

Diwali Sales Data Analysis

```
In [33]: # Import Python Libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
%matplotlib inline
import seaborn as sns
```

```
In [6]: # Import csv file
df=pd.read_csv('Diwali Sales Data.csv',encoding='unicode_escape')
```

```
In [7]: df.shape
# show rows and columns
```

Out[7]: (11251, 15)

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

Out[8]:

	User_ID	Cust_name	Product_ID	Gender	Age Group	Age	Marital_Status	State	Zone	Occupation	Product_Category	Orders	Amount	Status	unnamed
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	23877.0	NaN	NaN

```
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
```

Data Cleaning

```
In [13]: df.drop(['Status', 'unnamed1'],axis=1,inplace=True)
# Remove empty columns
```

```
In [14]: 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
```

```
In [15]: pd.isnull(df).sum()
# Check for null values
```

```
Out[15]: 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 [16]: df.dropna(inplace=True)
        # Remove null values
```

```
In [17]: df.shape
```

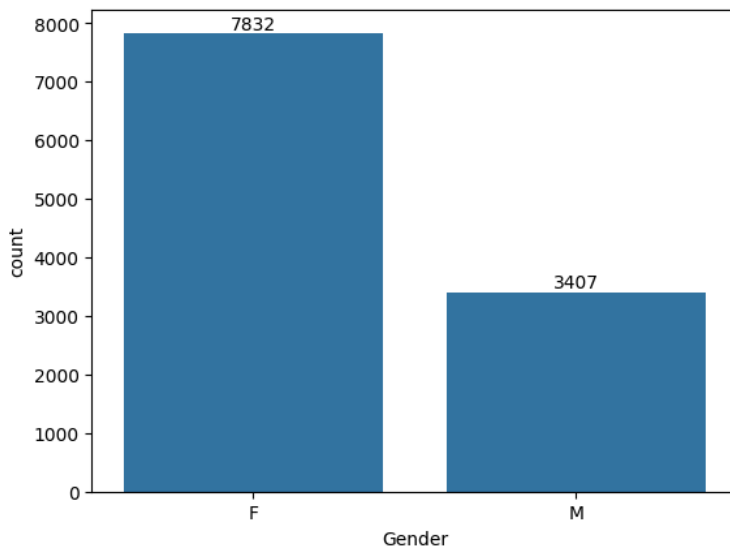
```
Out[17]: (11239, 13)
```

EDA(Exploratory Data Analysis)

```
In [18]: df.columns
```

```
Out[18]: Index(['User_ID', 'Cust_name', 'Product_ID', 'Gender', 'Age Group', 'Age',
               'Marital_Status', 'State', 'Zone', 'Occupation', 'Product_Category',
               'Orders', 'Amount'],
              dtype='object')
```

```
In [19]: # Plot a graph for Gender vs it's count
        ax=sns.countplot(x='Gender',data=df)
        for bars in ax.containers:
            ax.bar_label(bars)
```

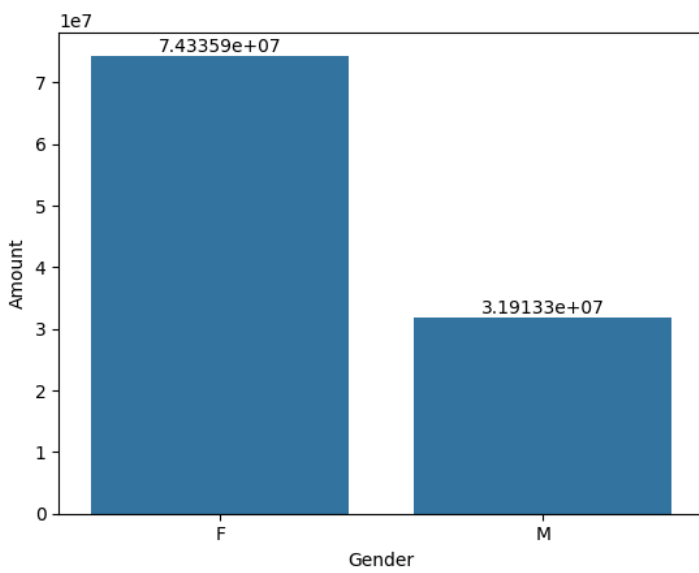


```
In [20]: # Plot a Bar chart for Gender Vs Total Amount
        df.groupby(['Gender'], as_index=False)['Amount'].sum().sort_values(by='Amount', ascending=False)
```

```
Out[20]:
```

