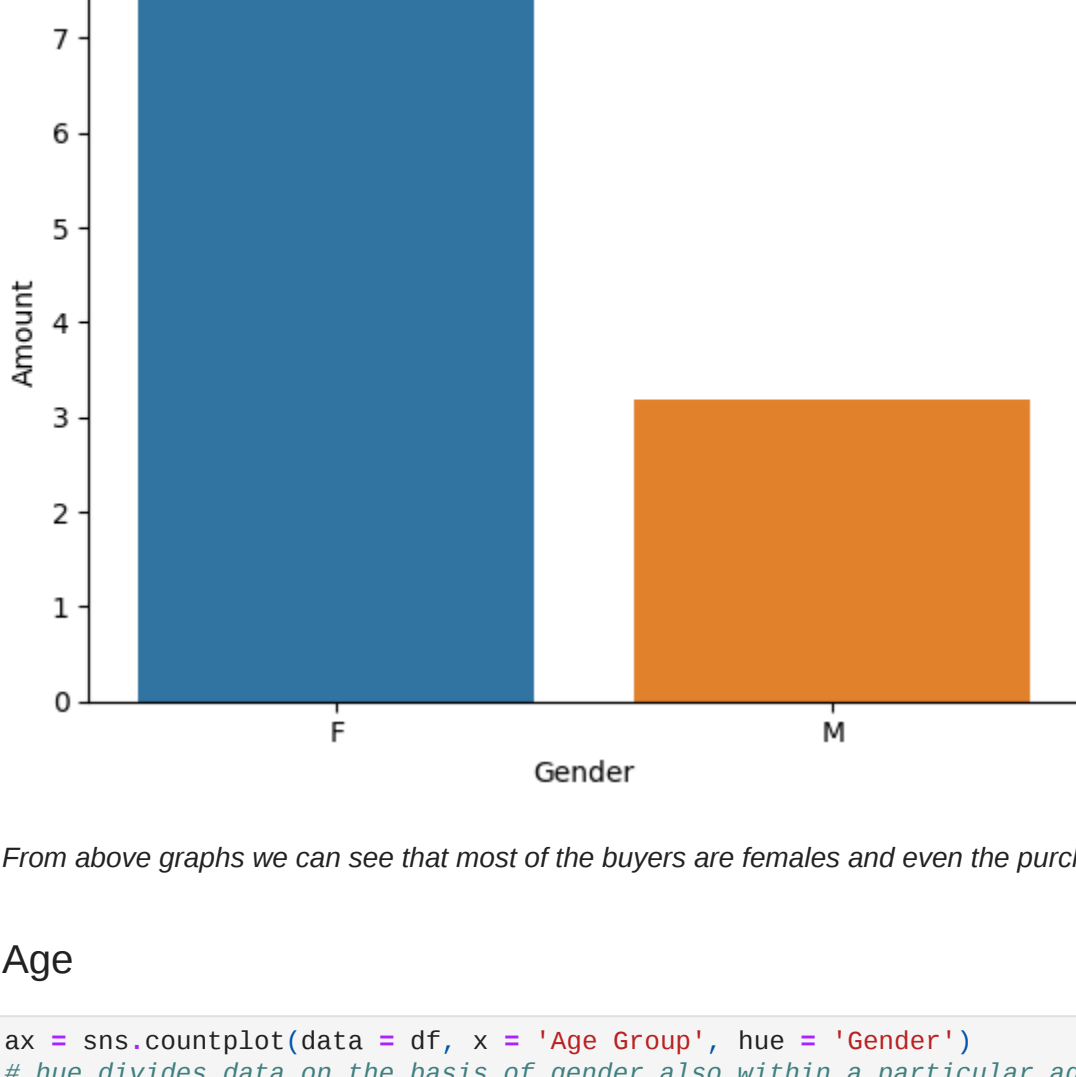
### Gender

```
In [19]: # plotting a bar chart for Gender and it's count
ax = sns.countplot(x = 'Gender', data = df)
for bars in ax.containers: # used for adding labels to the bars
    ax.bar_label(bars)
```



```
In [16]: # 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[16]: <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

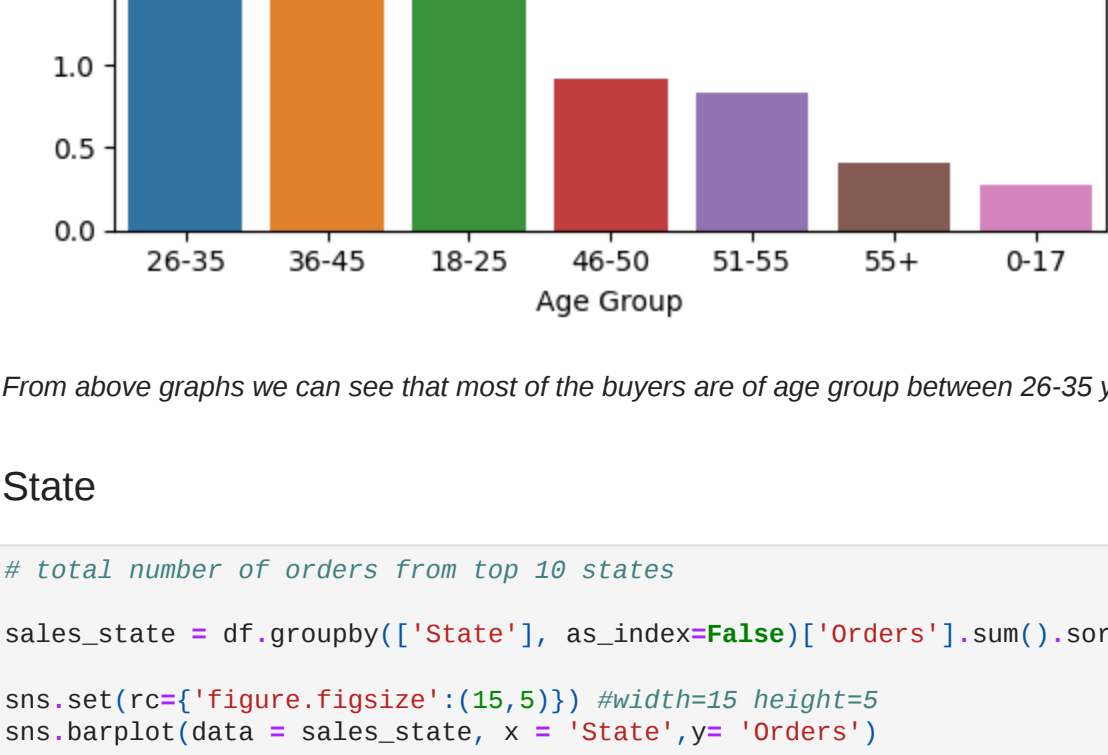
