

```
In [1]: import pandas as pd
```

```
In [2]: import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
import seaborn as sns
```

```
In [5]: import pandas as pd

df = pd.read_csv('C://Users//Nikki Chauhan//Downloads//Diwali Sales Data.csv', encoding='latin1')
```

In [6]: df

Out[6]:

	User_ID	Cust_name	Product_ID	Gender	Age Group	Age	Marital_Status	State	Zone	Occupation	Product_Categor
0	1002903	Sanskriti	P00125942	F	26-35	28	0	Maharashtra	Western	Healthcare	Aut
1	1000732	Kartik	P00110942	F	26-35	35	1	Andhra Pradesh	Southern	Govt	Aut
2	1001990	Bindu	P00118542	F	26-35	35	1	Uttar Pradesh	Central	Automobile	Aut
3	1001425	Sudevi	P00237842	M	0-17	16	0	Karnataka	Southern	Construction	Aut
4	1000588	Joni	P00057942	M	26-35	28	1	Gujarat	Western	Food Processing	Aut
...	...	...	...	...	...	...	...	...	...	...	.
11246	1000695	Manning	P00296942	M	18-25	19	1	Maharashtra	Western	Chemical	Offic
11247	1004089	Reichenbach	P00171342	M	26-35	33	0	Haryana	Northern	Healthcare	Veterinar
11248	1001209	Oshin	P00201342	F	36-45	40	0	Madhya Pradesh	Central	Textile	Offic
11249	1004023	Noonan	P00059442	M	36-45	37	0	Karnataka	Southern	Agriculture	Offic
11250	1002744	Brumley	P00281742	F	18-25	19	0	Maharashtra	Western	Healthcare	Offic

11251 rows × 15 columns



In [7]: df.shape

Out[7]: (11251, 15)

In [8]: `df.head(10)`

Out[8]:

	User_ID	Cust_name	Product_ID	Gender	Age Group	Age	Marital_Status	State	Zone	Occupation	Product_Category	Or
0	1002903	Sanskriti	P00125942	F	26-35	28	0	Maharashtra	Western	Healthcare		Auto
1	1000732	Kartik	P00110942	F	26-35	35	1	Andhra Pradesh	Southern	Govt		Auto
2	1001990	Bindu	P00118542	F	26-35	35	1	Uttar Pradesh	Central	Automobile		Auto
3	1001425	Sudevi	P00237842	M	0-17	16	0	Karnataka	Southern	Construction		Auto
4	1000588	Joni	P00057942	M	26-35	28	1	Gujarat	Western	Food Processing		Auto
5	1000588	Joni	P00057942	M	26-35	28	1	Himachal Pradesh	Northern	Food Processing		Auto
6	1001132	Balk	P00018042	F	18-25	25	1	Uttar Pradesh	Central	Lawyer		Auto
7	1002092	Shivangi	P00273442	F	55+	61	0	Maharashtra	Western	IT Sector		Auto
8	1003224	Kushal	P00205642	M	26-35	35	0	Uttar Pradesh	Central	Govt		Auto
9	1003650	Ginny	P00031142	F	26-35	26	1	Andhra Pradesh	Southern	Media		Auto

In [9]: 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 [11]: df.drop(['Status', 'unnamed1'],axis=1,inplace=True)
df
```

Out[11]:

	User_ID	Cust_name	Product_ID	Gender	Age Group	Age	Marital_Status	State	Zone	Occupation	Product_Categor
0	1002903	Sanskriti	P00125942	F	26-35	28	0	Maharashtra	Western	Healthcare	Aut
1	1000732	Kartik	P00110942	F	26-35	35	1	Andhra Pradesh	Southern	Govt	Aut
2	1001990	Bindu	P00118542	F	26-35	35	1	Uttar Pradesh	Central	Automobile	Aut
3	1001425	Sudevi	P00237842	M	0-17	16	0	Karnataka	Southern	Construction	Aut
4	1000588	Joni	P00057942	M	26-35	28	1	Gujarat	Western	Food Processing	Aut
...	...	...	...	...	...	...	...	...	...	...	.
11246	1000695	Manning	P00296942	M	18-25	19	1	Maharashtra	Western	Chemical	Offic
11247	1004089	Reichenbach	P00171342	M	26-35	33	0	Haryana	Northern	Healthcare	Veterinar
11248	1001209	Oshin	P00201342	F	36-45	40	0	Madhya Pradesh	Central	Textile	Offic
11249	1004023	Noonan	P00059442	M	36-45	37	0	Karnataka	Southern	Agriculture	Offic
11250	1002744	Brumley	P00281742	F	18-25	19	0	Maharashtra	Western	Healthcare	Offic

11251 rows × 13 columns



```
In [12]: pd.isnull(df).sum()
```

```
Out[12]: 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 [13]: df.dropna(inplace=True)
```

```
In [14]: 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           0
dtype: int64
```

In [15]: df.info()

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 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
```

In [17]: df['Amount']=df['Amount'].astype('int')

In [18]: df['Amount'].dtypes

Out[18]: dtype('int32')

In [19]: df.columns

Out[19]: Index(['User\_ID', 'Cust\_name', 'Product\_ID', 'Gender', 'Age Group', 'Age',  
 'Marital\_Status', 'State', 'Zone', 'Occupation', 'Product\_Category',  
 'Orders', 'Amount'],  
 dtype='object')

```
In [20]: df.describe()
```

```
Out[20]:
```

	User_ID	Age	Marital_Status	Orders	Amount
count	1.123900e+04	11239.000000	11239.000000	11239.000000	11239.000000
mean	1.003004e+06	35.410357	0.420055	2.489634	9453.610553
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

```
In [22]: df[['Age', 'Orders', 'Amount']].describe()
```

```
Out[22]:
```

