

FileEditViewKernelHelpSettingsHelp

JupyterLabPython 3 (ipykernel)

[3]:

import python libraries

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

[4]:

import csv file
df = pd.read_csv('Diwali Sales Data.csv', encoding= 'unicode_escape')

[5]:

df.shape

[5]:

(11251, 15)

[7]:

df.head(5)

[7]:

	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	
1	1000732	Kartik	P00110942	F	26-35	35	1	Andhra Pradesh	Southern	Govt	Auto	3	23934.0	NaN	
2	1001990	Bindu	P00118542	F	26-35	35	1	Uttar Pradesh	Central	Automobile	Auto	3	23924.0	NaN	
3	1001425	Sudevi	P00237842	M	0-17	16	0	Karnataka	Southern	Construction	Auto	2	23912.0	NaN	
4	1000588	Jani	P00057942	M	26-35	28	1	Gujarat	Western	Food Processing	Auto	2	23877.0	NaN	

[8]:

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

[9]:

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

[10]:

df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 11251 entries, 0 to 11250
Data columns (total 13 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

[10]:

df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 11251 entries, 0 to 11250
Data columns (total 13 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
dtypes: float64(1), int64(4), object(8)
memory usage: 1.1+ MB

[11]:

pd.isnull(df).sum()

[11]:

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

[12]:

df.dropna(inplace=True)

[13]:

<class 'pandas.core.frame.DataFrame'>
Index: 11239 entries, 0 to 11250
Data columns (total 13 columns):
Column Non-Null Count Dtype

0 User_ID 11239 non-null int64
1 Cust_name 11239 non-null object
2 Product_ID 11239 non-null object
3 Gender 11239 non-null object
4 Age Group 11239 non-null object
5 Age 11239 non-null int64
6 Marital_Status 11239 non-null int64
7 State 11239 non-null object
8 Zone 11239 non-null object
9 Occupation 11239 non-null object
10 Product_Category 11239 non-null object
11 Orders 11239 non-null int64
12 Amount 11239 non-null float64
dtypes: float64(1), int64(4), object(8)
memory usage: 1.2+ MB

[15]:

change data type
df['Amount'] = df['Amount'].astype('int')

[16]:

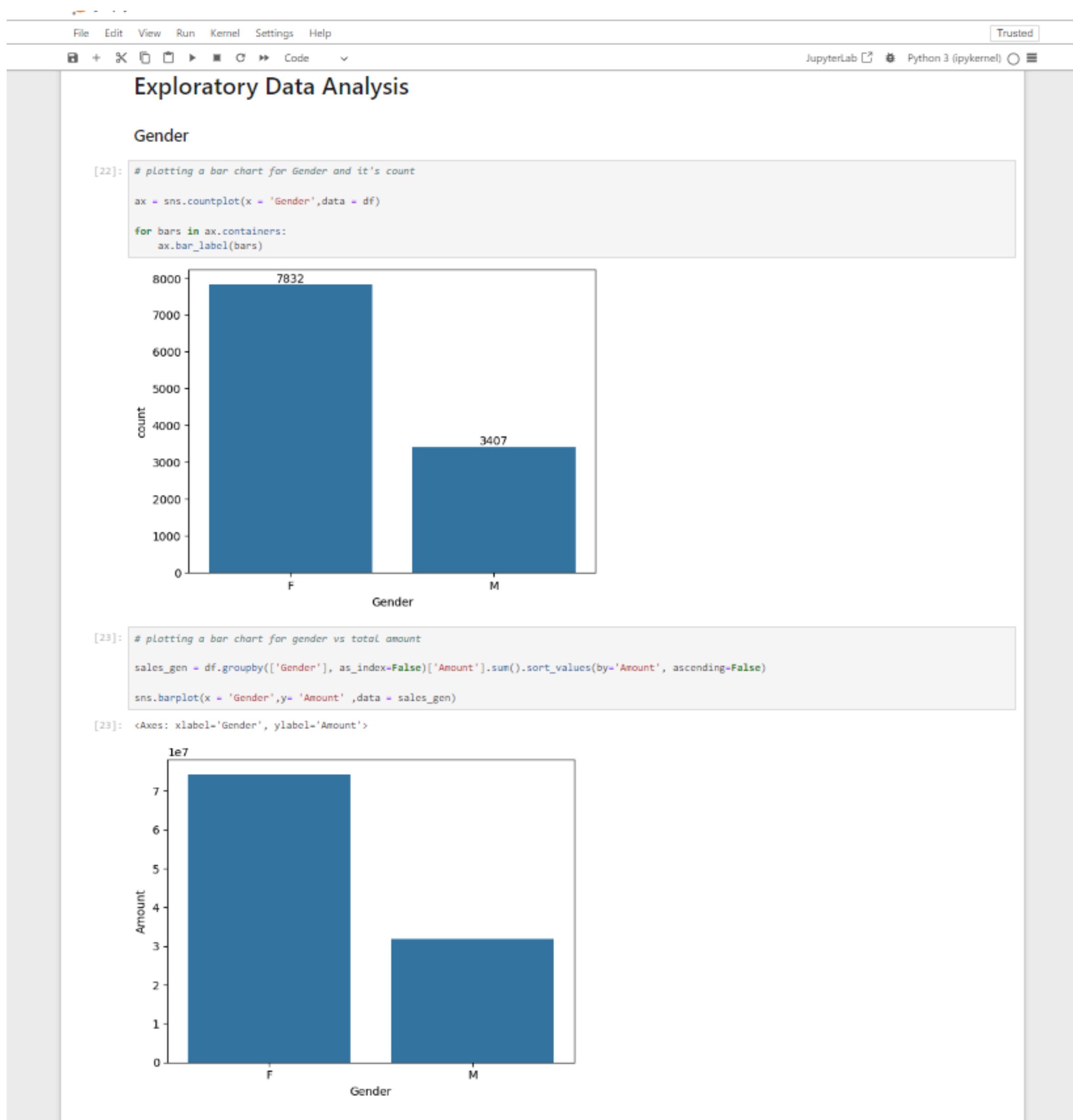
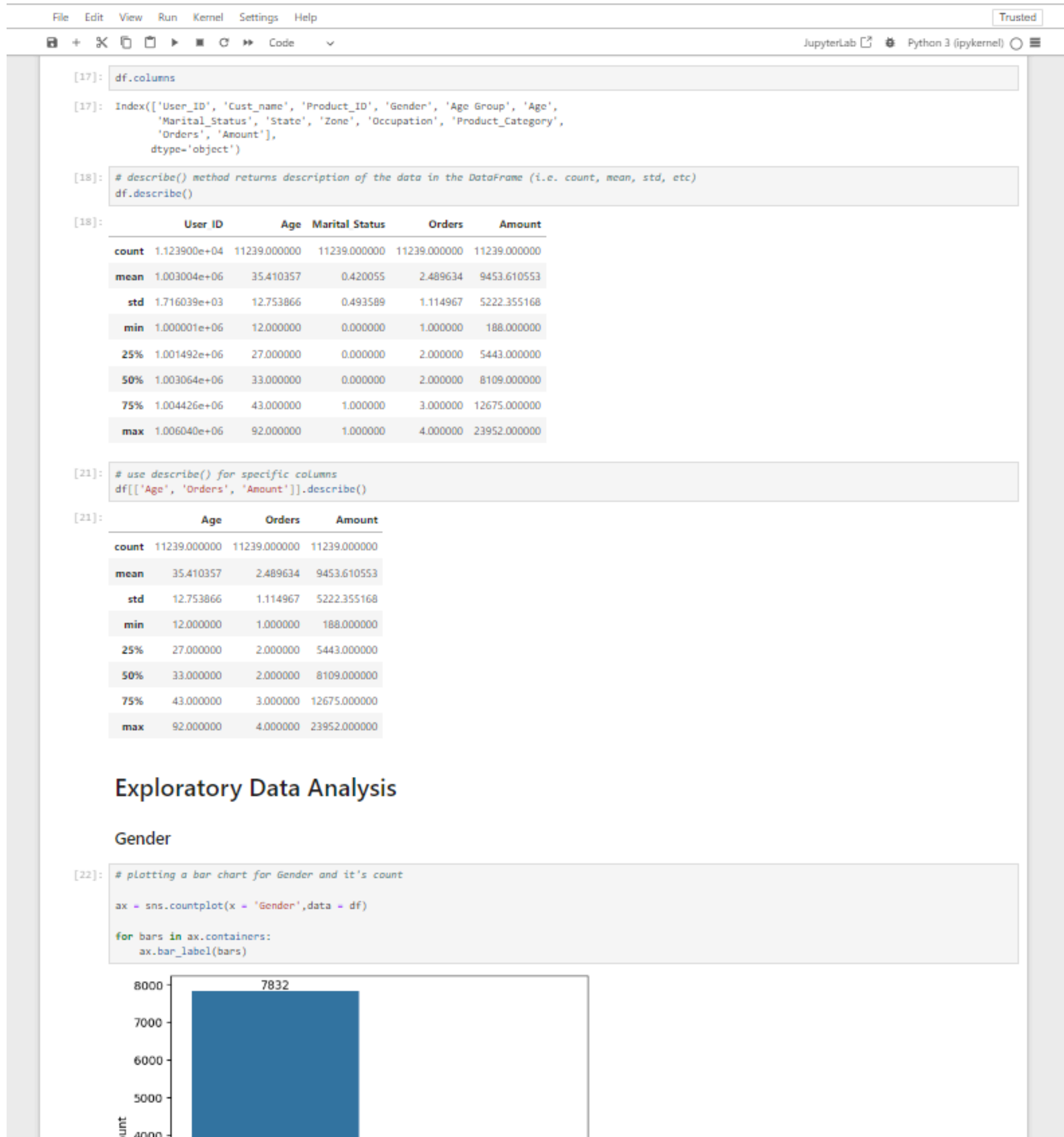
df['Amount'].dtypes

[16]:

dtype('int64')

[17]:

df.columns

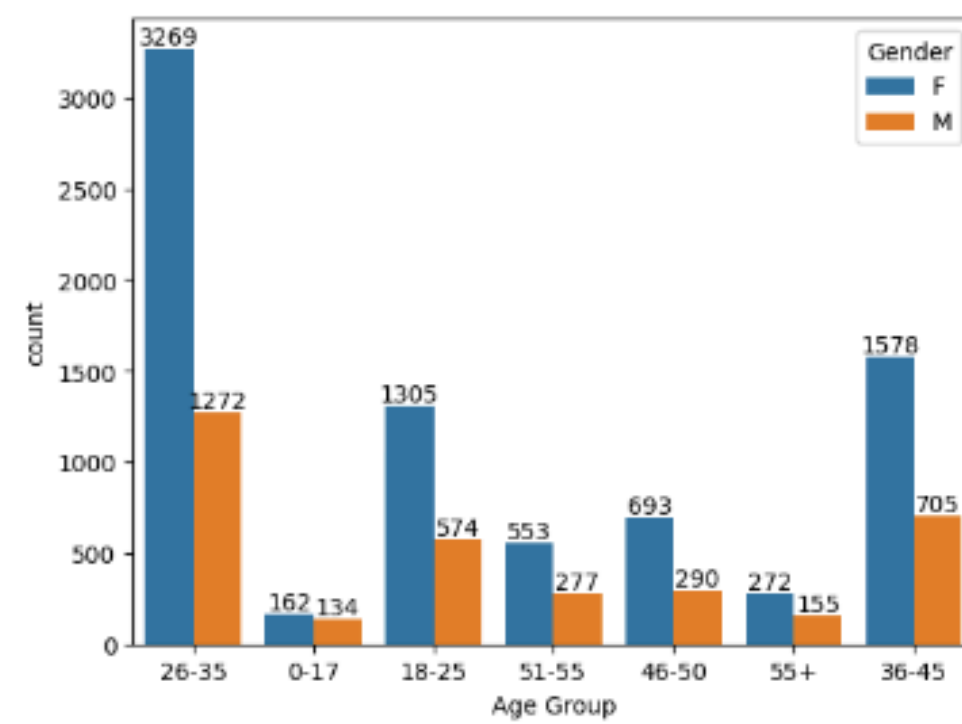


From above graphs we can see that most of the buyers are females and even the purchasing power of females are greater than men

Age

```
[26]: ax = sns.countplot(data = df, x = 'Age Group', hue = 'Gender')

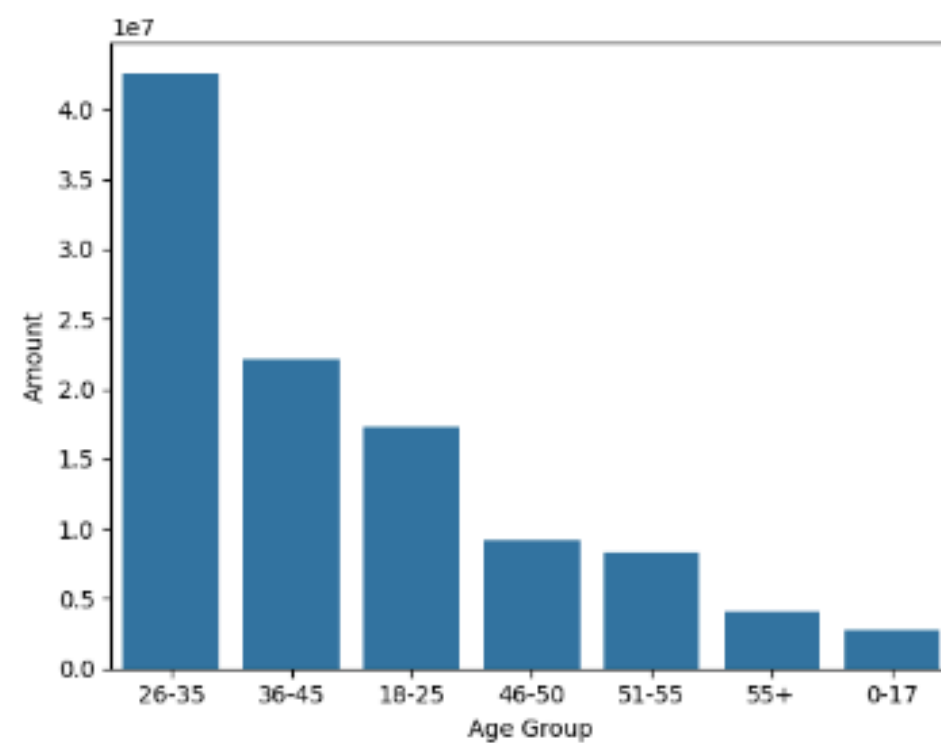
for bars in ax.containers:
    ax.bar_label(bars)
```



```
[27]: # 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)
```

```
[27]: <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

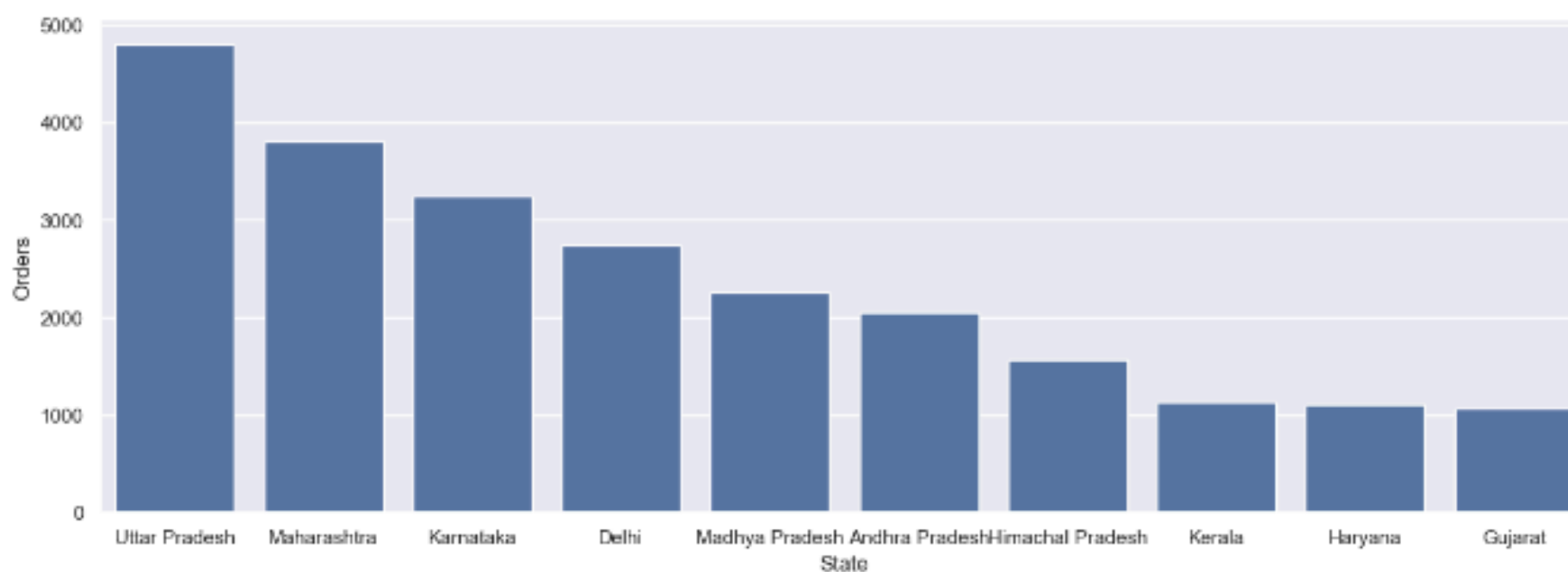
State

```
[29]: # 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')
```

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

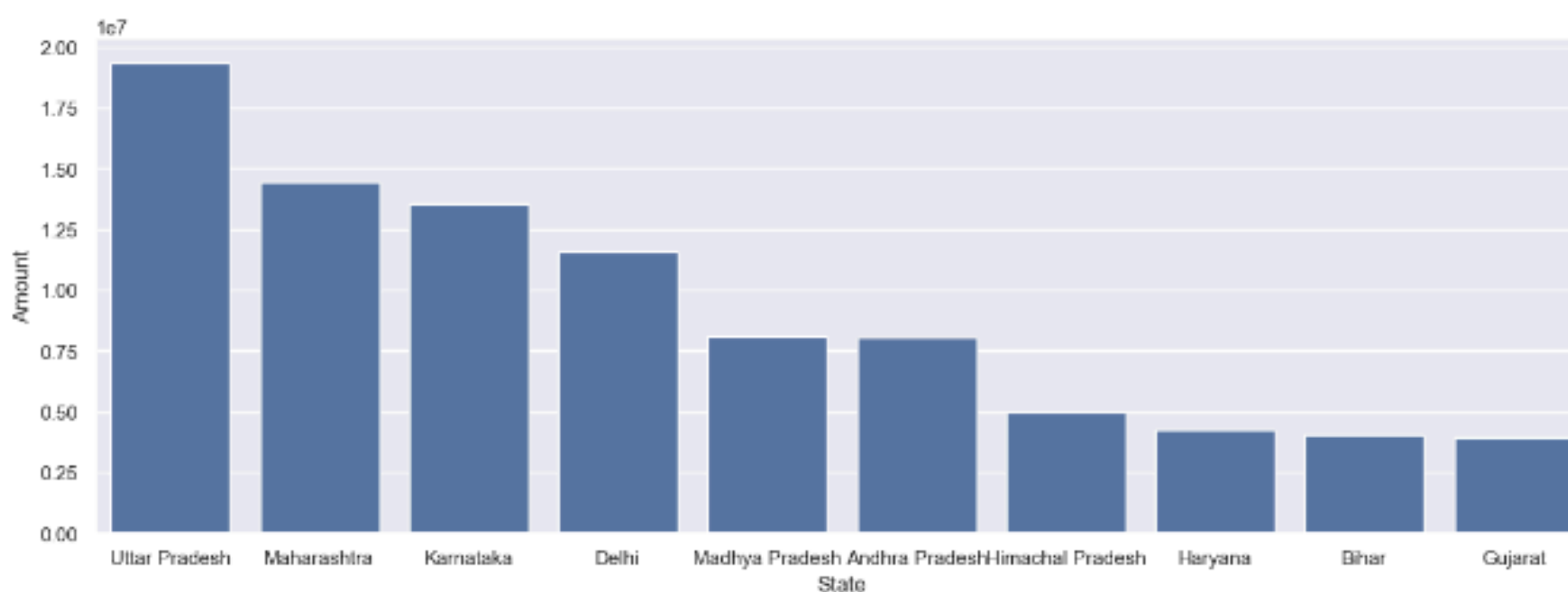


```
[30]: # total amount/sales from top 10 states

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')
```

```
[30]: <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

