

```
In [8]: # Import python libraries

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

In [9]: # Import csv file
df = pd.read_csv(r"C:\Users\Admin\Downloads\Python_Diwali_Sales_Analysis-main\Python_Diwali_Sales_Analysis-main\Diwali Sales Data.csv", encoding='unicode_escape')

In [10]: df.shape
Out[10]: (11251, 15)

In [11]: df.head()
Out[11]:
   User_ID  Cust_name  Product_ID  Gender  Age Group  Age  Marital_Status  State  Zone  Occupation  Product_Category  Orders  Amount  Status  unnamed1
0  1002903  Sanskriti  P0012942      F    26-35    28      0  Maharashtra  Western  Healthcare      Auto      1      23952.0  NaN      NaN
1  1000732  Karthik  P00110942      F    26-35    35      1  Andhra Pradesh  Southern  Govt      Auto      3      23934.0  NaN      NaN
2  1001990  Bindu  P00118542      F    26-35    35      1  Uttar Pradesh  Central  Automobile      Auto      3      23924.0  NaN      NaN
3  1001425  Sudeshi  P00237842      M    0-17    16      0  Karnataka  Southern  Construction      Auto      2      23912.0  NaN      NaN
4  1000588  Joni  P00057942      M    26-35    28      1  Gujarat  Western  Food Processing      Auto      2      23877.0  NaN      NaN

In [12]: 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: 1.3+ MB

In [13]: #drop unrelated/blank columns
df.drop(['Status', 'unnamed1'], axis=1, inplace=True)

In [14]: #check for null values
pd.isnull(df).sum()
Out[14]:
User_ID      0
Cust_name    0
Product_ID   0
Gender       0
Age Group    0
Age          0
Marital_Status  0
State        0
Zone         0
Occupation   0
Product_Category  0
Orders       0
Amount      12
dtype: int64

In [15]: # drop null values
df.dropna(inplace=True)

In [16]: # change data type
df['Amount'] = df['Amount'].astype('int')

In [17]: df['Amount'].dtypes
Out[17]:
dtype('int32')

In [18]: df.columns
Out[18]:
Index(['User_ID', 'Cust_name', 'Product_ID', 'Gender', 'Age Group', 'Age', 'Marital_Status', 'State', 'Zone', 'Occupation', 'Product_Category', 'Orders', 'Amount'],
      dtype='object')

In [19]: #rename column
df.rename(columns={'Marital_Status':'Shadi'})
Out[19]:
   User_ID  Cust_name  Product_ID  Gender  Age Group  Age  Shadi  State  Zone  Occupation  Product_Category  Orders  Amount
0  1002903  Sanskriti  P0012942      F    26-35    28      0  Maharashtra  Western  Healthcare      Auto      1      23952
1  1000732  Karthik  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  Sudeshi  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  P0026942      M    18-25    19      1  Maharashtra  Western  Chemical      Office      4      370
11247  1004089  Reichanbach  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  Noonan  P00059442      M    36-45    37      0  Karnataka  Southern  Agriculture      Office      3      206
11250  1002744  Burnley  P00281742      F    18-25    19      0  Maharashtra  Western  Healthcare      Office      3      188
11239 rows x 13 columns

In [20]: # describe() method returns description of the data in the DataFrame (i.e. count, mean, std, etc)
df.describe()
Out[20]:
   User_ID      Age  Marital_Status  Orders  Amount
count  11239.000000  11239.000000  11239.000000  11239.000000
mean    35.410357    2.486334  9453.610553
std    12.763866    1.114967  5222.365168
min    12.000000    0.000000  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

In [21]: # use describe() for specific columns
df[['Age', 'Orders', 'Amount']].describe()
Out[21]:
   Age  Orders  Amount
count  11239.000000  11239.000000
mean    35.410357    2.486334  9453.610553
std    12.763866    1.114967  5222.365168
min    12.000000    0.000000  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 [22]: # 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)

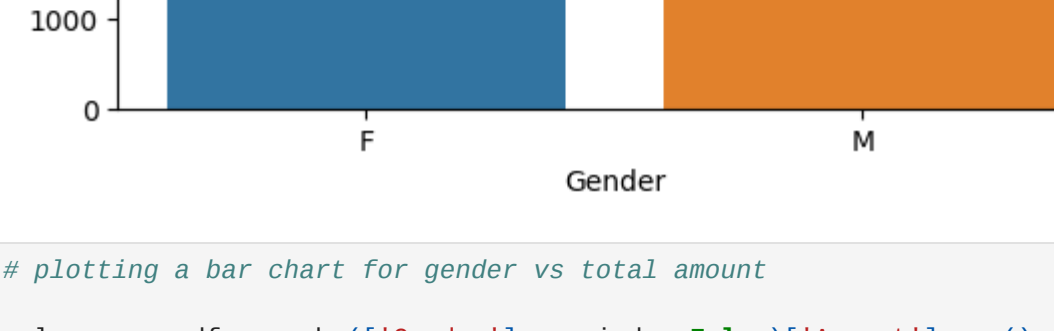
Out[22]:
8000
7000
6000
5000
4000
3000
2000
1000
0
F M
Gender
```