### Age

```
In [17]: ax = sns.countplot(data = df, x = 'Age Group', hue = 'Gender')
# hue divides data on the basis of gender also within a particular age group
for bars in ax.containers:
    ax.bar_label(bars)
```



```
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'>
```

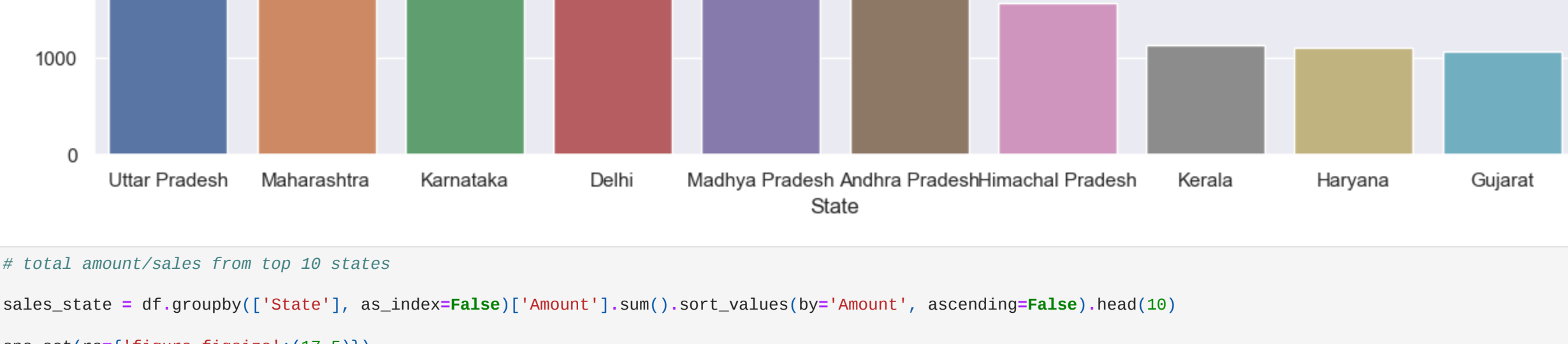


From above graphs we can see that most of the buyers are of age group between 26-35 yrs female