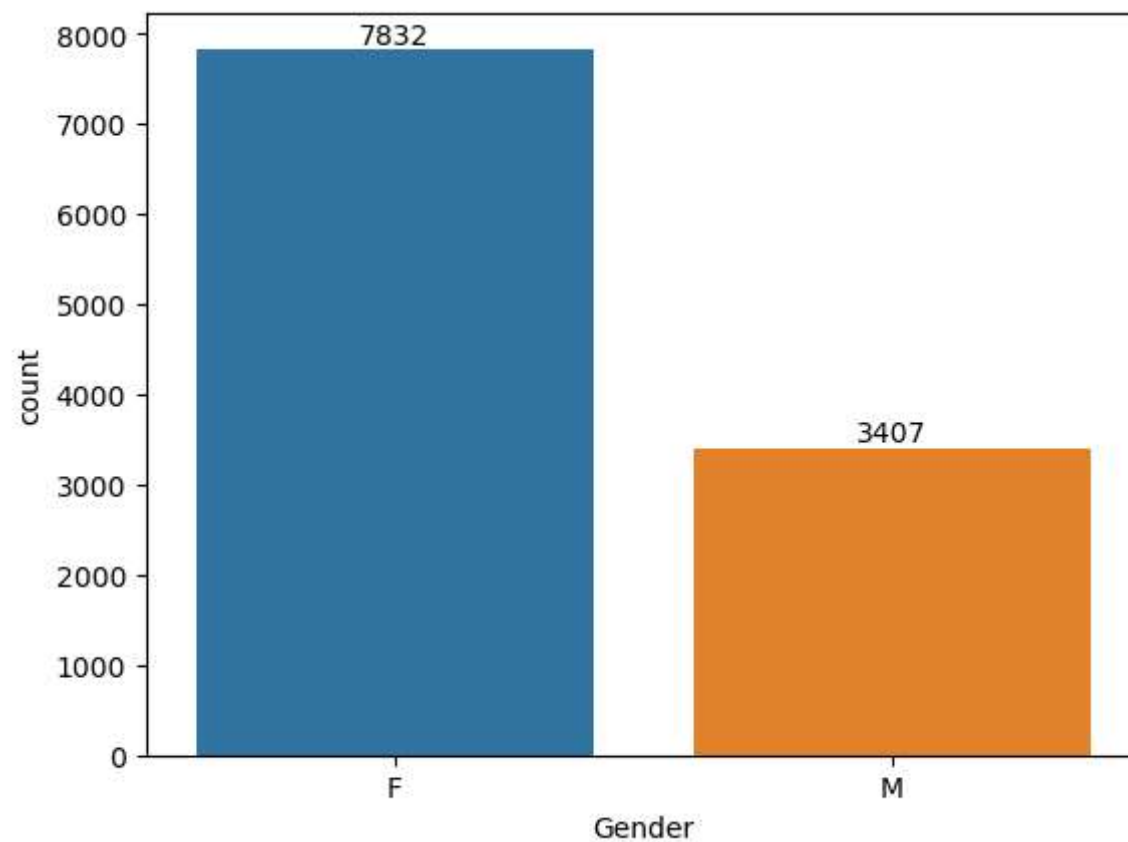
	Age	Orders	Amount
count	11239.000000	11239.000000	11239.000000
mean	35.410357	2.489634	9453.610553
std	12.753866	1.114967	5222.355168
min	12.000000	1.000000	188.000000
25%	27.000000	2.000000	5443.000000
50%	33.000000	2.000000	8109.000000
75%	43.000000	3.000000	12675.000000
max	92.000000	4.000000	23952.000000



```
In [23]: df.columns
```

```
Out[23]: Index(['User_ID', 'Cust_name', 'Product_ID', 'Gender', 'Age Group', 'Age',  
               'Marital_Status', 'State', 'Zone', 'Occupation', 'Product_Category',  
               'Orders', 'Amount'],  
              dtype='object')
```

```
In [24]: ax=sns.countplot(x='Gender',data=df)  
for bars in ax.containers:  
    ax.bar_label(bars)
```



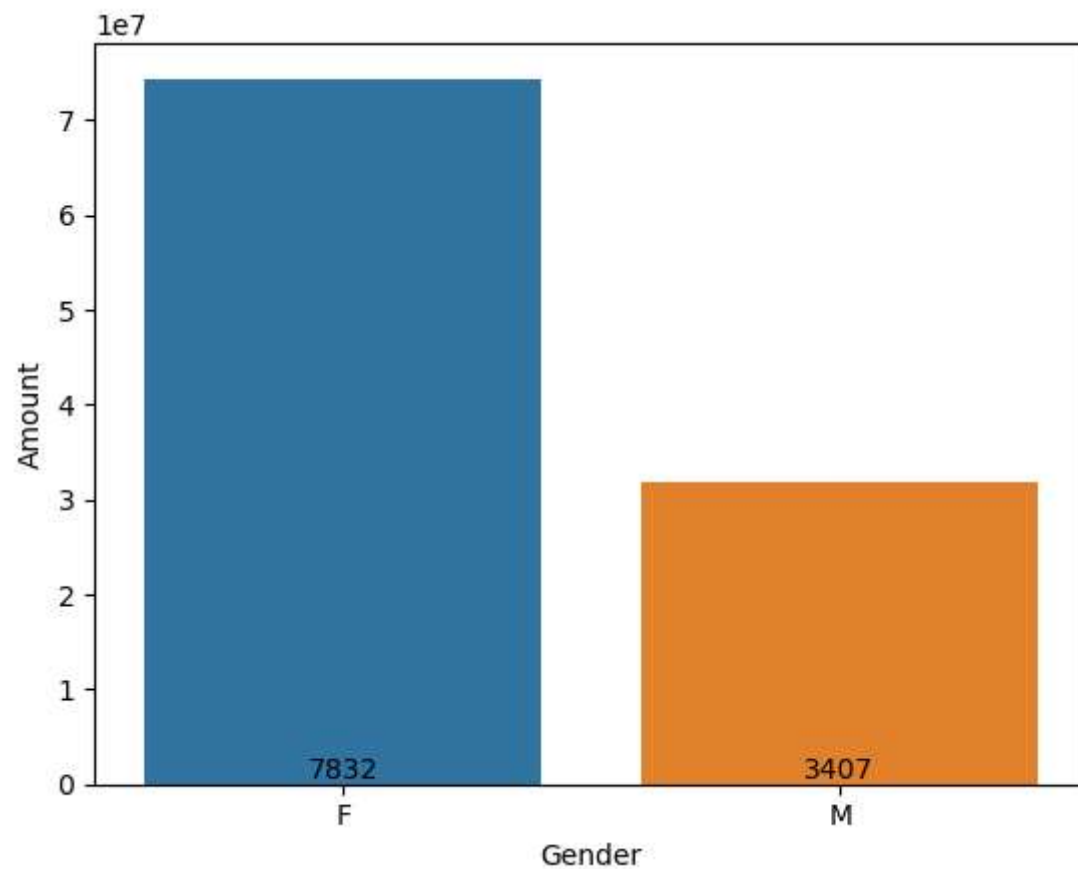
```
In [26]: df.groupby(['Gender'],as_index=False)['Amount'].sum().sort_values(by='Amount',ascending=True)
```