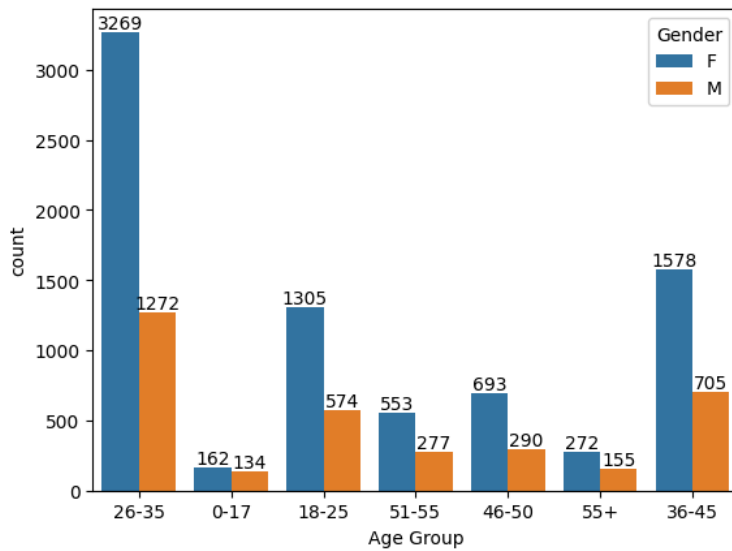
	Gender	Amount
0	F	74335856.43
1	M	31913276.00

```
In [22]: sales_gen=df.groupby(['Gender'], as_index=False)['Amount'].sum().sort_values(by='Amount', ascending=False)
        var_sales_gen=sns.barplot(x = 'Gender',y= 'Amount' ,data = sales_gen)
        for bars in var_sales_gen.containers:
            var_sales_gen.bar_label(bars)
```



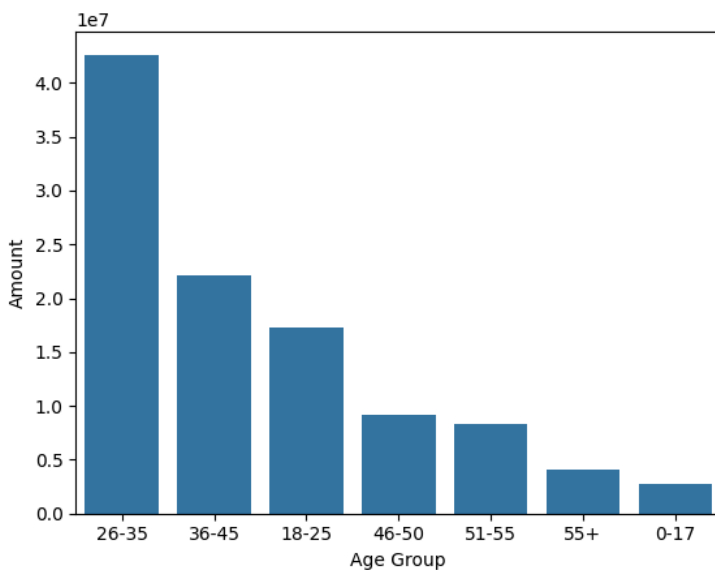
Insight-> From above we can see easily that most of the buyer are female and purchasing power of female is more.

```
In [26]: # Plot a graph for Age and it's count
var_age = sns.countplot(data = df, x = 'Age Group', hue = 'Gender')
for bars in var_age.containers:
    var_age.bar_label(bars)
```



```
In [27]: # Plot a graph for age- 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)
```

