

Exploratory Data Analysis of Diwali Sales | Python Project-2

Objective: Analysing Diwali sales data to improve customer experience and sales by making use of NumPy, Pandas, Matplotlib and Seaborn libraries

Step-1: Importing all the necessary libraries, uploading the csv file and checking the file contents

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

In [2]: #Reading the Diwali_Sales CSV file and assigning it to a DataFrame
df = pd.read_csv("D:\VProject\Portfolio\Python\Diwali Sales Analysis_Project-2\Diwali Sales Data.csv", encoding = 'unicode_escape')

In [3]: #Checking the total number of rows and columns
df.shape

Out[3]:
(11251, 15)

In [4]: #Checking all the columns and first 5 rows
df.head()

Out[4]:
   User_ID  Cust_name  Product_ID  Gender  Age  Marital_Status  State  Zone  Occupation  Product_Category  Orders  Amount  Status  unnamed1
0  1000903  Sankriti  P00125942  F  26-35  28  0  Maharashtra  Western  Healthcare  Auto  1  23962.0  NaN  NaN
1  1000732  Karik  P00110942  F  26-35  35  1  Andhra Pradesh  Southern  Govt  Auto  3  23934.0  NaN  NaN
2  1001090  Bindu  P00118542  F  26-35  35  1  Uttar Pradesh  Central  Automobile  Auto  3  23924.0  NaN  NaN
3  1001425  Sudini  P00217842  M  0-17  16  0  Karnataka  Southern  Construction  Auto  2  23912.0  NaN  NaN
4  1000588  Joni  P00057842  M  26-35  28  1  Gujarat  Western  Food Processing  Auto  2  23877.0  NaN  NaN

In [5]: #Digging deeper to know more about the columns, data types etc
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  object
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 [6]: #Fetching all the column names
df.columns

Out[6]:
Index(['User_ID', 'Cust_name', 'Product_ID', 'Gender', 'Age Group', 'Age',
       'Marital_Status', 'State', 'Zone', 'Occupation', 'Product_Category',
       'Orders', 'Amount', 'Status', 'unnamed1'],
      dtype='object')

In [7]: #Checking the rows/columns count
df.size

Out[7]:
(11251, 15)
```

Step-2: Data cleaning & Exploring the data using Pandas and NumPy

```
In [8]: #Dropping the empty columns
df.drop(['Status', 'unnamed1'], axis=1, inplace=True)

In [9]: #Cross verifying the contents of the file
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  object
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

In [10]: #Checking for the null values in the file
df.isnull()

Out[10]:
   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  False
1  False  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  False
3  False  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  False
...
11247  False  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  False
11249  False  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  False

11251 rows x 13 columns

In [11]: #Checking the total count of null values in all the columns
df.isnull().sum()

Out[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

In [12]: df.shape

Out[12]:
(11251, 13)

In [13]: #Dropping the null values
df.dropna(inplace=True)

In [14]: #Cross verifying the total rows and columns
df.shape

Out[14]:
(11239, 13)

In [15]: #Converting the datatype for our convenience
df['Amount'] = df['Amount'].astype(int)

In [16]: #Cross verifying the conversion of datatype
df['Amount'].dtype

Out[16]:
dtype('int32')

In [18]: #Checking the basic statistics of required columns
df[['Age', 'Orders', 'Amount']].describe()

Out[18]:
   Age  Orders  Amount
count 11239.000000  11239.000000  11239.000000
mean    35.410357    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  16785.000000
max      92.000000   4.000000  23952.000000
```

Step-3: Exploratory Data Analysis using Matplotlib and Seaborn

```
In [29]: #Checking the column names for creating various charts
df.columns

Out[29]:
Index(['User_ID', 'Cust_name', 'Product_ID', 'Gender', 'Age Group', 'Age',
       'Marital_Status', 'State', 'Zone', 'Occupation', 'Product_Category',
       'Orders', 'Amount'],
      dtype='object')