Out[26]:

	Gender	Amount
1	M	31913276
0	F	74335853

```
In [27]: sales_gen=df.groupby(['Gender'],as_index=False)['Amount'].sum().sort_values(by='Amount',ascending=False)
ax1=sns.barplot(x='Gender',y='Amount',data=sales_gen)
ax1.bar_label(bars)
```

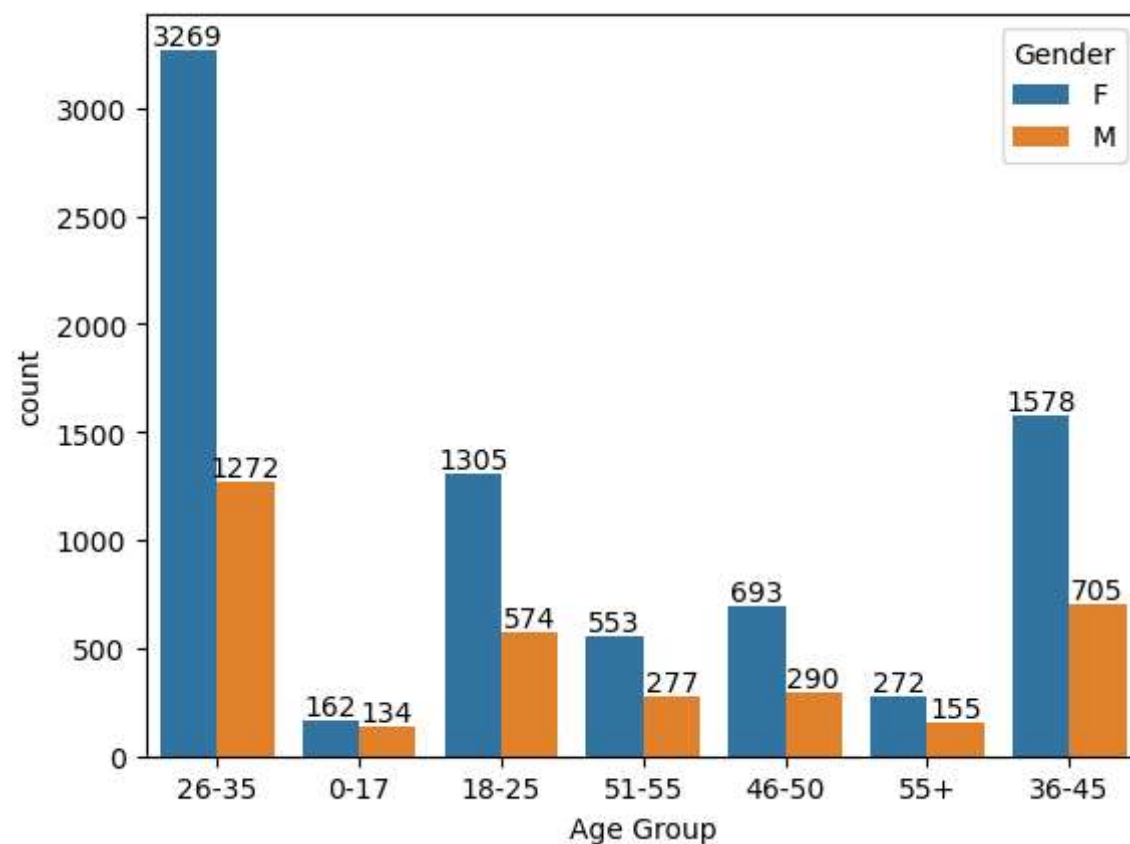
Out[27]: [Text(0, 0, '7832'), Text(0, 0, '3407')]



```
In [28]: df.columns
```

