

```
In [1]: # import python libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [2]: # import csv file
df = pd.read_csv('Biswall Sales Data.csv', encoding= 'unicode_escape')
```

```
In [3]: df.shape
Out[3]: (11251, 15)
```

```
In [4]: df.head()
```

Out[4]:

	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	23952.0	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	P00227842	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 [5]: df.tail()
```

	User_ID	Cust_name	Product_ID	Gender	Age Group	Age	Marital_Status	State	Zone	Occupation	Product_Category	Orders	Amount	Status	unnamed1	
	11246	1000695	Manning	P00296942	M	18-25	19	1	Maharashtra	Western	Chemical	Office	4	370.0	NaN	NaN
	11247	1004089	Reichenbach	P00171342	M	26-35	33	0	Haryana	Northern	Healthcare	Veterinary	3	367.0	NaN	NaN
	11248	1001209	Oshin	P00201342	F	36-45	40	0	Madhya Pradesh	Central	Textile	Office	4	213.0	NaN	NaN
	11249	1004023	Noonan	P00059442	M	36-45	37	0	Karnataka	Southern	Agriculture	Office	3	206.0	NaN	NaN
	11250	1002744	Bumley	P00281742	F	18-25	19	0	Maharashtra	Western	Healthcare	Office	3	188.0	NaN	NaN

```
In [6]: 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 [7]: #drop unrelated/blank columns
df.drop(['Status', 'unnamed1'], axis=1, inplace=True)
```

```
In [8]: pd.isnull(df).sum()
```

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 [9]: df.dropna(inplace=True)
```

```
In [10]: pd.isnull(df).sum()
```

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	0
dtype:	int64

```
In [11]: # change data type
df['Amount'] = df['Amount'].astype('int')
```

```
In [12]: df['Amount'].dtypes
Out[12]: dtype('int32')
```

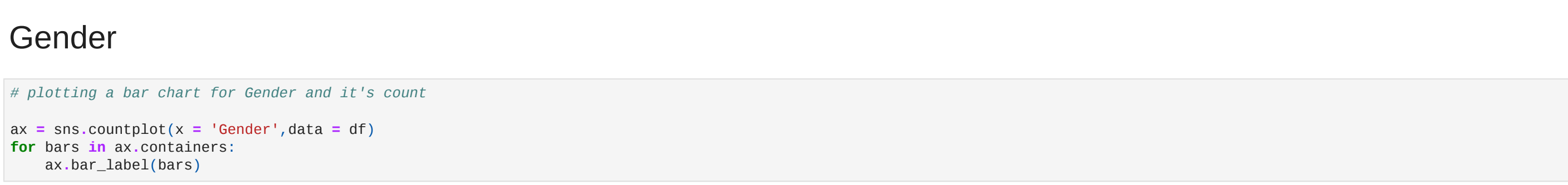
```
In [13]: df.columns
```

```
Out[13]: Index(['User_ID', 'Cust_name', 'Product_ID', 'Gender', 'Age_Group', 'Age', 'Marital_Status', 'State', 'Zone', 'Occupation', 'Product_Category', 'Orders', 'Amount'], dtype='object')
```

## Data Analysis

### Gender

```
In [14]: # plotting a bar chart for Gender and it's count
```



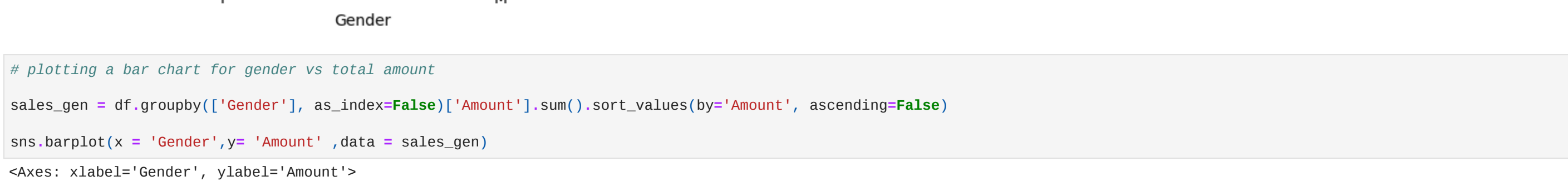
```
In [15]: # plotting a bar chart for gender vs total amount
```



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 [16]: ax = sns.countplot(data = df, x = 'Age_Group', hue = 'Gender')
```

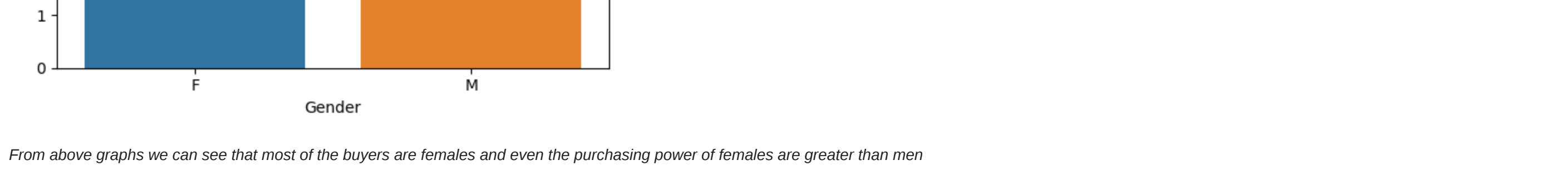


```
In [17]: # Total Amount vs Age_Group
sales_age = df.groupby(['Age_Group'], as_index=False)['Amount'].sum().sort_values(by='Amount', ascending=False)
```

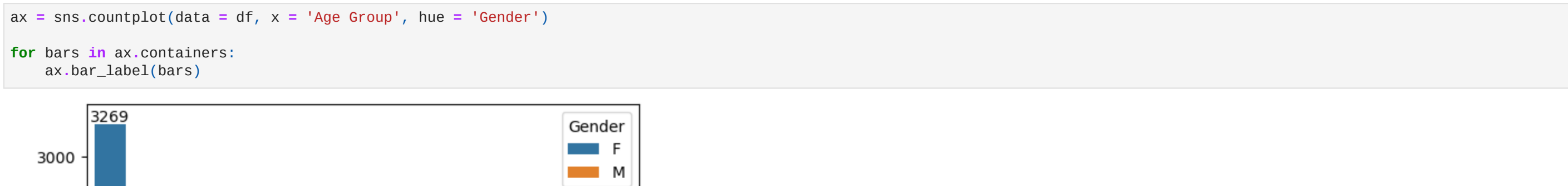


From above graphs we can see that most of the buyers are of age group between 26-35 yrs female

```
In [18]: # total number of orders from top 5 states
sales_state = df.groupby(['State'], as_index=False)['Orders'].sum().sort_values(by='Orders', ascending=False).head(5)
```



```
In [19]: # total amount/sales from top 5 states
sales_state = df.groupby(['State'], as_index=False)['Amount'].sum().sort_values(by='Amount', ascending=False).head(5)
```



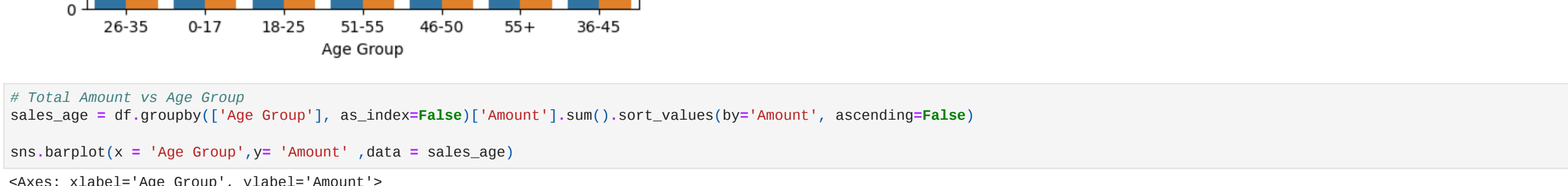
From above graphs we can see that most of the orders & total sales/amount are from Uttar Pradesh, Maharashtra and Karnataka respectively

```
In [20]: df.columns
```

```
Out[20]: Index(['User_ID', 'Cust_name', 'Product_ID', 'Gender', 'Age_Group', 'Age', 'Marital_Status', 'State', 'Zone', 'Occupation', 'Product_Category', 'Orders', 'Amount'], dtype='object')
```

### Marital Status

```
In [21]: ax = sns.countplot(data = df, x = 'Marital_Status')
```



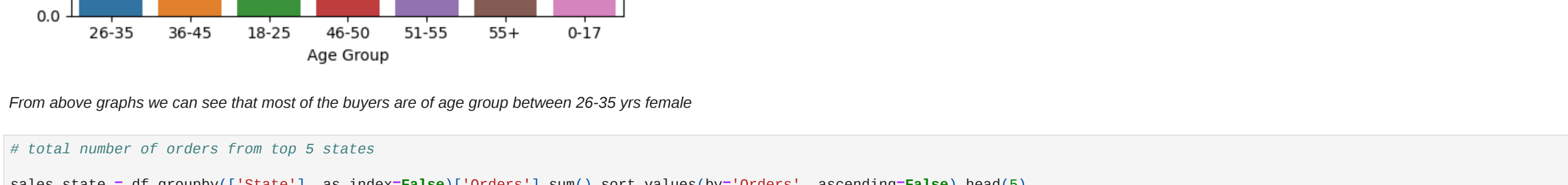
```
In [22]: sales_state = df.groupby(['Marital_Status', 'Gender'], as_index=False)['Amount'].sum().sort_values(by='Amount', ascending=False)
```



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':(15,5)})
ax = sns.countplot(data = df, x = 'Occupation')
```



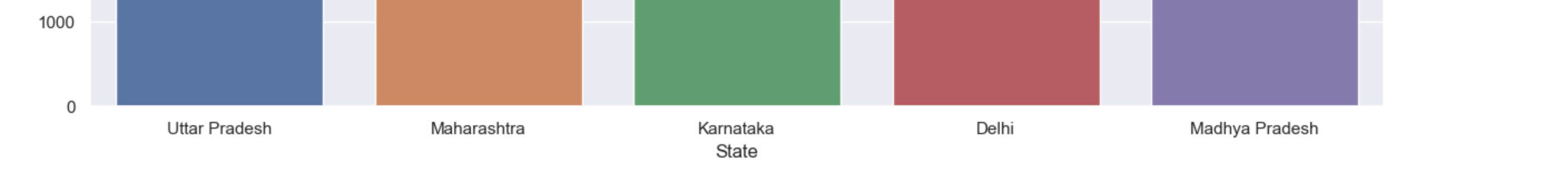
```
In [24]: sales_state = df.groupby(['Occupation'], as_index=False)['Amount'].sum().sort_values(by='Amount', ascending=False)
```



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

### Product Category

```
In [25]: sns.set(rc={'figure.figsize':(25,9)})
ax = sns.countplot(data = df, x = 'Product_Category')
```

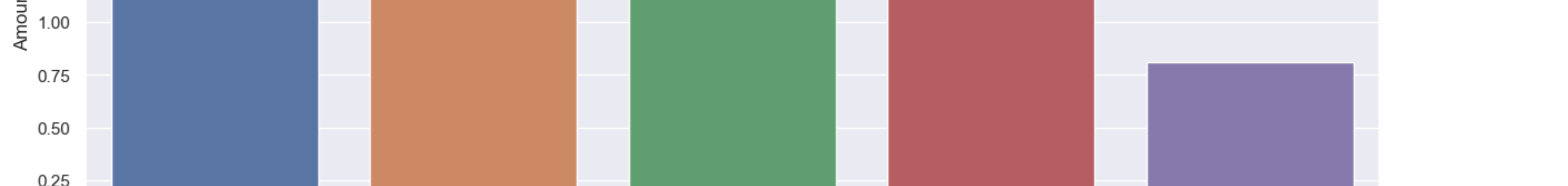


```
In [26]: sales_state = df.groupby(['Product_Category'], as_index=False)['Amount'].sum().sort_values(by='Amount', ascending=False).head(10)
```

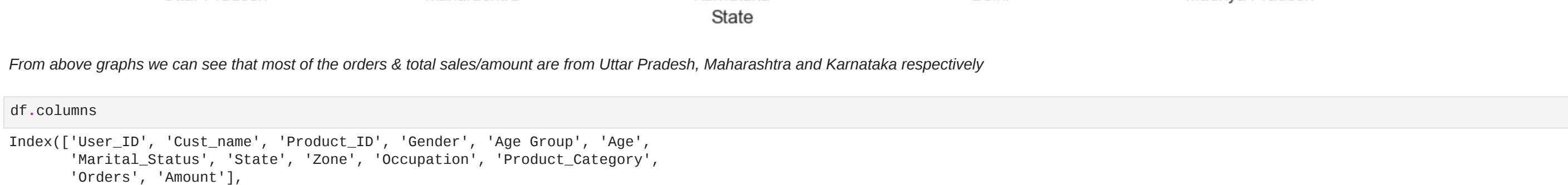


From above graphs we can see that most of the sold products are from Food, Clothing and Electronics category

```
In [27]: sales_state = df.groupby(['Product_ID'], as_index=False)['Orders'].sum().sort_values(by='Orders', ascending=False).head(5)
```



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



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

```
In [ ]:
```