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 [23]: # 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[23]:
<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 [24]: ax = sns.countplot(data = df, x = 'Age Group', hue = 'Gender')

for bars in ax.containers:
    ax.bar_label(bars)

Out[24]:
3269
2772
102
134
1305
574
553
277
693
290
272
155
1578
705
26-35 0-17 18-25 51-55 46-50 55+ 36-45
Gender
Count
```

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

State

```
In [25]: # 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[25]:
<Axes: xlabel='Age Group', ylabel='Amount'>
```



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

State

```
In [26]: # 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[26]:
<Axes: xlabel='State', ylabel='Orders'>
```

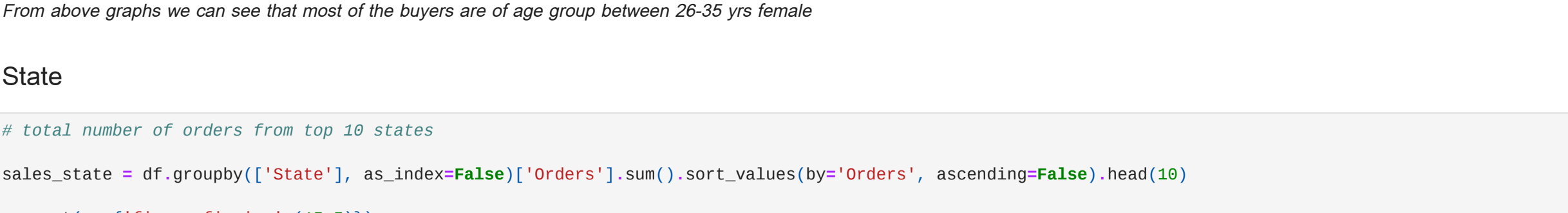


From above graphs we can see that most of the orders & total sales/amount are from Uttar Pradesh, Maharashtra and Karnataka respectively

State

```
In [27]: # 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':(15,5)})
sns.barplot(data = sales_state, x = 'State', y= 'Amount')

Out[27]:
<Axes: xlabel='State', ylabel='Amount'>
```



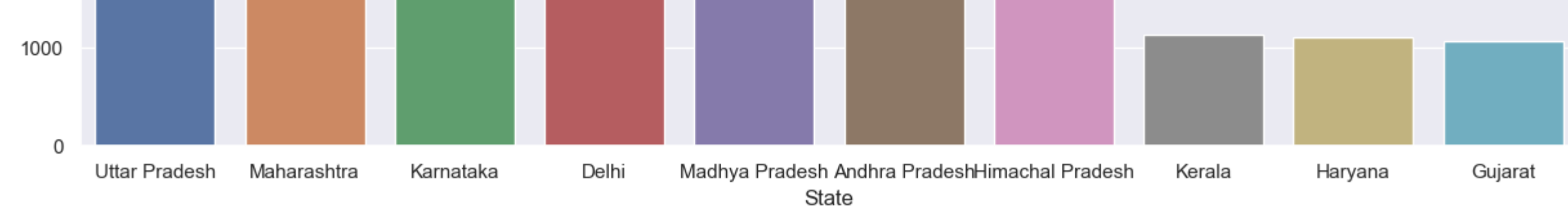
From above graphs we can see that most of the orders & total sales/amount are from Uttar Pradesh, Maharashtra and Karnataka respectively

Marital Status

```
In [28]: ax = sns.countplot(data = df, x = 'Marital_Status')

sns.set(rc={'figure.figsize':(7,5)})
for bars in ax.containers:
    ax.bar_label(bars)

Out[28]:
6518
4721
0 1
Marital_Status
```

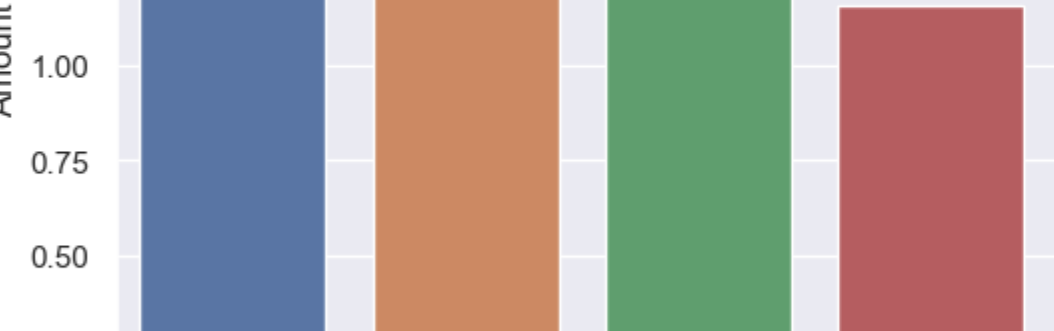


From above graphs we can see that most of the buyers are married (women) and they have high purchasing power

Occupation

```
In [29]: sales_state = df.groupby(['Marital_Status', 'Gender'], as_index=False)['Amount'].sum().sort_values(by='Amount', ascending=False)
sns.set(rc={'figure.figsize':(10,5)})
sns.barplot(data = sales_state, x = 'Marital_Status', y= 'Amount', hue='Gender')

Out[29]:
<Axes: xlabel='Marital_Status', ylabel='Amount'>
```



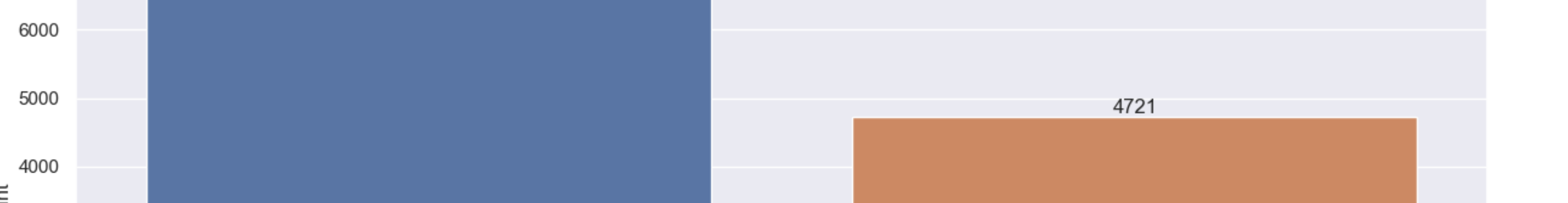
From above graphs we can see that most of the buyers are married (women) and they have high purchasing power

Occupation

```
In [30]: sns.set(rc={'figure.figsize':(20,5)})
ax = sns.countplot(data = df, x = 'Occupation')

for bars in ax.containers:
    ax.bar_label(bars)

Out[30]:
1498
854
565
414
423
531
637
1137
501
1583
1310
703
283
349
541
Healthcare Govt Automobile Construction Food Processing Lawyer Media Banking Occupation Retail IT Sector Aviation Hospitality Agriculture Textile Chemical
```



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

Product Category