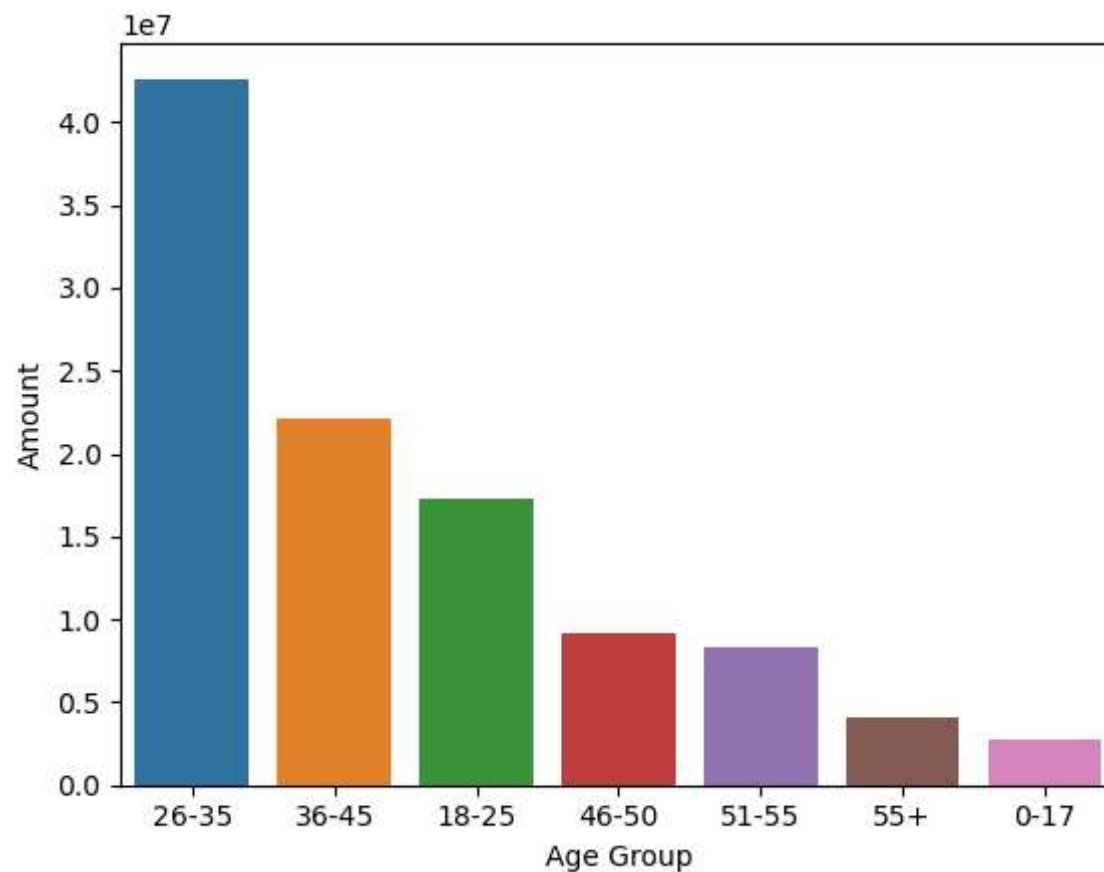
```
Out[28]: Index(['User_ID', 'Cust_name', 'Product_ID', 'Gender', 'Age Group', 'Age',  
              'Marital_Status', 'State', 'Zone', 'Occupation', 'Product_Category',  
              'Orders', 'Amount'],  
              dtype='object')
```

```
In [29]: ax=sns.countplot(data=df,x='Age Group',hue='Gender')  
for bars in ax.containers:  
    ax.bar_label(bars)
```



```
In [30]: 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[30]: <AxesSubplot:xlabel='Age Group', ylabel='Amount'>
```



