

```
In [12]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [13]: df=pd.read_csv('https://raw.githubusercontent.com/rishabhmisra/Python_Diwali_Sales_Analysis/refs/heads/main/Diwali20Sales20Data.csv',encoding='unicode_escape')
```

```
In [14]: df.shape
Out[14]: (11251, 15)
```

0

1002903

Sanskriti

P00125942

F

26-35

28

0

Maharashtra

Western

Healthcare

Auto

1

23952.0

NaN

NaN

1

1000732

Kartik

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

Sudevi

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

23977.0

NaN

NaN

In [31]:

df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 11251 entries, 0 to 11250
Data columns (total 15 columns):
 #   Column                Non-Null Count  Dtype
---  -
```

```
In [13]: df.info()
<class 'pandas.core.frame.DataFrame'>
Int64Index: 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  object
 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
dtypes: float64(1), int64(4), object(8)
memory usage: 3.1+ MB
```

```
In [16]: pd.isnull(df)
```

```
Out[16]:
```

User_ID	Cust_Name	Product_ID	Gender	Age Group	Age	Marital_Status	State	Zone	Occupation	Product_Category	Orders	Amount
0	False	False	False	False	False	False	False	False	False	False	False	False
1	False	False	False	False	False	False	False	False	False	False	False	False
2	False	False	False	False	False	False	False	False	False	False	False	False
3	False	False	False	False	False	False	False	False	False	False	False	False
4	False	False	False	False	False	False	False	False	False	False	False	False
...
11246	False	False	False	False	False	False	False	False	False	False	False	False
11247	False	False	False	False	False	False	False	False	False	False	False	False
11248	False	False	False	False	False	False	False	False	False	False	False	False
11249	False	False	False	False	False	False	False	False	False	False	False	False
11250	False	False	False	False	False	False	False	False	False	False	False	False

```
In [12]: df.columns
Out[12]: Index(['User_ID', 'Cust_name', 'Product_ID', 'Gender', 'Age Group', 'Age', 'Marital_Status', 'State', 'Zone', 'Occupation', 'Product_Category', 'Orders', 'Amount', 'unnamed1'], dtype='object')
```

```
In [14]: # check for null values
pd.isnull(df.amm)
```

```
Out[14]:
```

User_ID	Cust_name	Product_ID	Gender	Age Group	Age	Marital_Status	State	Zone	Occupation	Product_Category	Orders	Amount
0	False	0	0	0	0	0	0	0	0	0	0	0
1	False	0	0	0	0	0	0	0	0	0	0	0
2	False	0	0	0	0	0	0	0	0	0	0	0
3	False	0	0	0	0	0	0	0	0	0	0	0
4	False	0	0	0	0	0	0	0	0	0	0	0
...
11246	False	0	0	0	0	0	0	0	0	0	0	0
11247	False	0	0	0	0	0	0	0	0	0	0	0
11248	False	0	0	0	0	0	0	0	0	0	0	0
11249	False	0	0	0	0	0	0	0	0	0	0	0
11250	False	0	0	0	0	0	0	0	0	0	0	0

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

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

```
In [16]: df.shape
Out[16]: (11239, 13)
```

```
In [18]: #initiallize list of state
data_test=[['madhav',11],['Gopi',15],['Keshav'],['Laila',16]]
```

```
#Create the pandas DataFrame using list
df_test=pd.DataFrame(data_test,columns=['Name','Age'])
df_test
```

```
Out[18]:
```

Name	Age
0	madhav 11.0
1	Gopi 15.0
2	Keshav NaN
3	Laila 16.0

```
In [19]: df_test.dropna()
```

```
Out[19]:
```

Name	Age
0	madhav 11.0
1	Gopi 15.0
3	Laila 16.0

```
In [20]: df_test
```

```
Out[20]:
```

Name	Age
0	madhav 11.0
1	Gopi 15.0
2	Keshav NaN
3	Laila 16.0

```
In [21]: #change data type
df['Amount']=df['Amount'].astype('int')
```

```
In [22]: df['Amount'].dtypes
Out[22]: dtype('int32')
```

```
In [23]: #rename columns
df.rename(columns={'Marital_Status':'Shadi'})
```

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	P00236942	M	18-25	19	1	Maharashtra	Western	Chemical	Office	4	370
11247	1004029	Raichenbach	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	Nonnan	P00059442	M	36-45	37	0	Karnataka	Southern	Agriculture	Office	3	206

```
In [24]: df.describe()
```

```
Out[24]:
```

	User_ID	Cust_Name	Product_ID	Gender	Age Group	Age	Marital_Status	Orders	Amount
count	1123900e+04	11239.000000	11239.000000	11239.000000	11239.000000	11239.000000	11239.000000	11239.000000	11239.000000
mean	1.003004e+06	35.410357	0.420055	2.489634	9453.610563				
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				

```
Out[25]:
```

	Age	Orders	Amount
count	11239.000000	11239.000000	11239.000000
mean	35.410357	2.489634	9453.610563
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 countplot



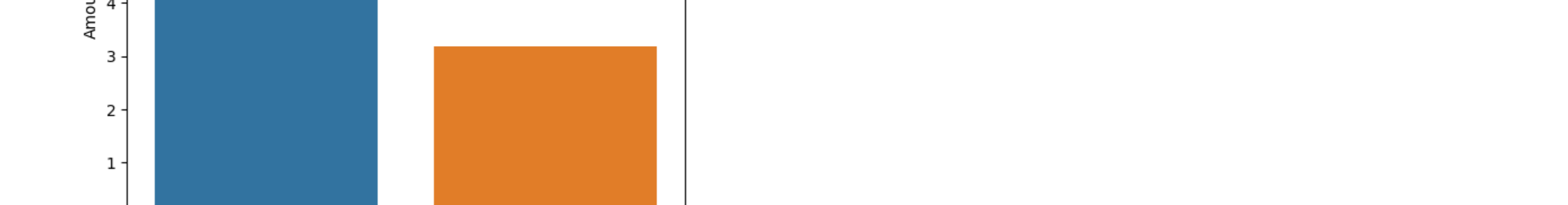
```
In [27]: sales_gender.groupby(['Gender'],as_index=False)['Amount'].sum().sort_values(by='Amount',ascending=False)
sns.barplot(x='Gender',y='Amount',data=sales_gender)
```



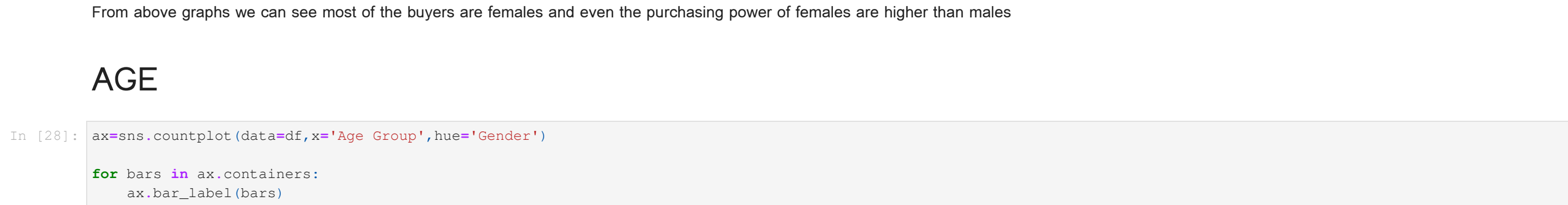
From above graphs we can see most of the buyers are females and even the purchasing power of females are higher than males

AGE

```
In [28]: ax=sns.countplot(data=df,x='Age Group',hue='Gender')
for bars in ax.containers:
    ax.bar_label(bars)
```



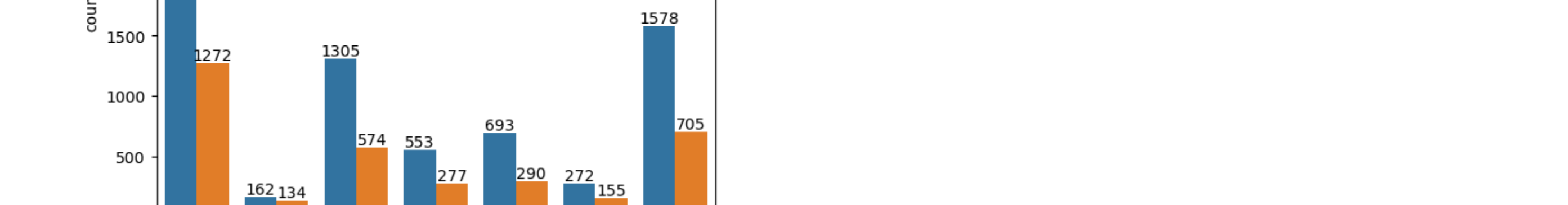
```
In [29]: # Total Amount vs Age Group
sales_age=df.groupby(['Age Group'],as_index=False)['Amount'].sum().sort_values(by='Amount',ascending=False)
sns.barplot(x='Age Group',y='Amount',data=sales_age)
```



From above we can see that most buyers are between in range of 26-35

State

```
In [30]: #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.barplot(data=sales_state,x='State',y='Orders')
```



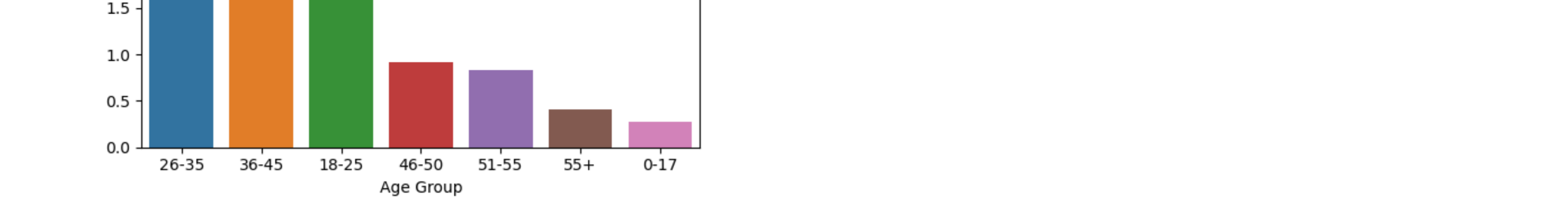
```
In [31]: # total sales from top 10 states
sales_state=df.groupby(['State'],as_index=False)['Amount'].sum().sort_values(by='Amount',ascending=False).head(10)
sns.set(rc={'figure.figsize':(15,5)})
sns.barplot(data=sales_state,x='State',y='Amount')
```



From above graphs we can see that most orders are from Uttar Pradesh,Maharashtra,Karnataka respectively but total sales/amount is from UP,Karnataka then Maharashtra

Marital Status

```
In [40]: ax=sns.countplot(data=df, x='Marital_Status')
sns.set(rc={'figure.figsize':(6,5)})
for bars in ax.containers:
    ax.bar_label(bars)
```



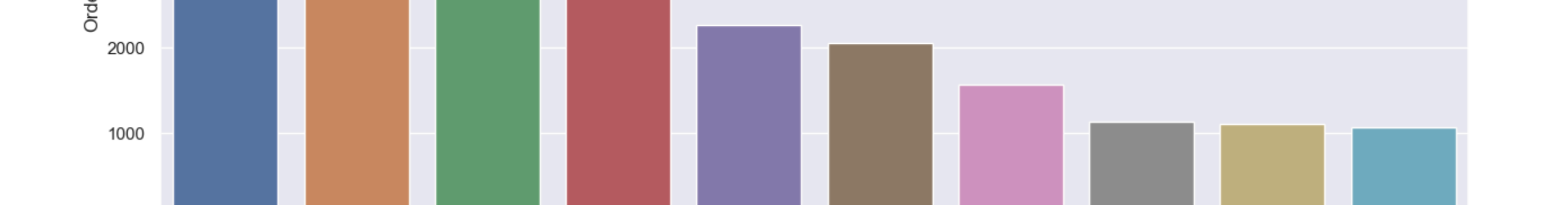
```
In [41]: sales_state=df.groupby(['Marital_Status','Gender'],as_index=False)['Amount'].sum().sort_values(by='Amount',ascending=False)
sns.set(rc={'figure.figsize':(6,5)})
sns.barplot(data=sales_state,x='Marital_Status',y='Amount')
```



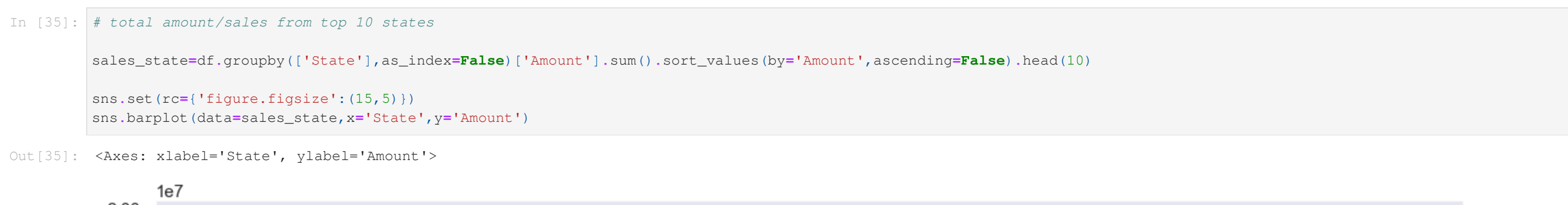
From above graphs we can see that most of the buyers are married (women) and they have high purchasing power

Occupation

```
In [42]: sns.set(rc={'figure.figsize':(20,5)})
ax = sns.countplot(data = df, x = "Occupation")
for bars in ax.containers:
    ax.bar_label(bars)
```



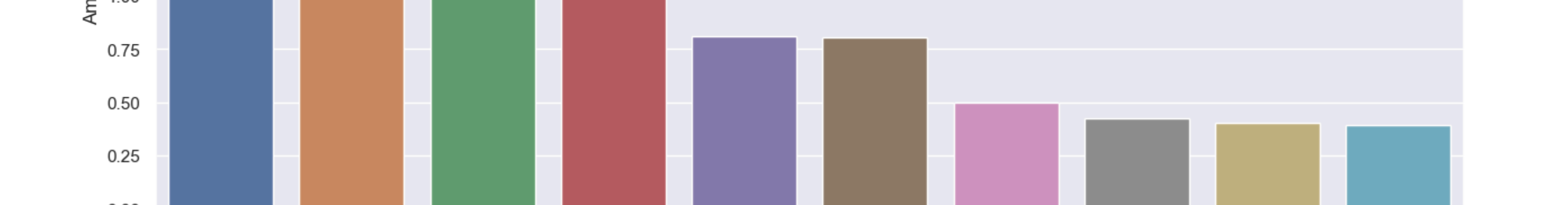
```
In [43]: sales_state = df.groupby(['Occupation'], as_index=False)['Amount'].sum().sort_values(by='Amount', ascending=False)
sns.set(rc={'figure.figsize':(20,5)})
sns.set(rc={'figure.figsize':(10,5)})
sns.barplot(data = sales_state, x = "Occupation", y = "Amount")
```



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

Product Category

```
In [44]: 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 [45]: 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")
```



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

```
In [46]: sales_state = df.groupby(['Product_ID'], as_index=False)['Orders'].sum().sort_values(by='Orders', ascending=False).head(10)
df.groupby(['Product_ID','Orders']).sum().nlargest(10).sort_values(ascending=False).plot(kind='bar')
```

```
Out[46]: <Axes: xlabel='Product_ID', ylabel='Orders'>
```



top 10 most sold products (same thing as above)

```
fig1, ax1 = plt.subplots(figsize=(12,7))
df.groupby(['Product_ID','Orders']).sum().nlargest(10).sort_values(ascending=False).plot(kind='bar')
```

```
Out[47]: <Axes: xlabel='Product_ID'>
```