Out[27]: <Axes: xlabel='Age Group', ylabel='Amount'>

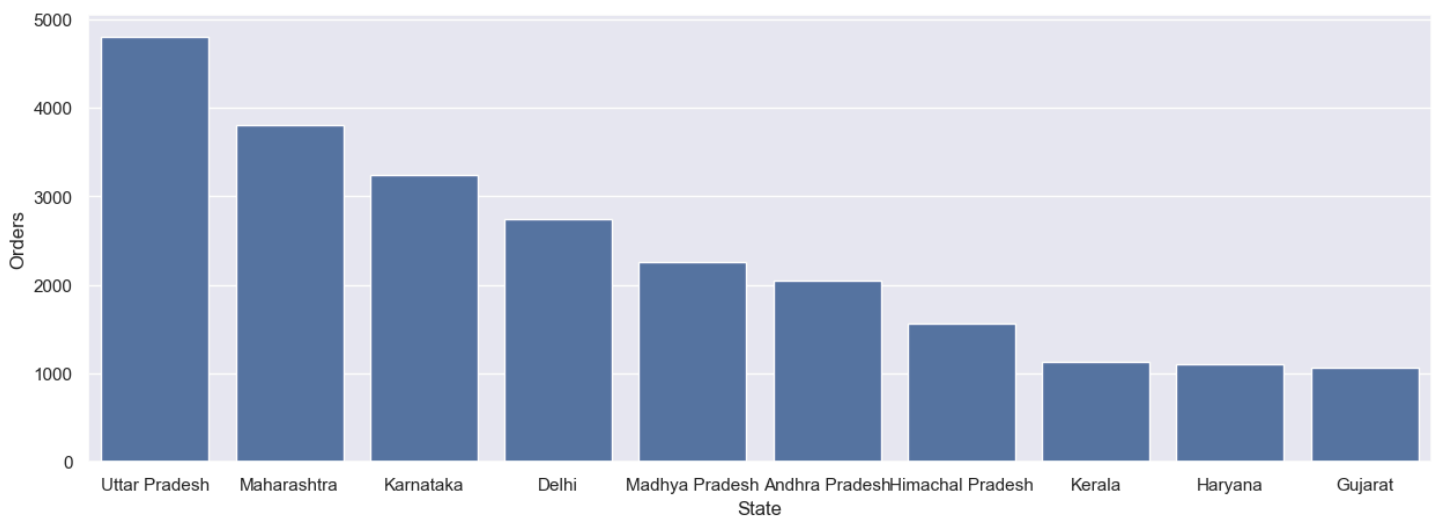


Insight-> From above graphs we can see that most of the buyers are of age group between 26-35 yrs female.

```
In [28]: # Plot a graph for State- 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')
```