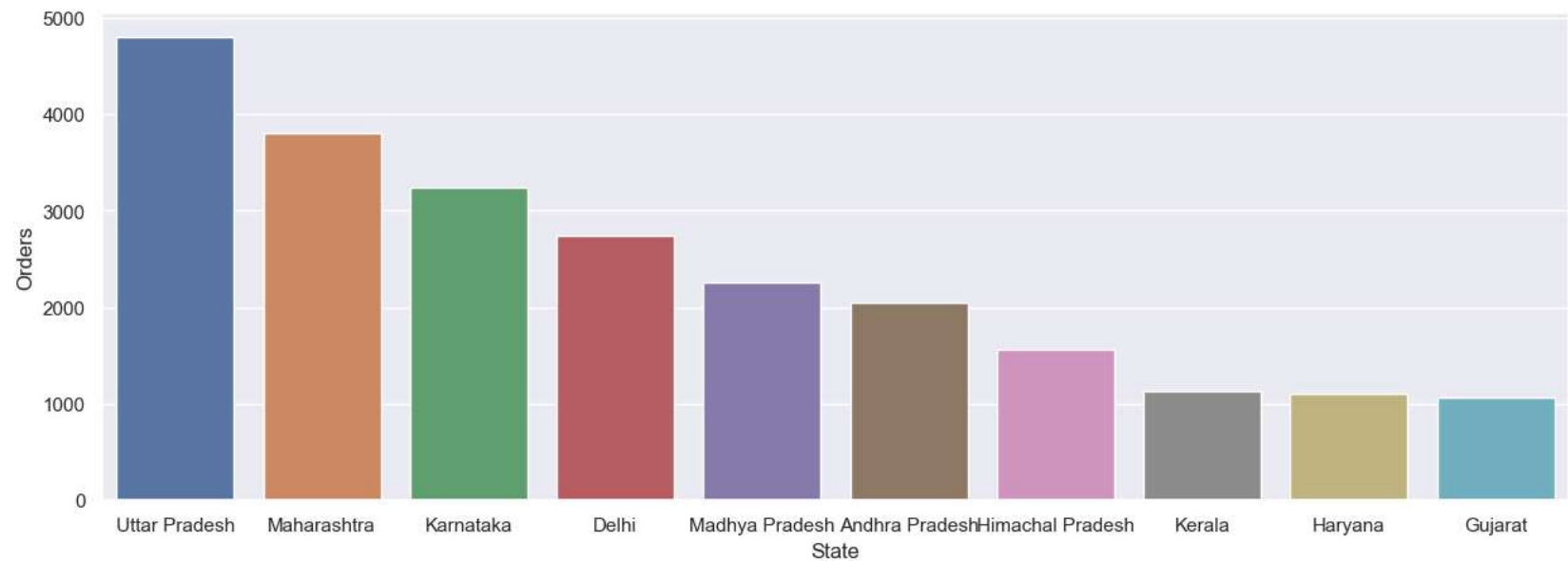
```
In [31]: df.columns
```

```
Out[31]: Index(['User_ID', 'Cust_name', 'Product_ID', 'Gender', 'Age Group', 'Age',
               'Marital_Status', 'State', 'Zone', 'Occupation', 'Product_Category',
               'Orders', 'Amount'],
              dtype='object')
```

```
In [32]: sales_state=df.groupby(['State'],as_index=False)['Orders'].sum().sort_values(by='Orders',ascending=False).head(10)

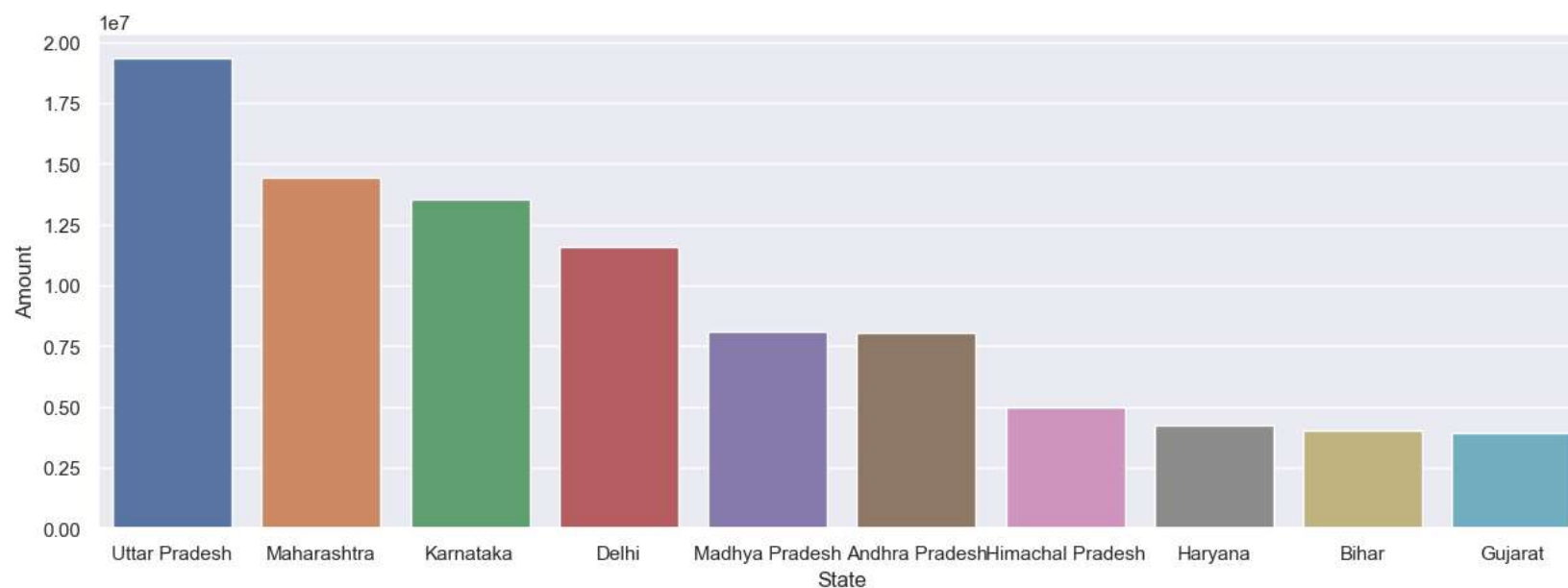
sns.set(rc={'figure.figsize':(15,5)})
sns.barplot(data=sales_state,x='State',y='Orders')
```

Out[32]: <AxesSubplot:xlabel='State', ylabel='Orders'>



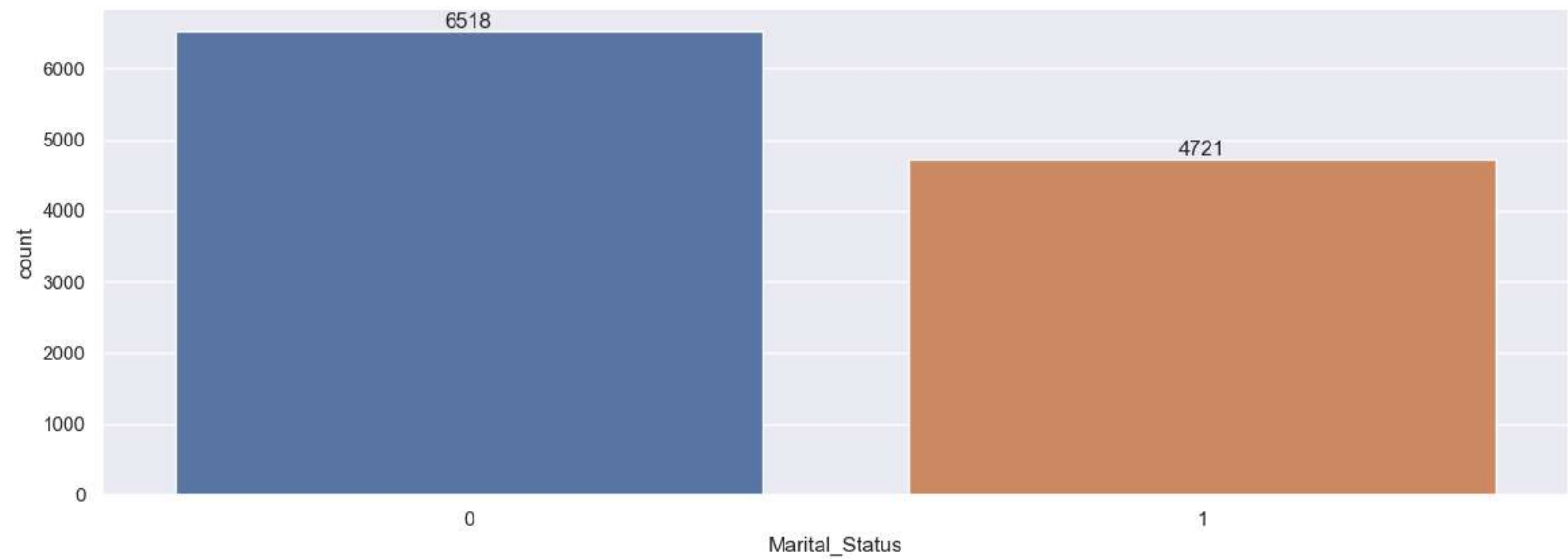
```
In [33]: sales_state=df.groupby(['State'],as_index=False)['Amount'].sum().sort_values(by='Amount',ascending=False).head(10)
sns.set(rc={'figure.figsize':(15,5)})
sns.barplot(data=sales_state,x='State',y='Amount')
```