Analysing the Gender data

In [57]: # Creating a count plot for Gender column
ax = sns.countplot(data=df, x='Gender')

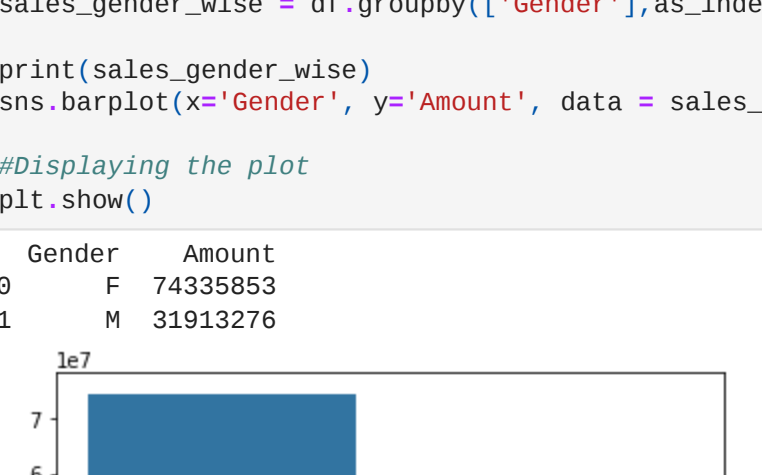
# Adding of labels and chart title
plt.xlabel('Gender')
plt.ylabel('Total Count')
plt.title('Total buyers in each gender')

#Inserting total value above the bars
for bars in ax.containers:
    ax.bar_label(bars)

#Displaying the plot
plt.show()

Total buyers in each gender

Gender
F 74335883
M 31913276
```

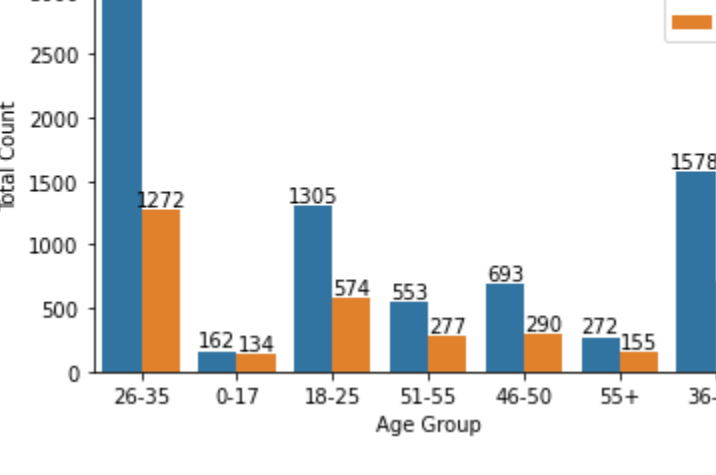


```
In [21]: #Grouping and aggregating values based on Gender
sales_gender_wise = df.groupby(['Gender'], as_index = False) ['Amount'].sum().sort_values(by='Amount', ascending=False)

print(sales_gender_wise)
sns.barplot(x='Gender', y='Amount', data = sales_gender_wise)

#Displaying the plot
plt.show()

Gender  Amount
0  F 74335883
1  M 31913276
```



Insight 1: From the above both graphs – we understand that females buyers are double the male buyers. And also females buyers spend double the amount compared to male buyers

Analysing the Age Group data

```
In [26]: # Creating a count plot for Age Group column
ax = sns.countplot(data=df, x='Age Group', hue = 'Gender')