Out[28]: <Axes: xlabel='State', ylabel='Orders'>

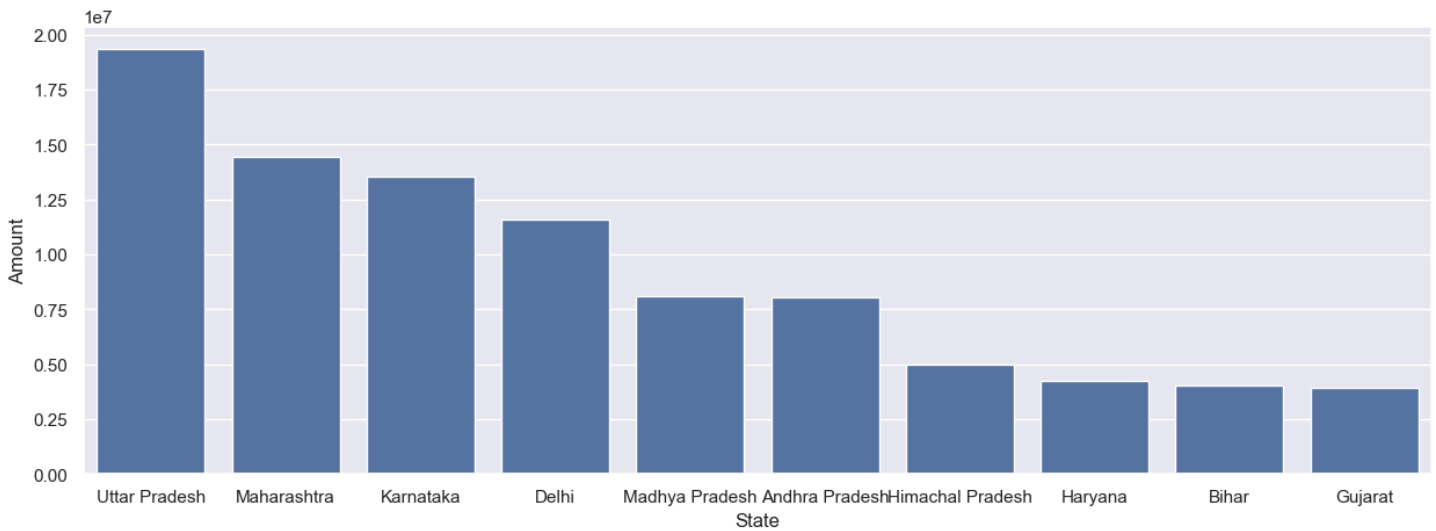


In [29]: # Plot a Graph for State-total amount/sales from top 10 states

```
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')
```

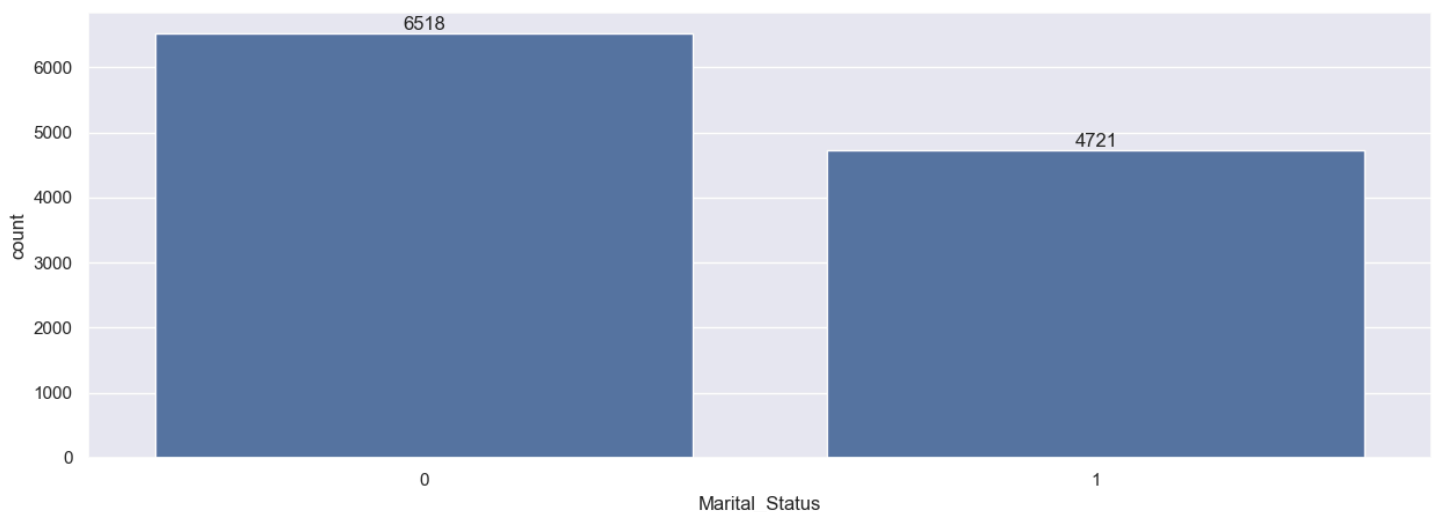
Out[29]: <Axes: xlabel='State', ylabel='Amount'>



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

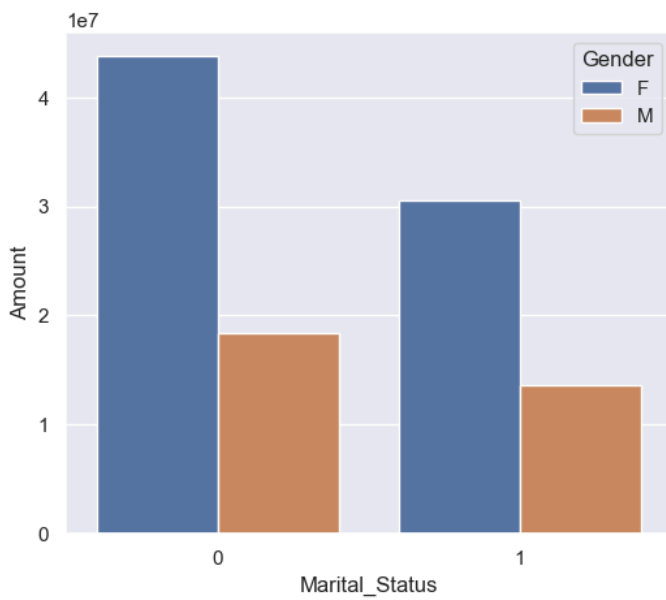
And also noticed that orders of kerala is more compared to Haryana but purchasing power of Haryana is more compared to kerala

```
In [30]: # Plot a graph for Marital status
ax = sns.countplot(data = df, x = 'Marital_Status')
sns.set(rc={'figure.figsize':(7,5)})
for bars in ax.containers:
    ax.bar_label(bars)
```