Out[33]: <AxesSubplot:xlabel='State', ylabel='Amount'>



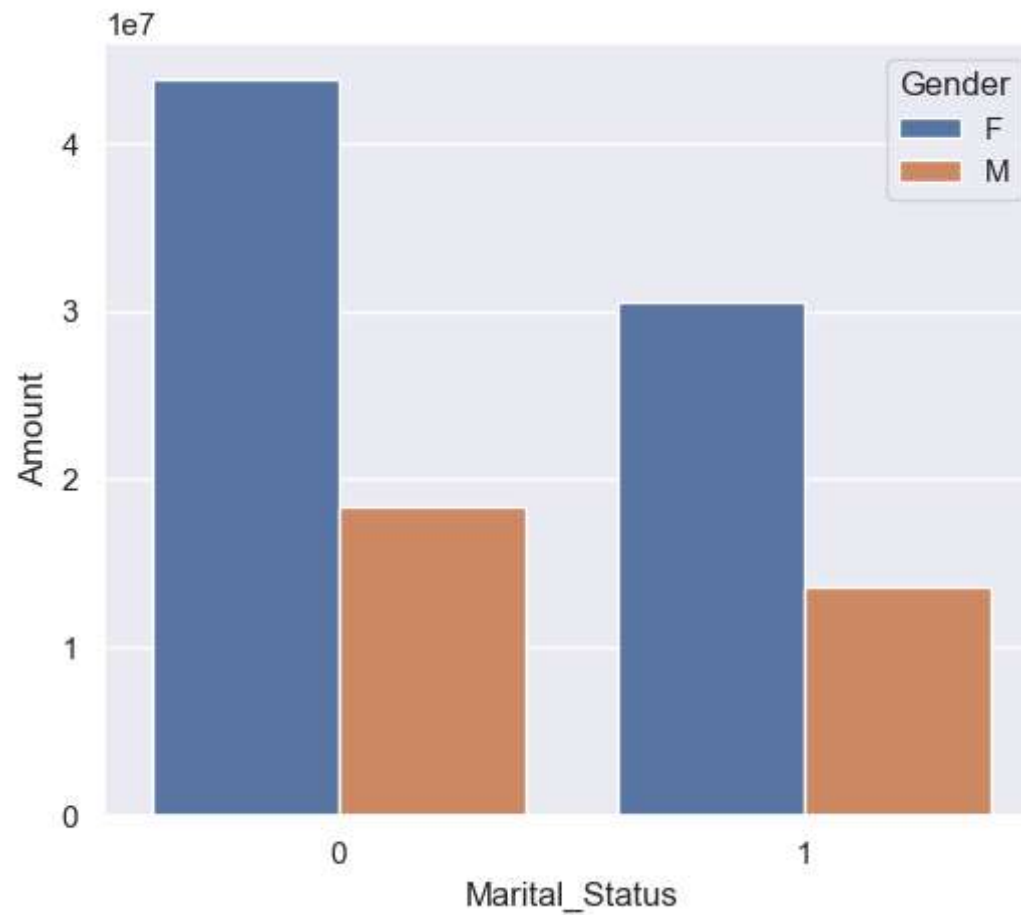
```
In [37]: ax = sns.countplot(data=df, x='Marital_Status')
sns.set(rc={'figure.figsize':(6,5)})

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',  
sns.set(rc={'figure.figsize':(6,5)})  
sns.barplot(data=sales_state,x='Marital_Status',y='Amount',hue='Gender')
```