### State

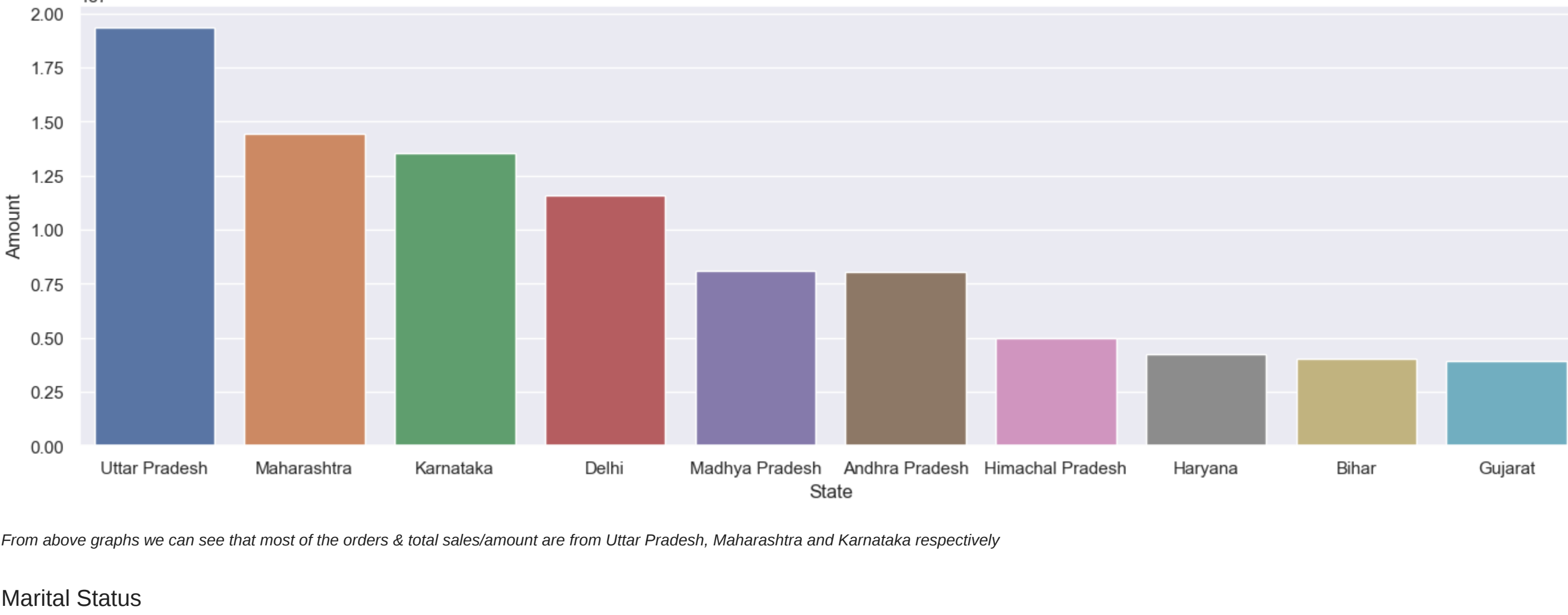
```
In [19]: # Total number of orders from top 10 states
sales_state = df.groupby(['State'], as_index=False)['Orders'].sum().sort_values(by='Orders', ascending=False).head(10)
sns.set(rc={'figure.figsize': (15,5)})
sns.barplot(data = sales_state, x = 'State', y= 'Orders')
```

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



```
In [34]: # total amount/sales from top 10 states
sales_state = df.groupby(['State'], as_index=False)['Amount'].sum().sort_values(by='Amount', ascending=False).head(10)
sns.set(rc={'figure.figsize': (17,5)})
sns.barplot(data = sales_state, x = 'State', y= 'Amount')
```

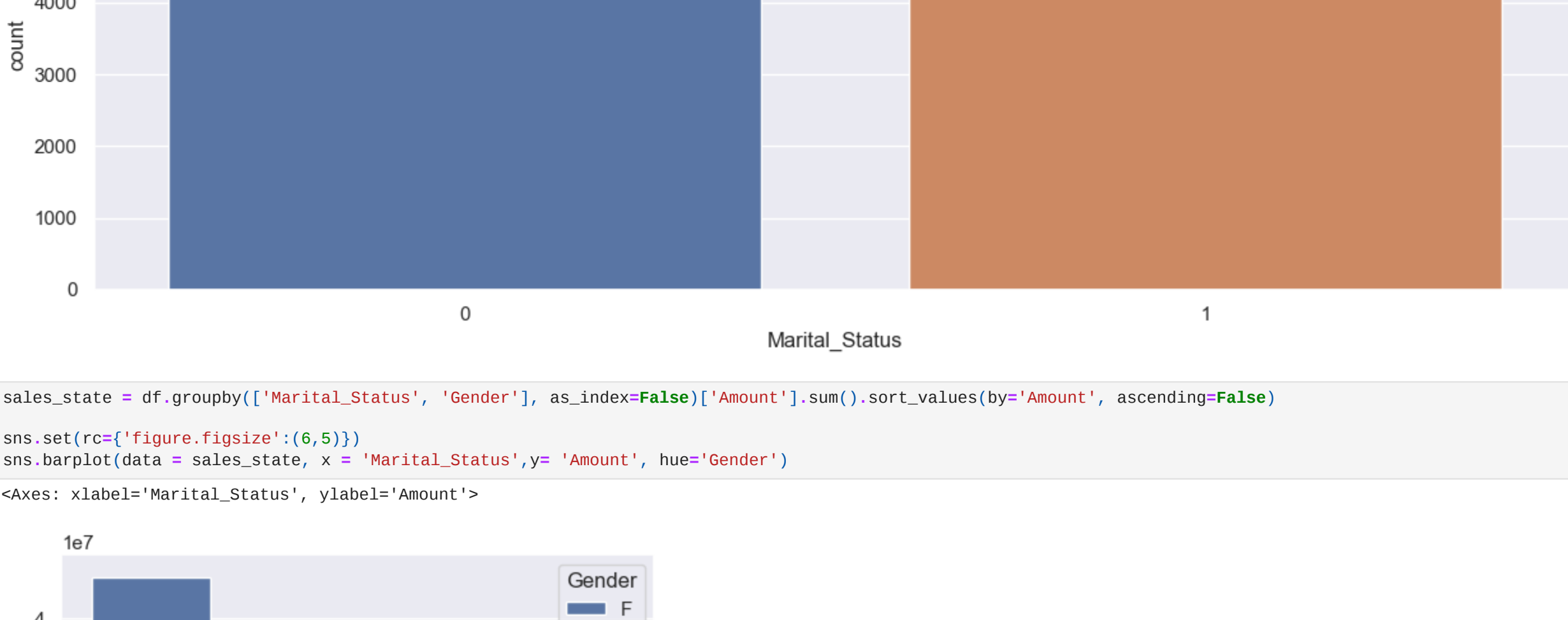
```
Out[34]: <Axes: xlabel='State', ylabel='Amount'>
```



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

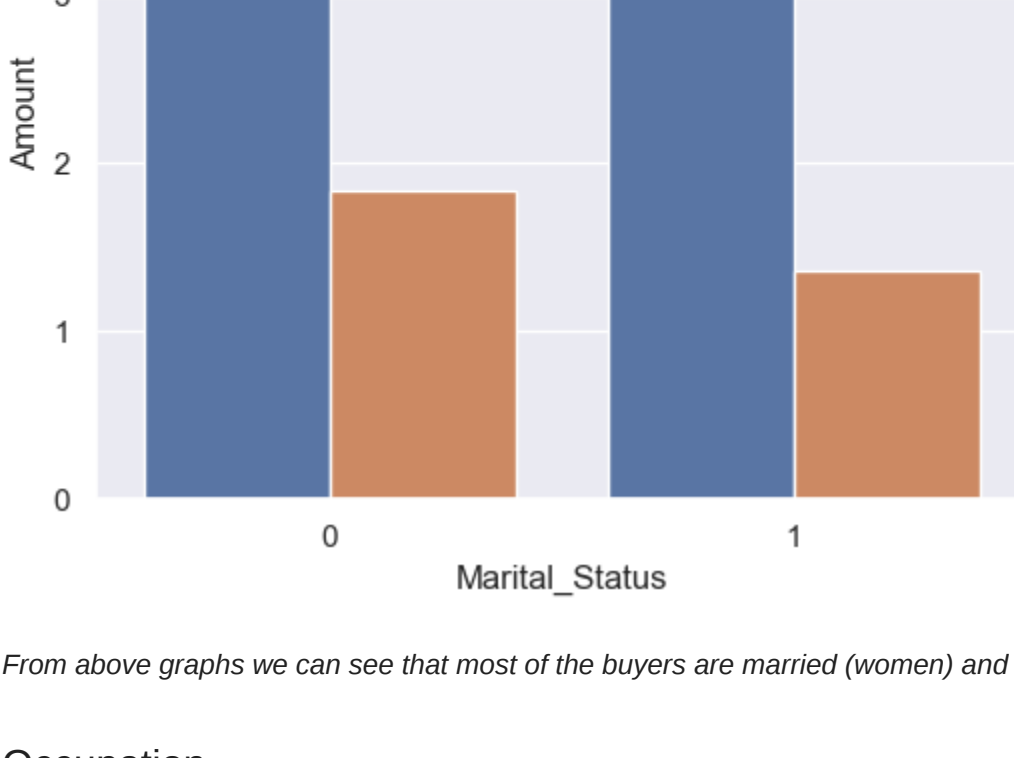
### Marital Status

```
In [21]: ax = sns.countplot(data = df, x = 'Marital_Status')
sns.set(rc={'figure.figsize': (7,5)})
for bars in ax.containers:
    ax.bar_label(bars)
```



```
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')
```

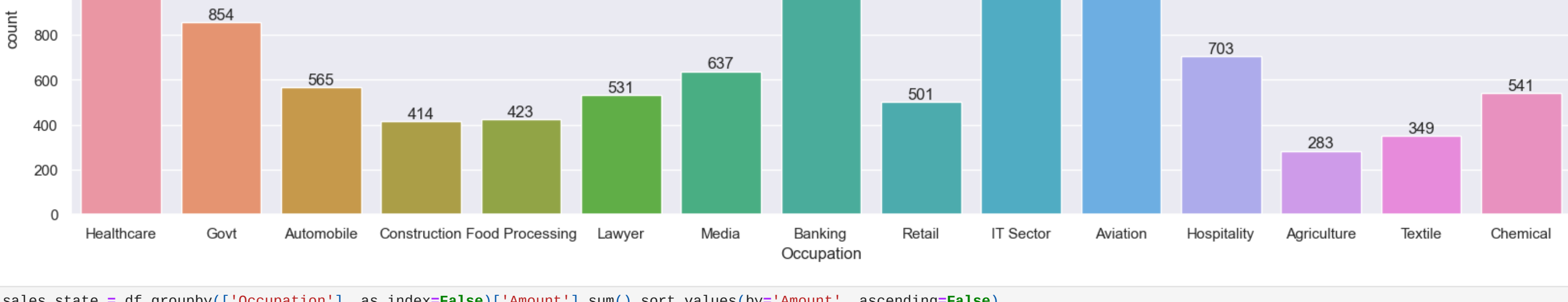
```
Out[22]: <Axes: xlabel='Marital_Status', ylabel='Amount'>
```



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

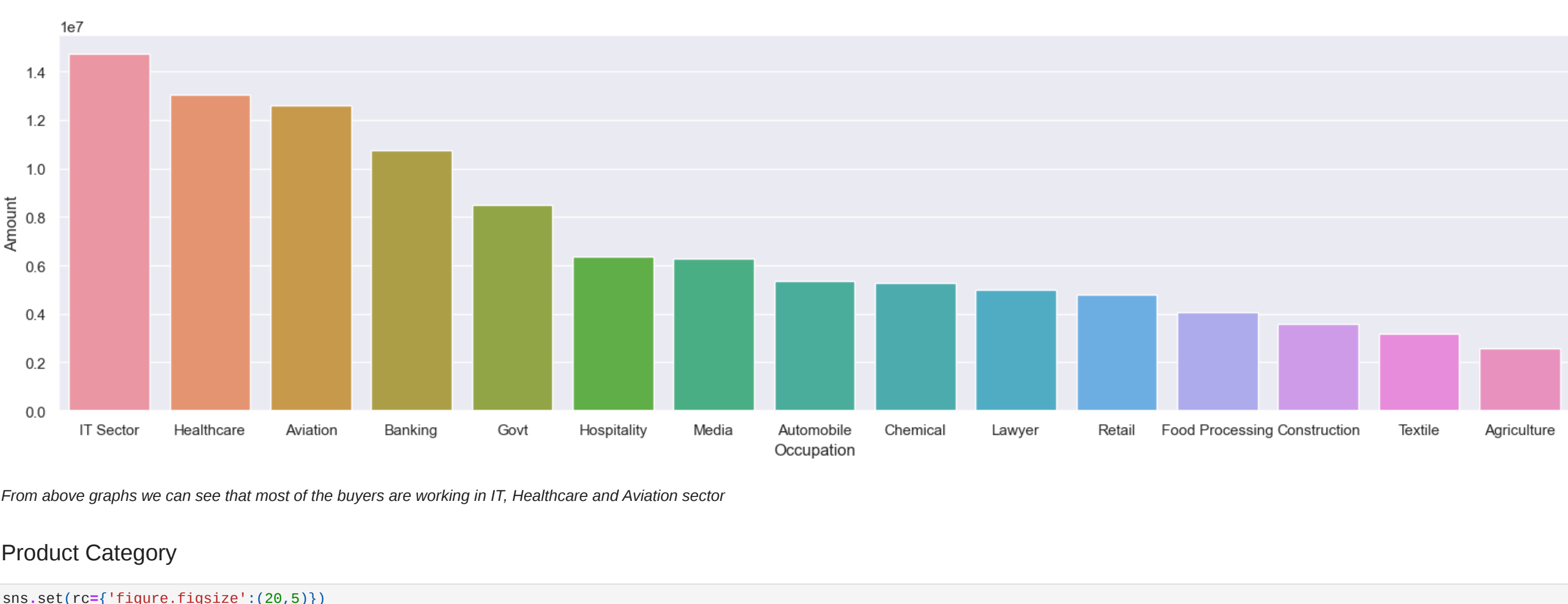
### Occupation

```
In [23]: sns.set(rc={'figure.figsize': (20,5)})
ax = sns.countplot(data = df, x = 'Occupation')
for bars in ax.containers:
    ax.bar_label(bars)
```



```
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')
```

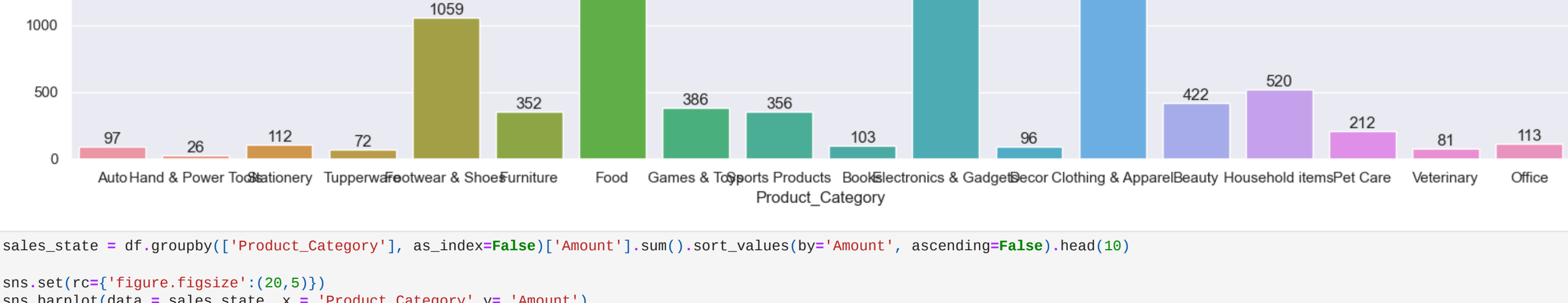
```
Out[24]: <Axes: xlabel='Occupation', ylabel='Amount'>
```



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

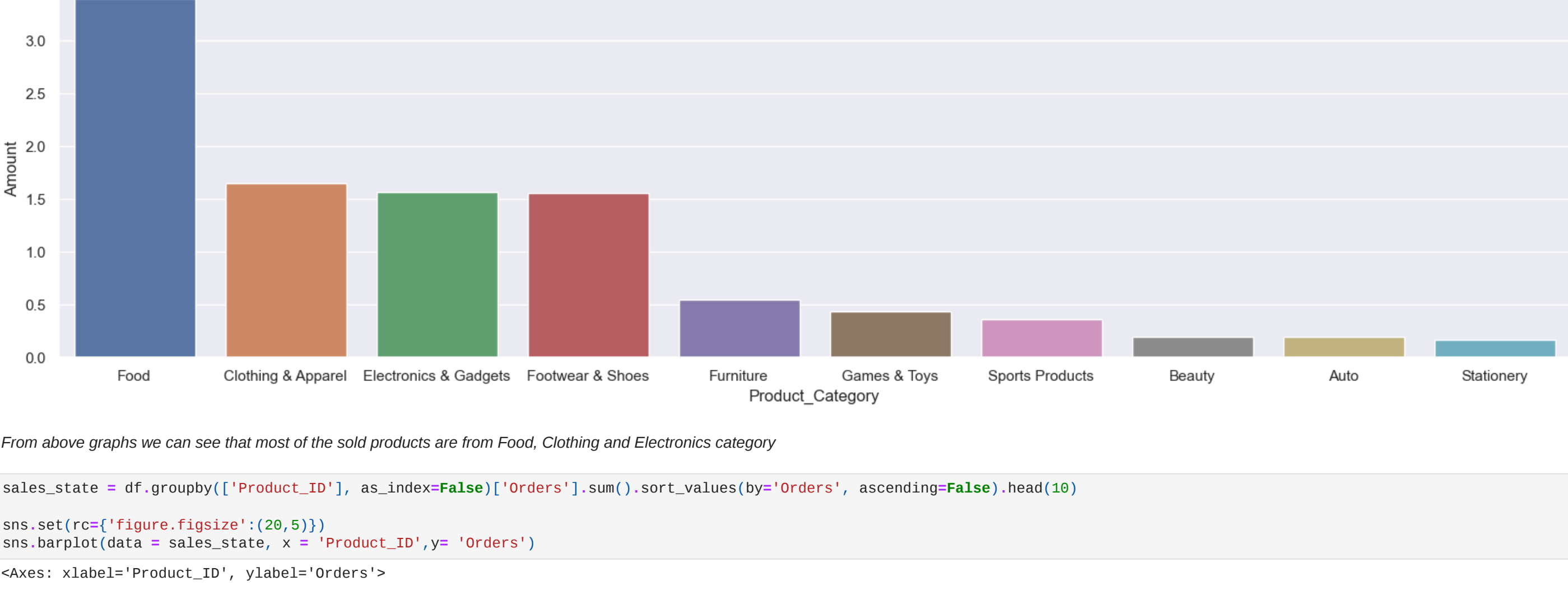
### Product Category

```
In [29]: sns.set(rc={'figure.figsize': (20,5)})
ax = sns.countplot(data = df, x = 'Product_Category')
for bars in ax.containers:
    ax.bar_label(bars)
```



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

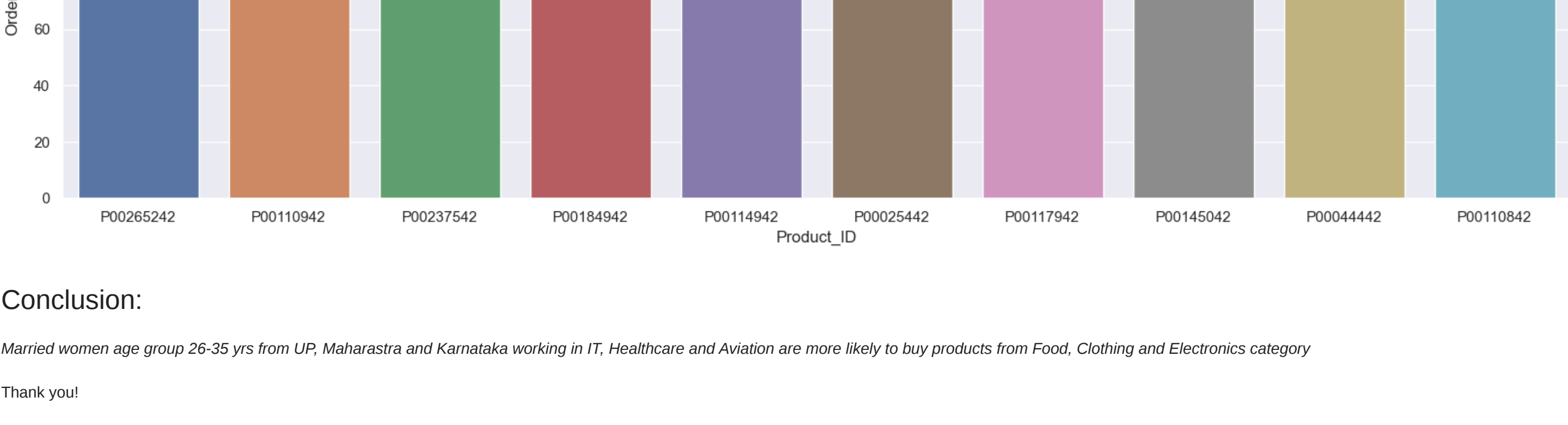
```
Out[26]: <Axes: xlabel='Product_Category', ylabel='Orders'>
```



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

```
In [27]: sales_state = df.groupby(['Product_ID'], as_index=False)['Orders'].sum().sort_values(by='Orders', ascending=False).head(10)
sns.set(rc={'figure.figsize': (20,5)})
sns.barplot(data = sales_state, x = 'Product_ID', y= 'Orders')
```

```
Out[27]: <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!