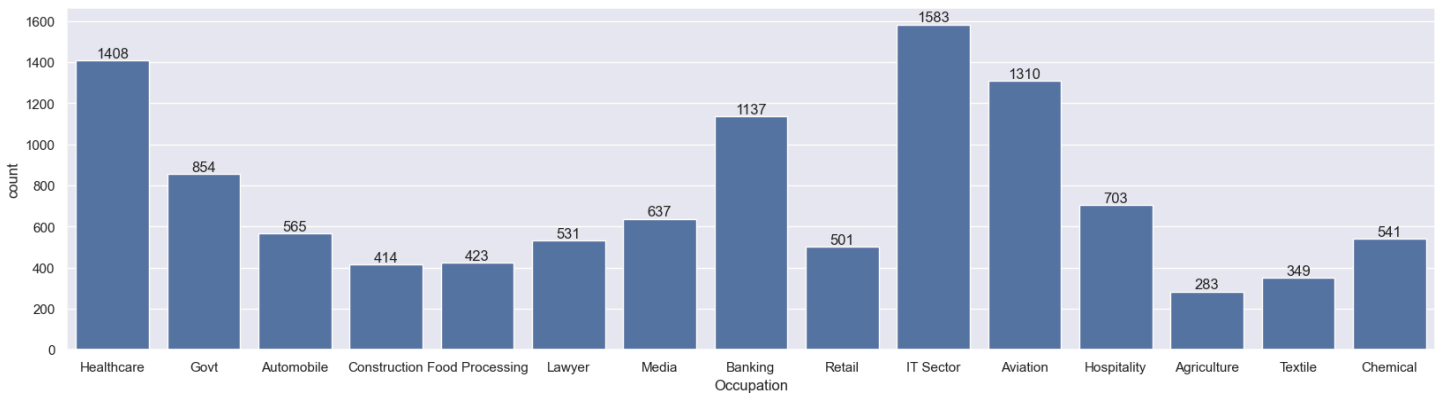
```
In [31]: sales_state = df.groupby(['Marital_Status', 'Gender'], as_index=False)['Amount'].sum().sort_values(by='Amount', ascending=False)
sns.set(rc={'figure.figsize':(6,5)})
sns.barplot(data = sales_state, x = 'Marital_Status',y= 'Amount', hue='Gender')
```

Out[31]: <Axes: xlabel='Marital_Status', ylabel='Amount'>



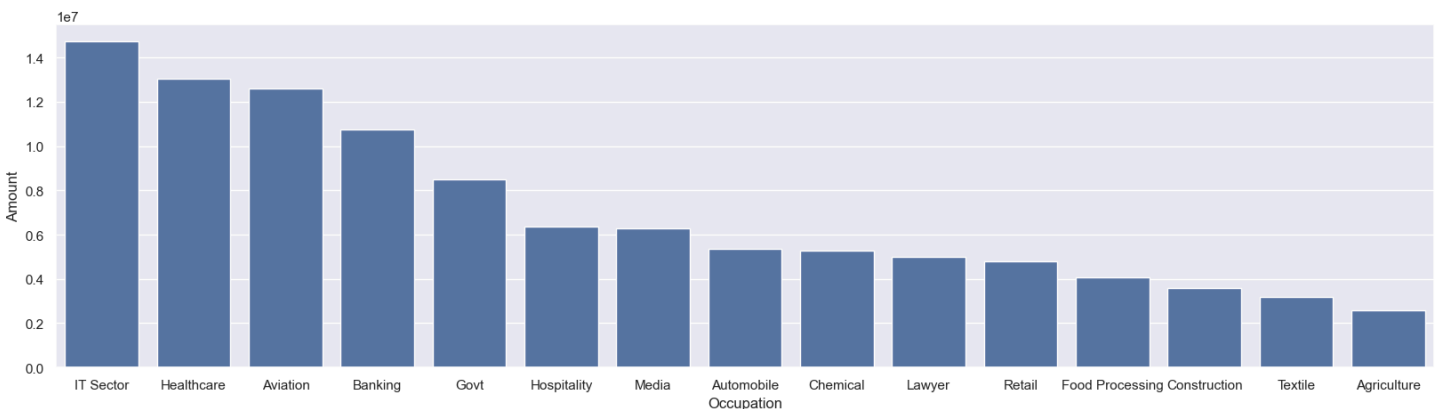
Insight-> From above graphs we can see that most of the buyers are married (women) and they have high purchasing power.

```
In [32]: # Plot a graph for Occupation.
sns.set(rc={'figure.figsize':(20,5)})
ax = sns.countplot(data = df, x = 'Occupation')
for bars in ax.containers:
    ax.bar_label(bars)
```



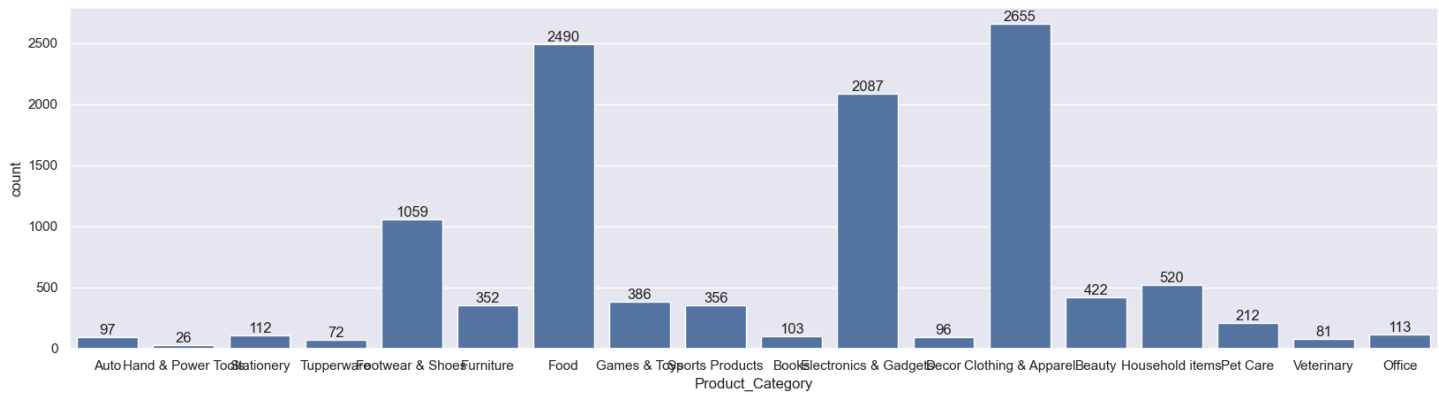
```
In [34]: sales_state = df.groupby(['Occupation'], as_index=False)['Amount'].sum().sort_values(by='Amount', ascending=False)
sns.set(rc={'figure.figsize':(20,5)})
sns.barplot(data = sales_state, x = 'Occupation', y= 'Amount')
```

Out[34]: <Axes: xlabel='Occupation', ylabel='Amount'>



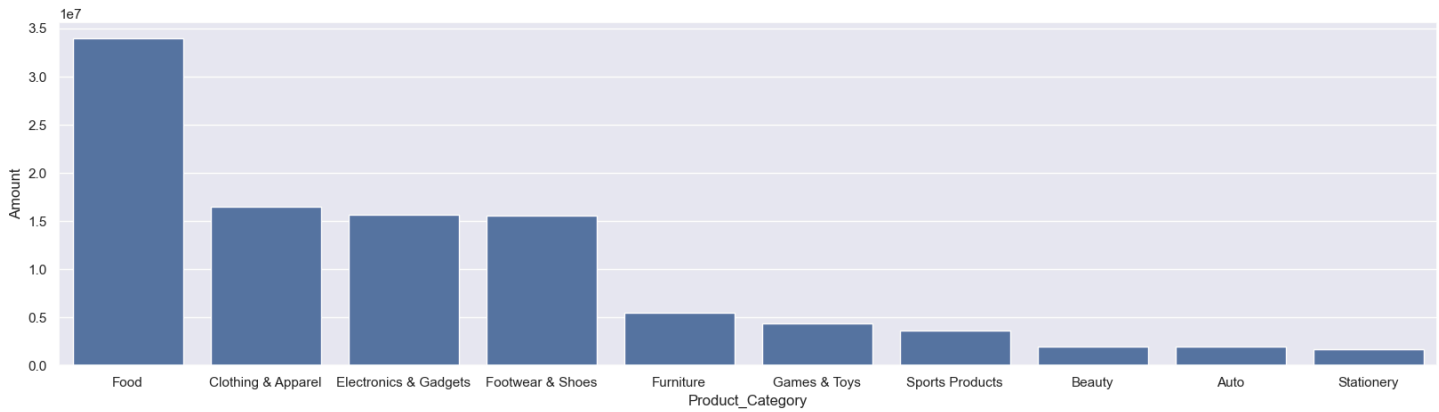
Insight-> From above graphs we can see that most of the buyers are working in IT, Healthcare and Aviation sector.

```
In [35]: # Plot a graph for Product category
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 [36]: 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')
```

```
Out[36]: <Axes: xlabel='Product_Category', ylabel='Amount'>
```

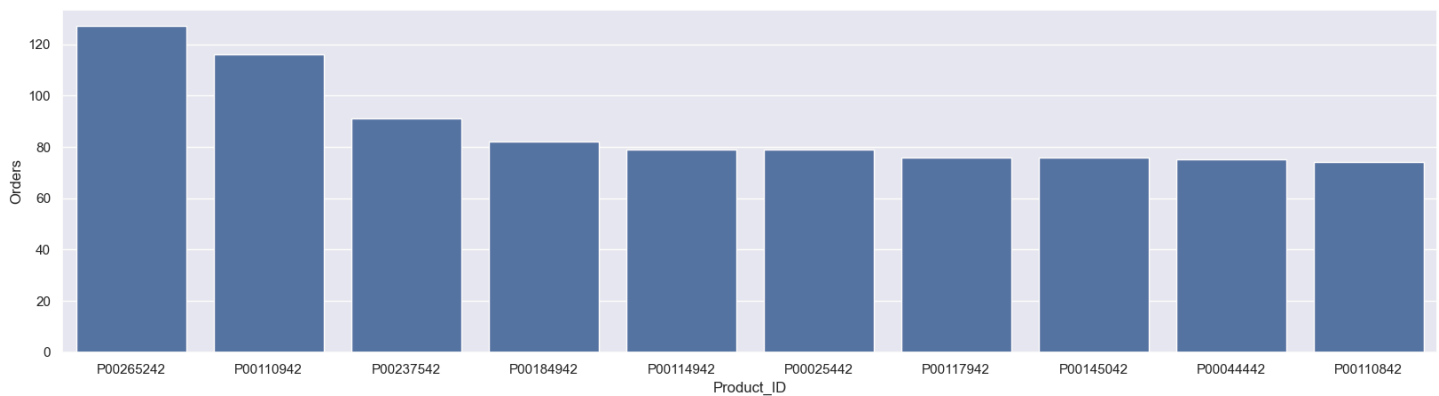


Insight-> From above graphs we can see that most of the sold products are from Food, Clothing and Electronics category.

Note--Rank of Food(Orders) is 8th but Rank of food(Purchasing power) is 1st.

```
In [37]: sales_state = df.groupby(['Product_ID'], as_index=False)['Orders'].sum().sort_values(by='Orders', ascending=False).head(10)
sns.set(rc={'figure.figsize':(20,5)})
sns.barplot(data = sales_state, x = 'Product_ID',y= 'Orders')
```

```
Out[37]: <Axes: xlabel='Product_ID', ylabel='Orders'>
```



Conclusion

From above all plots, we can conclude that-Married women age group 26-35 yrs from UP, Maharastra and Karnataka working in IT, Healthcare and Aviation are more likely to buy products from Food, Clothing and Electronics category.

```
In [ ]:
```