Out[41]: <AxesSubplot:xlabel='Marital\_Status', ylabel='Amount'>

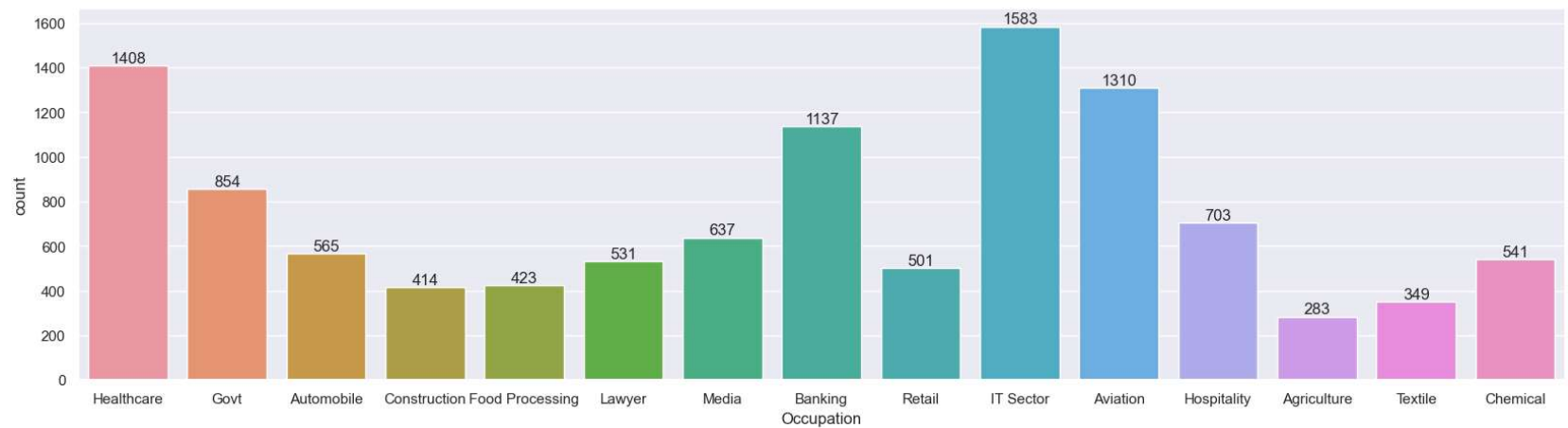




```
In [42]: df.columns
```

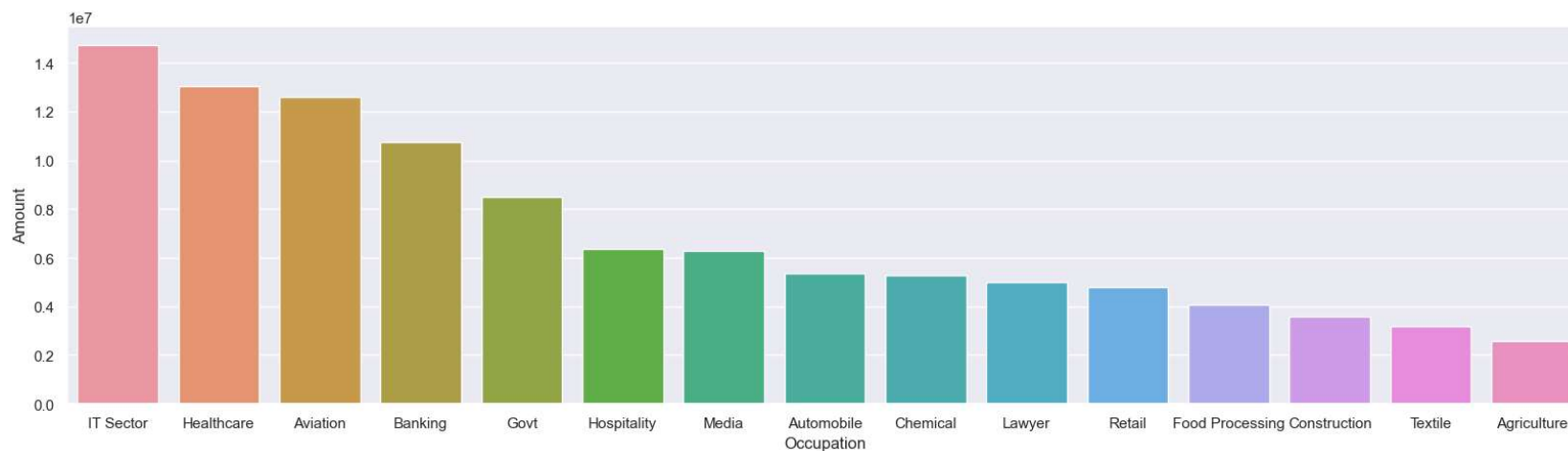
```
Out[42]: Index(['User_ID', 'Cust_name', 'Product_ID', 'Gender', 'Age Group', 'Age',  
              'Marital_Status', 'State', 'Zone', 'Occupation', 'Product_Category',  
              'Orders', 'Amount'],  
              dtype='object')
```

```
In [44]: sns.set(rc={'figure.figsize':(20,5)})  
ax=sns.countplot(data=df,x='Occupation')  
for bars in ax.containers:  
    ax.bar_label(bars)
```

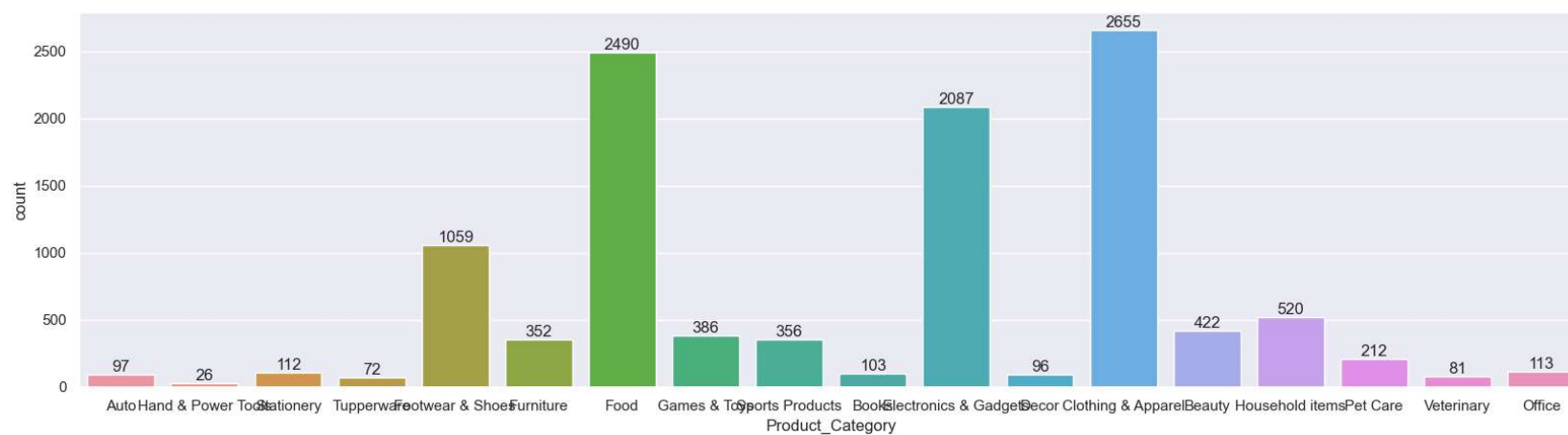


```
In [45]: sales_state=df.groupby(['Occupation',],as_index=False)['Amount'].sum().sort_values(by='Amount',ascending=False)
sns.set(rc={'figure.figsize':(20,5)})
sns.barplot(data=sales_state,x='Occupation',y='Amount')
```