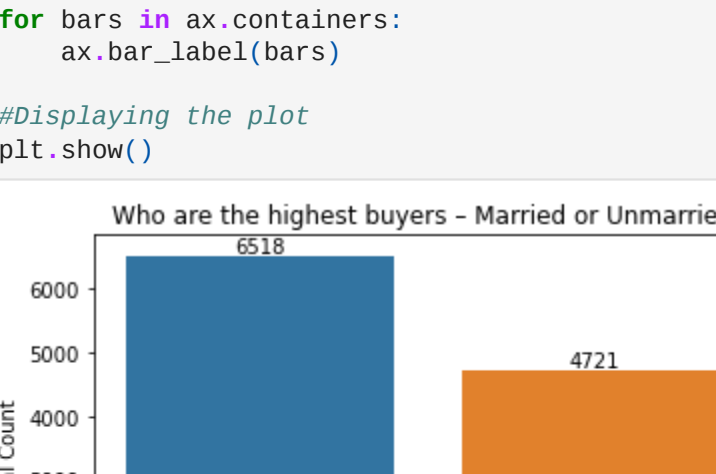
# Adding of labels and chart title
plt.xlabel('Age Group')
plt.ylabel('Total Count')
plt.title('Which Age group has highest buyers?')

#Inserting total value above the bars
for bars in ax.containers:
    ax.bar_label(bars)

#Displaying the plot
plt.show()

Which Age group has highest buyers?

Age Group  Amount
0  26-35  42621442
1  36-45  22144994
2  18-25  17240732
3  46-50  8201844
4  51-55  8261477
5  50+  4800987
6  0-17  2699653
```



Insight 2: Based on the above two charts – we understand that the 'Age group 26-35' contains highest number of buyers and especially female buyers are more compared to male buyers

Analysing the Marital Status data

```
In [56]: # Creating a count plot for Marital_Status column
ax = sns.countplot(data=df, x='Marital_Status')

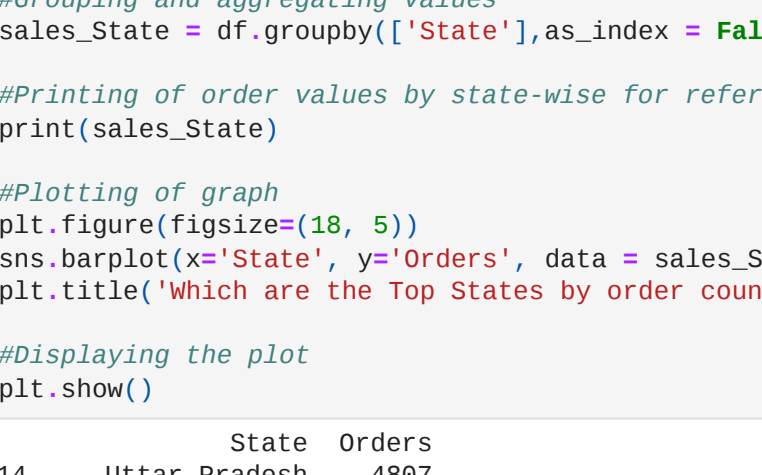
# Adding of labels and chart title
plt.xlabel('Marital_Status')
plt.ylabel('Total Count')
plt.title('Who are the highest buyers - Married or Unmarried?')

#Inserting total value above the bars
for bars in ax.containers:
    ax.bar_label(bars)

#Displaying the plot
plt.show()

Who are the highest buyers - Married or Unmarried?

Marital_Status  Amount
0  8518
1  4721
```



Insight 3: From the above graphs – we understand that the highest buyers are married people (especially women)

Analysing the State data

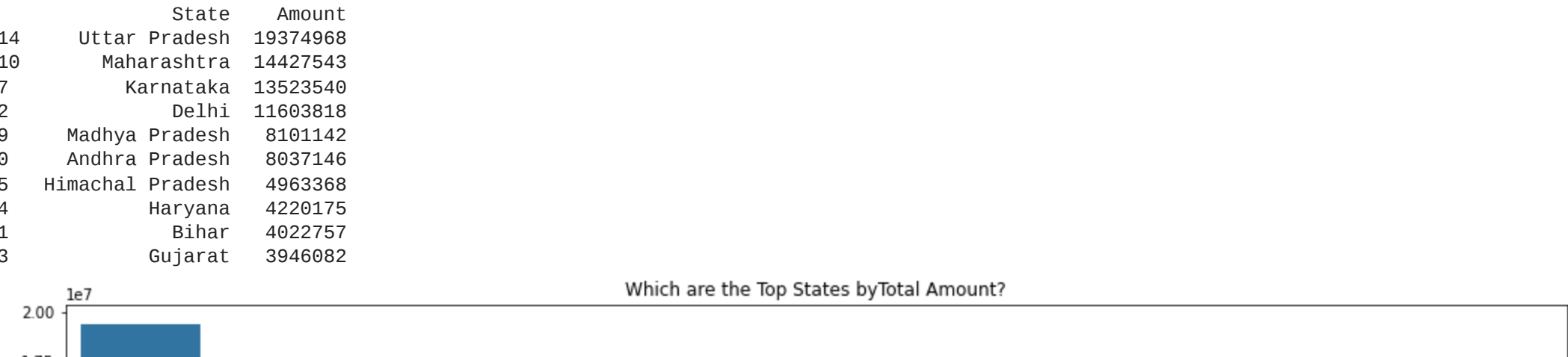
```
In [54]: #Checking top 10 states based on order count
#Grouping and aggregating values
sales_State = df.groupby(['State'], as_index = False) ['Orders'].sum().sort_values(by='Orders', ascending=False).head(10)

#Printing of order values by state-wise for reference
print(sales_State)

#Plotting of graph
plt.figure(figsize=(18, 5))
sns.barplot(x='State', y='Orders', data = sales_State)
plt.title('Which are the Top States by order count?')

#Displaying the plot
plt.show()

State  Orders
14  Uttar Pradesh  4807
10  Maharashtra  3810
7  Karnataka  3240
2  Delhi  2740
9  Madhya Pradesh  2552
0  Andhra Pradesh  2891
5  Himachal Pradesh  1565
9  Kerala  1137
4  Haryana  1189
3  Gujarat  1086
```



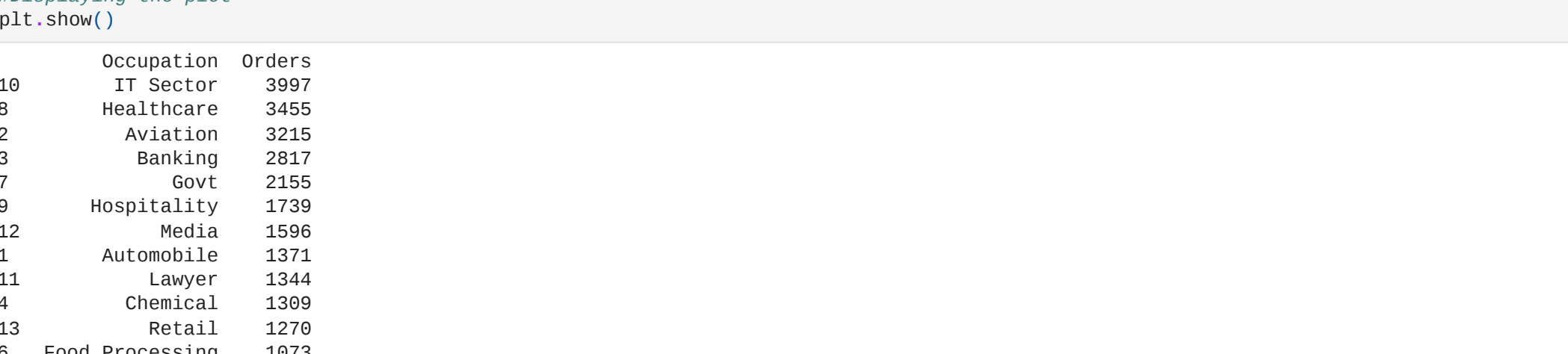
```
In [53]: #Checking top 10 states based on amount
#Grouping and aggregating values
sales_State = df.groupby(['State'], as_index = False) ['Amount'].sum().sort_values(by='Amount', ascending=False).head(10)

#Printing of order values by state-wise for reference
print(sales_State)

#Plotting of graph
plt.figure(figsize=(18, 5))
sns.barplot(x='State', y='Amount', data = sales_State)
plt.title('Which are the Top States by total Amount?')

#Displaying the plot
plt.show()

State  Amount
14  Uttar Pradesh  19374968
10  Maharashtra  14427542
7  Karnataka  13523548
2  Delhi  11683818
9  Madhya Pradesh  8181142
0  Andhra Pradesh  8937146
5  Himachal Pradesh  4963386
9  Haryana  4239175
1  Bihar  4822757
3  Gujarat  3946882
```



Insight 3: From the above graphs – we understand that Uttar Pradesh, Maharashtra and Karnataka states have highest order count and amount

Analysing the Occupation data

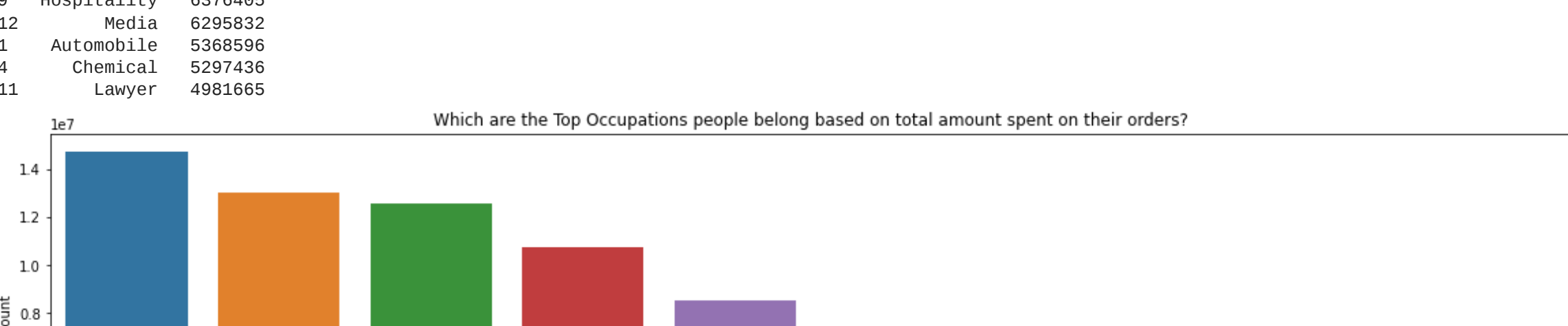
```
In [65]: #Checking top occupations people belong based on their order count
#Grouping and aggregating values
sales_Occupation = df.groupby(['Occupation'], as_index = False) ['Orders'].sum().sort_values(by='Orders', ascending=False)

#Printing of order values by state-wise for reference
print(sales_Occupation)

#Plotting of graph
plt.figure(figsize=(20, 5))
sns.barplot(x='Occupation', y='Orders', data = sales_Occupation)
plt.title('Which are the Top Occupations people belong by their order count?')

#Displaying the plot
plt.show()

Occupation  Orders
10  IT Sector  3997
8  Healthcare  3455
2  Aviation  2125
3  Banking  2817
7  Govt  2355
9  Hospitality  1739
12  Automobile  1596
1  Lawyer  1373
4  Chemical  1344
13  Retail  1270
6  Food Processing  1073
5  Construction  1025
14  Textile  893
0  Agriculture  722
```



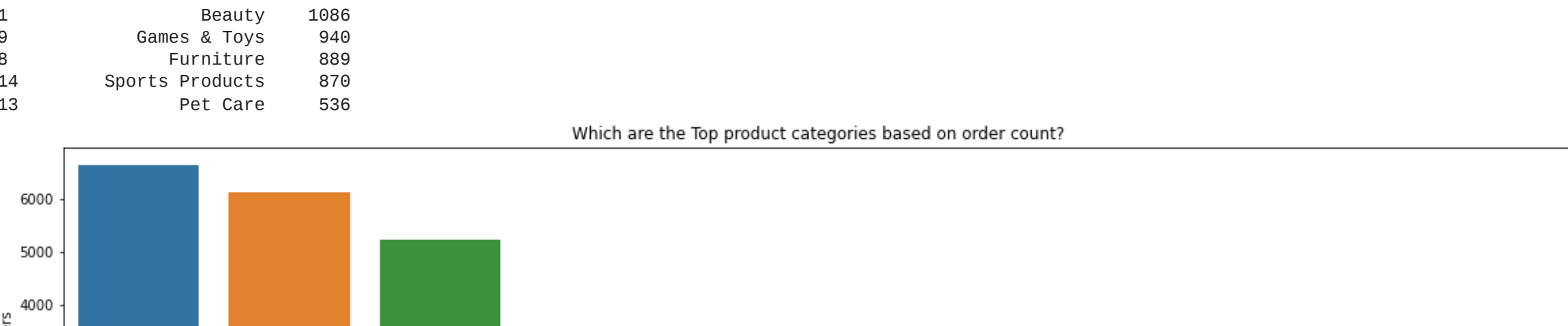
```
In [66]: #Checking top 10 occupations people belong based on total amount spent on their orders
sales_Occupation = df.groupby(['Occupation'], as_index = False) ['Amount'].sum().sort_values(by='Amount', ascending=False).head(10)

#Printing of order values by state-wise for reference
print(sales_Occupation)

#Plotting of graph
plt.figure(figsize=(20, 5))
sns.barplot(x='Occupation', y='Amount', data = sales_Occupation)
plt.title('Which are the Top Occupations people belong based on total amount spent on their orders?')

#Displaying the plot
plt.show()

Occupation  Amount
10  IT Sector  14755979
8  Healthcare  13934988
2  Aviation  12662298
3  Banking  10778010
7  Govt  8572212
9  Hospitality  6376405
12  Media  6258532
1  Automobile  5385906
4  Chemical  5297436
11  Lawyer  4981605
```



Insight 3: From the above graphs – we understand that most of the buyers belong to IT sector, Healthcare and Aviation sectors

Analysing the Product Category data

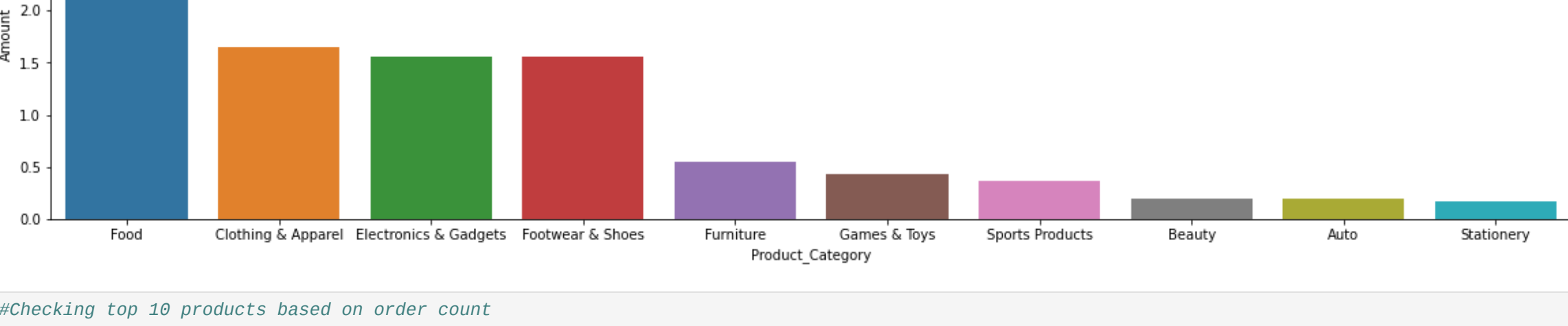
```
In [70]: #Checking top 10 product categories based on order count
#Grouping and aggregating values
sales_Product_Category = df.groupby(['Product_Category'], as_index = False) ['Orders'].sum().sort_values(by='Orders', ascending=False).head(10)

#Printing of order values by state-wise for reference
print(sales_Product_Category)

#Plotting of graph
plt.figure(figsize=(20, 5))
sns.barplot(x='Product_Category', y='Orders', data = sales_Product_Category)
plt.title('Which are the Top product categories based on order count?')

#Displaying the plot
plt.show()

Product_Category  Orders
3  Clothing & Apparel  6034
6  Food  6110
5  Electronics & Gadgets  5236
7  Footwear & Shoes  2646
11  Household Items  1331
1  Furniture  1886
9  Games & Toys  940
8  Sports Products  889
14  Pet Care  870
13  Stationery  536
```



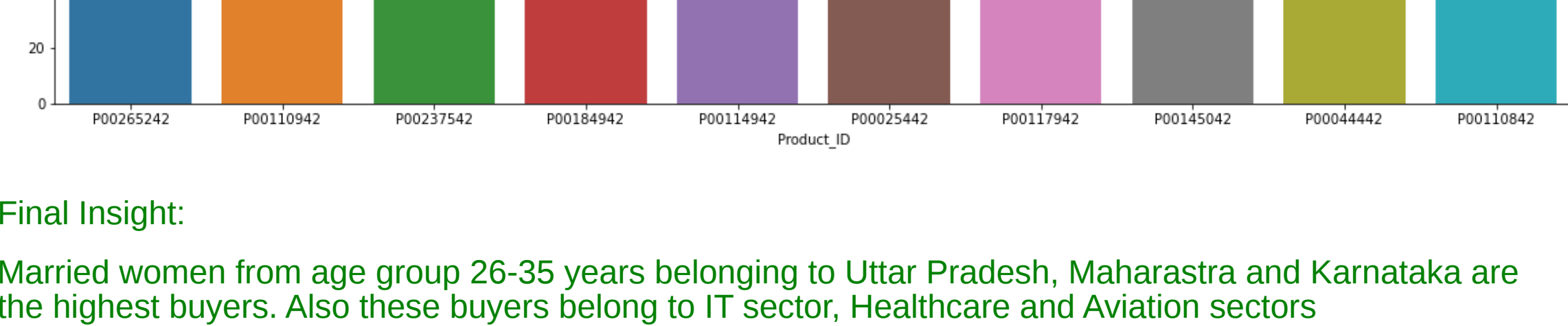
```
In [71]: #Checking top 10 product categories based on amount
#Grouping and aggregating values
sales_Product_Category = df.groupby(['Product_Category'], as_index = False) ['Amount'].sum().sort_values(by='Amount', ascending=False).head(10)

#Printing of order values by state-wise for reference
print(sales_Product_Category)

#Plotting of graph
plt.figure(figsize=(20, 5))
sns.barplot(x='Product_Category', y='Amount', data = sales_Product_Category)
plt.title('Which are the Top product categories based on amount?')

#Displaying the plot
plt.show()

Product_Category  Amount
3  Clothing & Apparel  3933383
6  Food  3693383
5  Electronics & Gadgets  15643846
7  Footwear & Shoes  15575269
9  Furniture  15400951
9  Games & Toys  4331694
14  Sports Products  3635933
1  Furniture  1934844
0  Auto  1958669
15  Stationery  1676651
```



Insight 3: From the above graphs – we understand that most of the buyers belong to IT sector, Healthcare and Aviation sectors

Married women from age group 26-35 years belonging to Uttar Pradesh, Maharashtra and Karnataka are the highest buyers. Also these buyers belong to IT sector, Healthcare and Aviation sectors