```
In [31]: 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[31]:
<Axes: xlabel='Occupation', ylabel='Amount'>
```

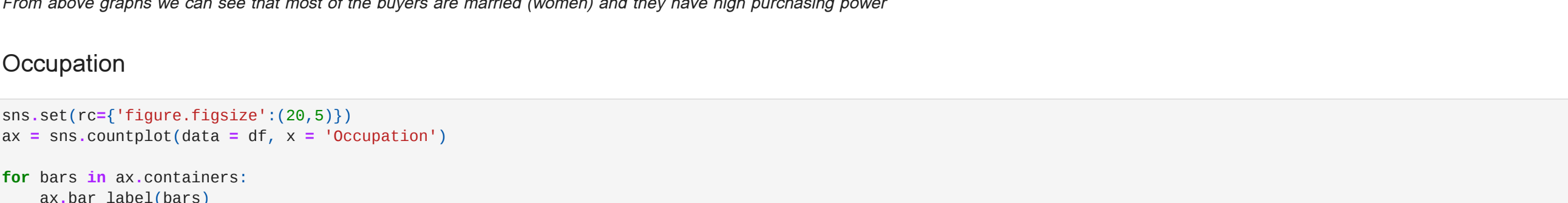


From above graphs we can see that most of the buyers are working in IT, Healthcare and Aviation sector

Product Category

```
In [32]: 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[32]:
<Axes: xlabel='Product_ID', ylabel='Orders'>
```

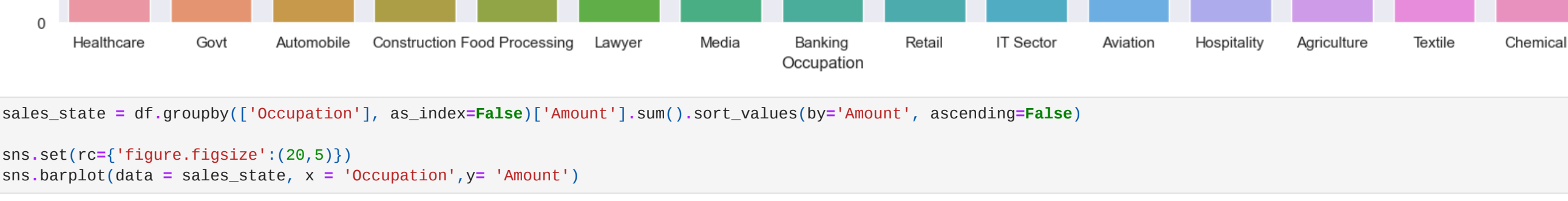


From above graphs we can see that most of the sold products are from Food, Clothing and Electronics category

Product Category

```
In [33]: 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[33]:
<Axes: xlabel='Product_ID', ylabel='Orders'>
```

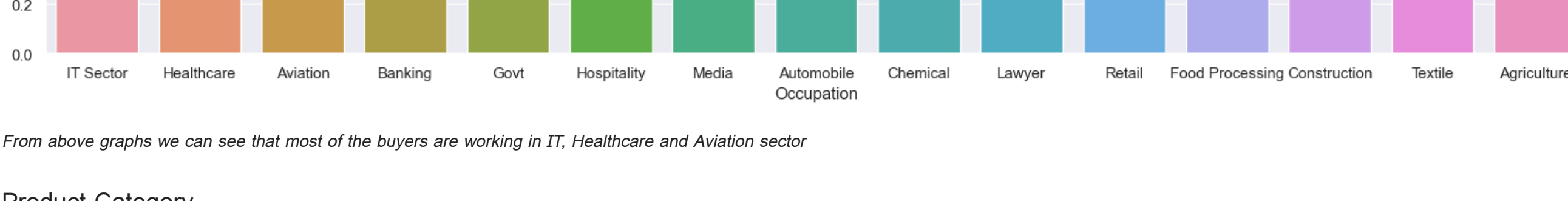


From above graphs we can see that most of the sold products are from Food, Clothing and Electronics category

Product Category

```
In [34]: # top 10 most sold products (same thing as above)
figt, ax1 = plt.subplots(figsize=(20,7))
df.groupby('Product_ID')['Orders'].sum().nlargest(10).sort_values(ascending=False).plot(kind='bar')

Out[34]:
<Axes: xlabel='Product_ID'>
```



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!