Out[45]: <AxesSubplot:xlabel='Occupation', ylabel='Amount'>

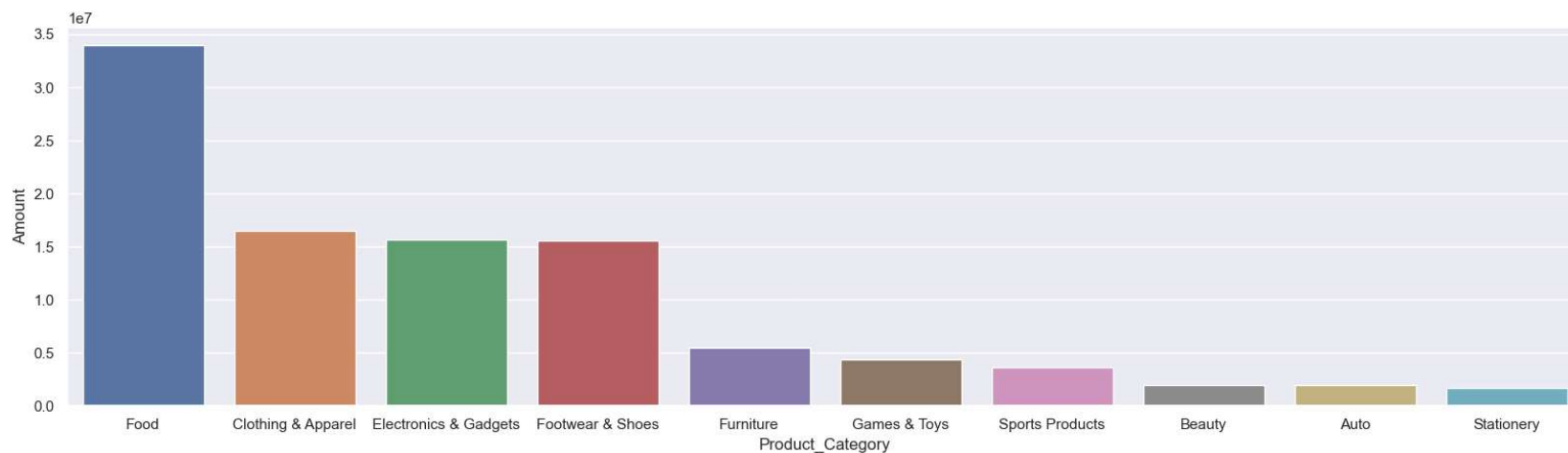


```
In [46]: sns.set(rc={'figure.figsize':(20,5)})
ax=sns.countplot(data=df,x='Product_Category')
for bars in ax.containers:
    ax.bar_label(bars)
```



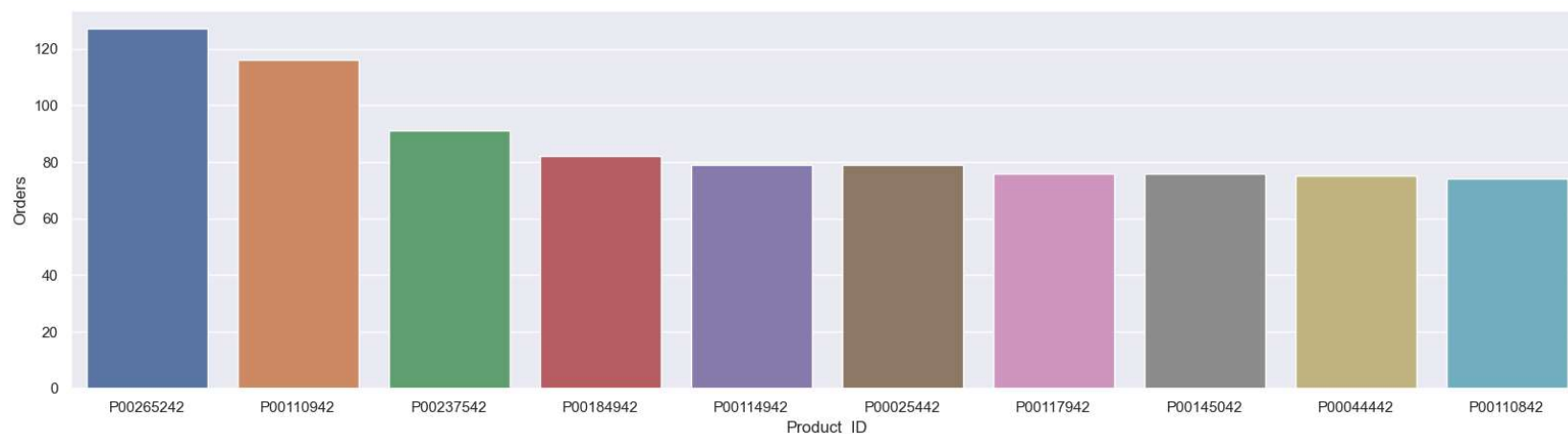
```
In [47]: sales_state=df.groupby(['Product_Category'],as_index=False)['Amount'].sum().sort_values(by='Amount',ascending=True)
sns.set(rc={'figure.figsize':(20,5)})
sns.barplot(data=sales_state,x='Product_Category',y='Amount')
```

Out[47]: <AxesSubplot:xlabel='Product\_Category', ylabel='Amount'>



```
In [48]: sales_state=df.groupby(['Product_ID'],as_index=False)['Orders'].sum().sort_values(by='Orders',ascending=False)
sns.set(rc={'figure.figsize':(20,5)})
sns.barplot(data=sales_state,x='Product_ID',y='Orders')
```

Out[48]: <AxesSubplot:xlabel='Product\_ID', ylabel='Orders'>



In [ ]: