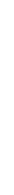


DIWALI SALES REPORT





AGENDA

Problem statement

Research Questions

Data overview

Methodology

EDA

Conclusion



PROBLEM STATEMENT

To determine which customer segments, regions and product categories contribute most to Diwali sales and use these insights to improve festive marketing and business decisions

RESEARCH QUESTIONS

- Who are the major buyers during Diwali (their gender , age and marital status)
- Which state contribute most to the sales
- Which occupation contribute most to the Diwali sales
- Which products generates the highest revenue
- What are the top 10 most sold products
- How can these insights will help the company in targeted marketing and its inventory optimization .

DATA OVERVIEW

For this analysis we have used the Diwali sales data , it consist of various columns like 'User_ID', 'Cust_name', 'Product_ID', 'Gender', 'Age Group', 'Age', 'Marital_Status', 'State', 'Zone', 'Occupation', 'Product_Category', 'Orders', 'Amount' . But for this analysis we will only use relevant columns .

```
df.head()
```

	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	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

Next steps:

[Generate code with df](#)[New interactive sheet](#)

df.columns

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

RELEVANT COLUMNS USED FOR THIS EDA

To focus only on relevant columns we have used Data cleaning procedures to concentrate only on essential columns for this analysis .

Columns we will use for this research are 'User_ID', 'Cust_name', 'Product_ID', 'Gender', 'Age Group', 'Age', 'Marital_Status', 'State', 'Zone', 'Occupation', 'Product_Category', 'Orders', 'Amount' .

	User_ID	Cust_name	Product_ID	Gender	Age Group	Age	Shaadi	State	Zone	Occupation	Product_Category	Orders	Amount
0	1002903	Sanskriti	P00125942	F	26-35	28	0	Maharashtra	Western	Healthcare	Auto	1	23952
1	1000732	Kartik	P00110942	F	26-35	35	1	Andhra Pradesh	Southern	Govt	Auto	3	23934
2	1001990	Bindu	P00118542	F	26-35	35	1	Uttar Pradesh	Central	Automobile	Auto	3	23924
3	1001425	Sudevi	P00237842	M	0-17	16	0	Karnataka	Southern	Construction	Auto	2	23912
4	1000588	Joni	P00057942	M	26-35	28	1	Gujarat	Western	Food Processing	Auto	2	23877
...
11246	1000695	Manning	P00296942	M	18-25	19	1	Maharashtra	Western	Chemical	Office	4	370
11247	1004089	Reichenbach	P00171342	M	26-35	33	0	Haryana	Northern	Healthcare	Veterinary	3	367
11248	1001209	Oshin	P00201342	F	36-45	40	0	Madhya Pradesh	Central	Textile	Office	4	213
11249	1004023	Noonan	P00059442	M	36-45	37	0	Karnataka	Southern	Agriculture	Office	3	206
11250	1002744	Brumley	P00281742	F	18-25	19	0	Maharashtra	Western	Healthcare	Office	3	188

11239 rows × 13 columns

METHODOLOGY

1. Data loading

Loaded Diwali Sales dataset using Pandas for analysis.

2. Data Inspection

Reviewed structure, data types, and basic statistics of the dataset.

3. Data Cleaning

Removed irrelevant columns , changed column name(marital status -> shaadi) and removed null values

4. Exploratory Data Analysis (EDA)

Used visualizations to study customer behavior, pattern of sales and category trends.

5. Insight Generation

Identified top customer , top states, and best performing product categories.

6. Conclusion

Concluded which age group from which state and profession is contributing more to which product

JOURNEY INSIGHTS

Understanding the raw dataset :

The process began by exploring the structure of the dataset using shape, head, and info. This will help us to identify unnecessary columns , their datatypes , missing values (null values) , number of rows and columns in a dataset .

```
df.shape
```

```
(11251, 15)
```

```
df.head()
```

	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	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

```
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
```


Cleaning and preparing the data :

Irrelevant columns like 'Status', 'unnamed1' were removed , Null values from Amount column is also removed , datatype of Amount column is changed to int64 from float64 , renamed the 'Marital_Status' column to 'Shaadi' . Now our dataset is ready for analysis .

```
#rename column
df.rename(columns= {'Marital_Status':'Shaadi'})
```

	User_ID	Cust_name	Product_ID	Gender	Age Group	Age	Shaadi	State	Zone	Occupation	Product_Category	Orders	Amount
0	1002903	Sanskriti	P00125942	F	26-35	28	0	Maharashtra	Western	Healthcare	Auto	1	23952
1	1000732	Kartik	P00110942	F	26-35	35	1	Andhra Pradesh	Southern	Govt	Auto	3	23934
2	1001990	Bindu	P00118542	F	26-35	35	1	Uttar Pradesh	Central	Automobile	Auto	3	23924
3	1001425	Sudevi	P00237842	M	0-17	16	0	Karnataka	Southern	Construction	Auto	2	23912
4	1000588	Joni	P00057942	M	26-35	28	1	Gujarat	Western	Food Processing	Auto	2	23877
...
11246	1000695	Manning	P00296942	M	18-25	19	1	Maharashtra	Western	Chemical	Office	4	370
11247	1004089	Reichenbach	P00171342	M	26-35	33	0	Haryana	Northern	Healthcare	Veterinary	3	367
11248	1001209	Oshin	P00201342	F	36-45	40	0	Madhya Pradesh	Central	Textile	Office	4	213
11249	1004023	Noonan	P00059442	M	36-45	37	0	Karnataka	Southern	Agriculture	Office	3	206
11250	1002744	Brumley	P00281742	F	18-25	19	0	Maharashtra	Western	Healthcare	Office	3	188

11239 rows × 13 columns

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

```
#check for null values
pd.isnull(df).sum()
```

	0
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

```
# drop null values
df.dropna(inplace=True)
```

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

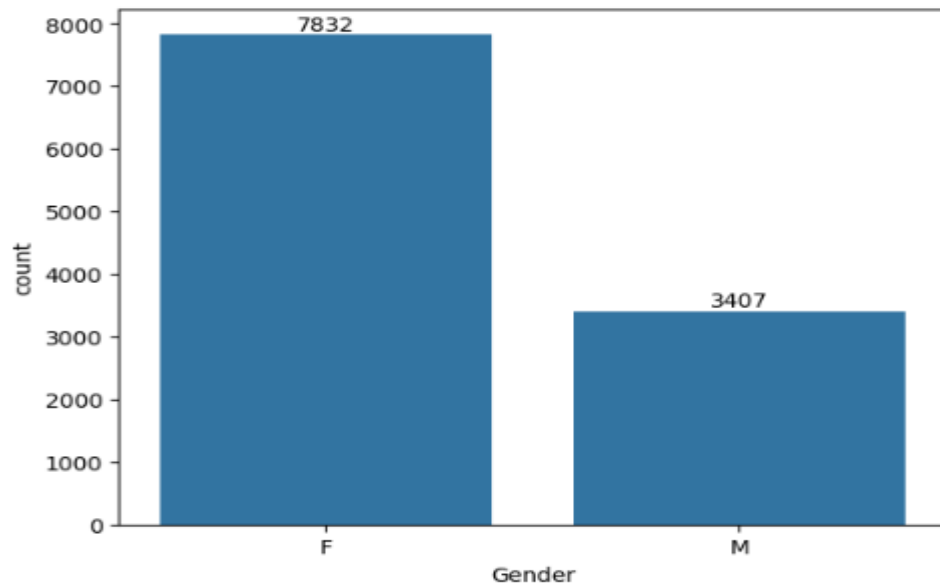
EXPLORATORY DATA ANALYSIS (EDA)

Plotting a bar chart for gender and its count and Gender VS total amount to understand which gender is contributing more and who have more purchasing power .

```
# plotting a bar chart for Gender and it's count
```

```
ax = sns.countplot(x = 'Gender',data = df)
```

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

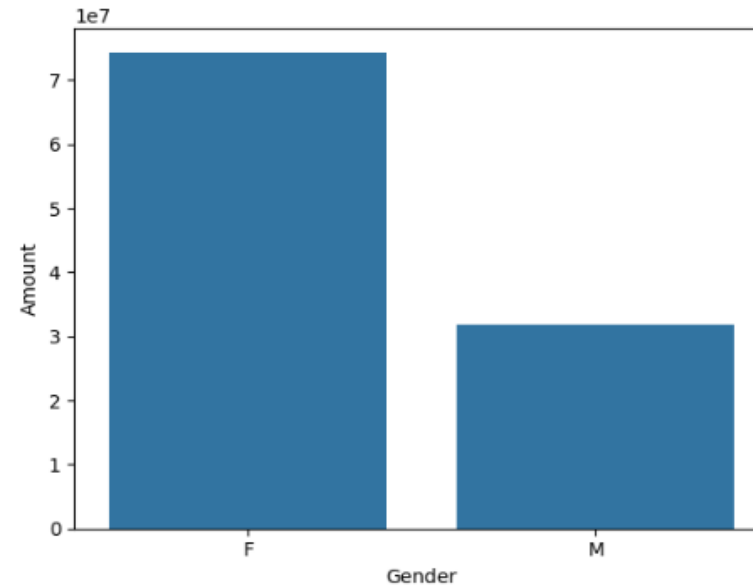


```
# plotting a bar chart for gender vs total amount
```

```
sales_gen = df.groupby(['Gender'], as_index=False)['Amount'].sum().sort_values(by='Amount', ascending=False)
```

```
sns.barplot(x = 'Gender',y= 'Amount' ,data = sales_gen)
```

```
<Axes: xlabel='Gender', ylabel='Amount'>
```

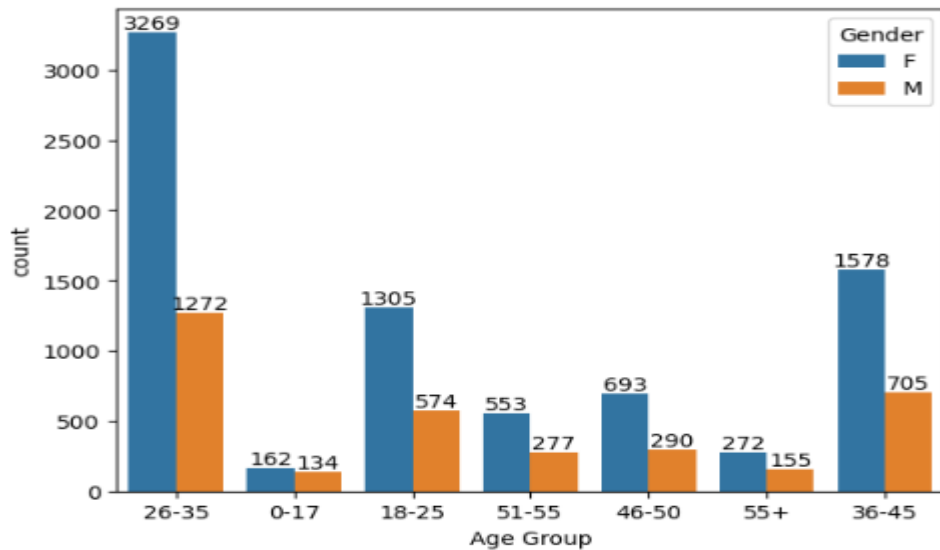


** From above graphs we can see that most of the buyers are females and even the purchasing power of females are greater than men*

Now using a count plot to examine which age group has the highest number of buyers and how purchasing frequency differ between male and female After this a bar plot for total amount VS age group is plotted to understand which age group from which gender is contributing most to revenue . These visualization will reveal which is most active and highest spending customer segment .

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

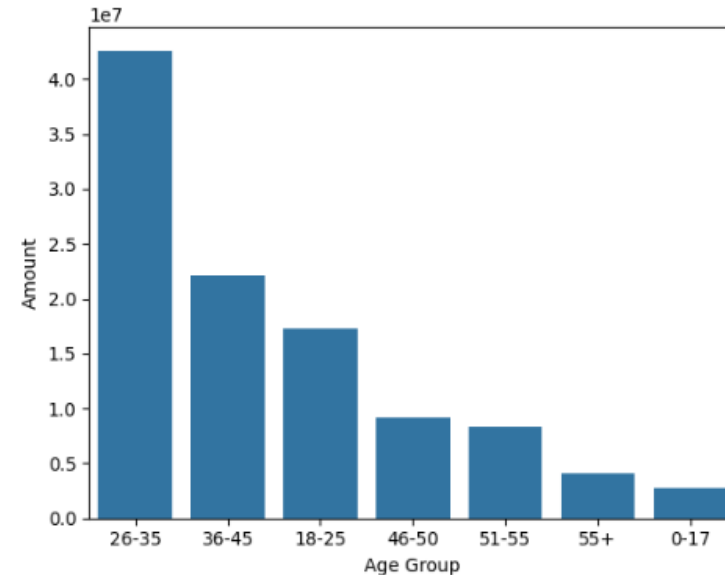
for bars in ax.containers:
    ax.bar_label(bars)
```



```
# 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)
```

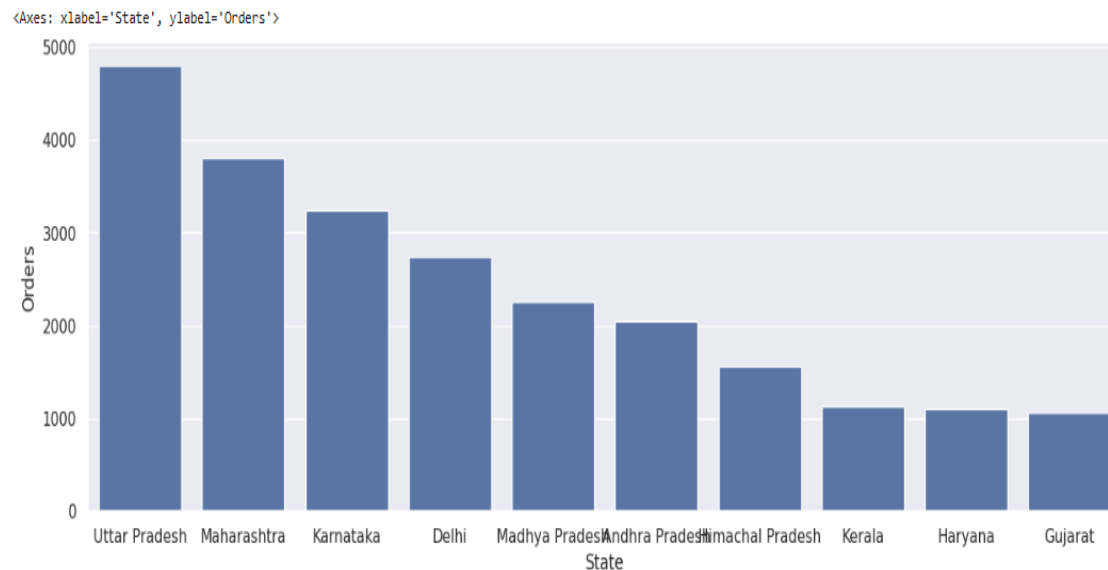
<Axes: xlabel='Age Group', ylabel='Amount'>



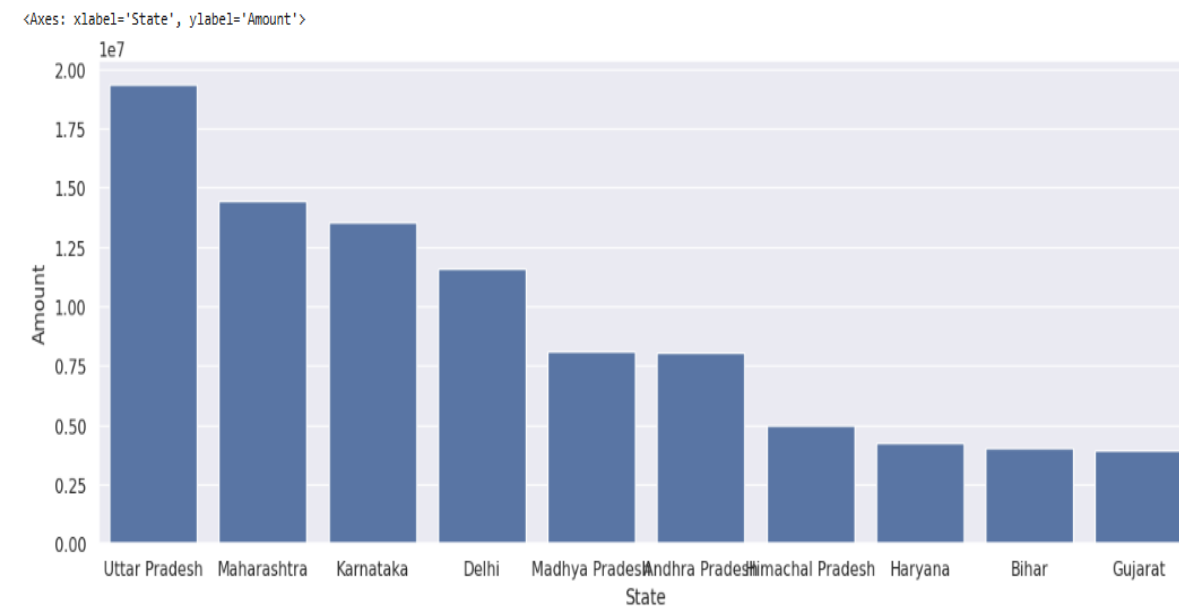
** From above graphs we can see that most of the buyers are of age group between 26-35 yrs female*

Analysing the top 10 states with highest number of orders , then analysing top 10 states generating the highest total revenue . Together this analysis will reveal which region is contributing most to Diwali sales .

```
# 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')
```



```
# 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')
```



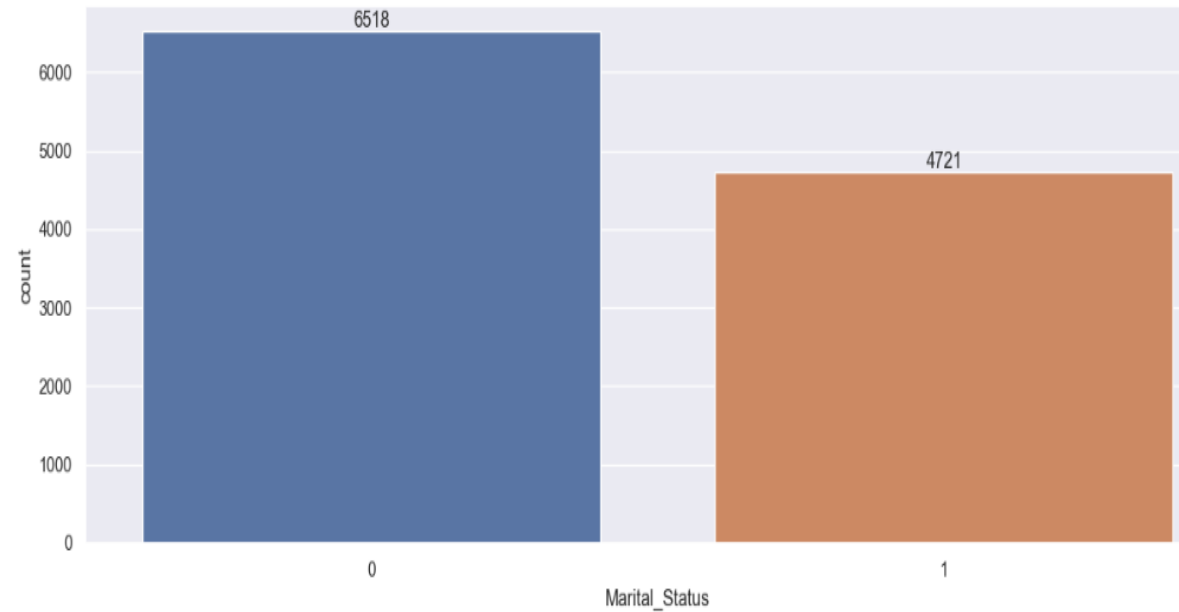
** From above graphs we can see that most of the orders & total sales/amount are from Uttar Pradesh, Maharashtra and Karnataka respectively*

Analysing which marital group places more orders and then identifying which gender marital segment is contributing the highest .

```
ax = sns.countplot(data = df, x = 'Marital_Status')
```

```
sns.set(rc={'figure.figsize':(7,5)})
```

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

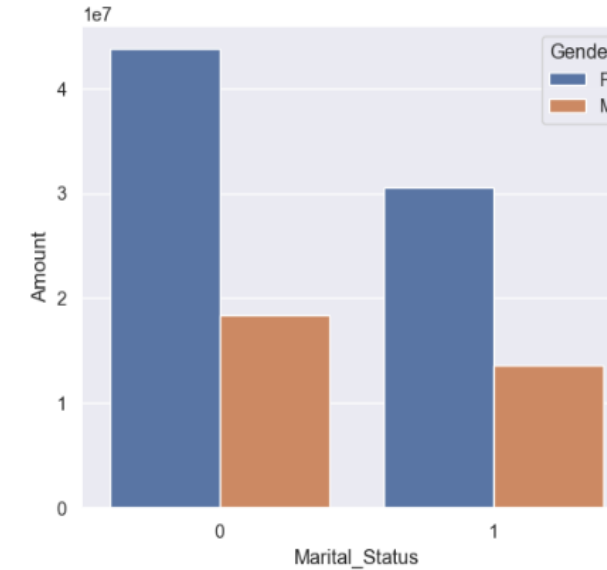


```
sales_state = df.groupby(['Marital_Status', 'Gender'], as_index=False)['Amount'].sum().sort_values(by='Amount', ascending=False)
```

```
sns.set(rc={'figure.figsize':(6,5)})
```

```
sns.barplot(data = sales_state, x = 'Marital_Status', y= 'Amount', hue='Gender')
```

<Axes: xlabel='Marital_Status', ylabel='Amount'>

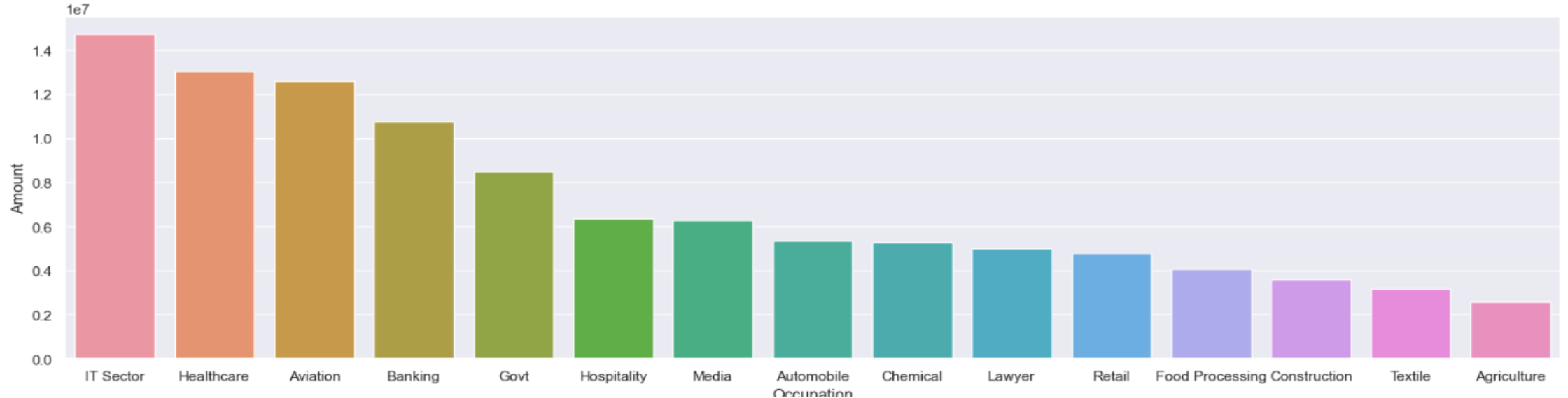


** From above graphs we can see that most of the buyers are married (women) and they have high purchasing power .*

Examining which occupation is generating most number of orders then analysing total revenue contribution by each occupation . This will help us to understand which is high value and high engagement customer segments .

```
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')
```

<Axes: xlabel='Occupation', ylabel='Amount'>



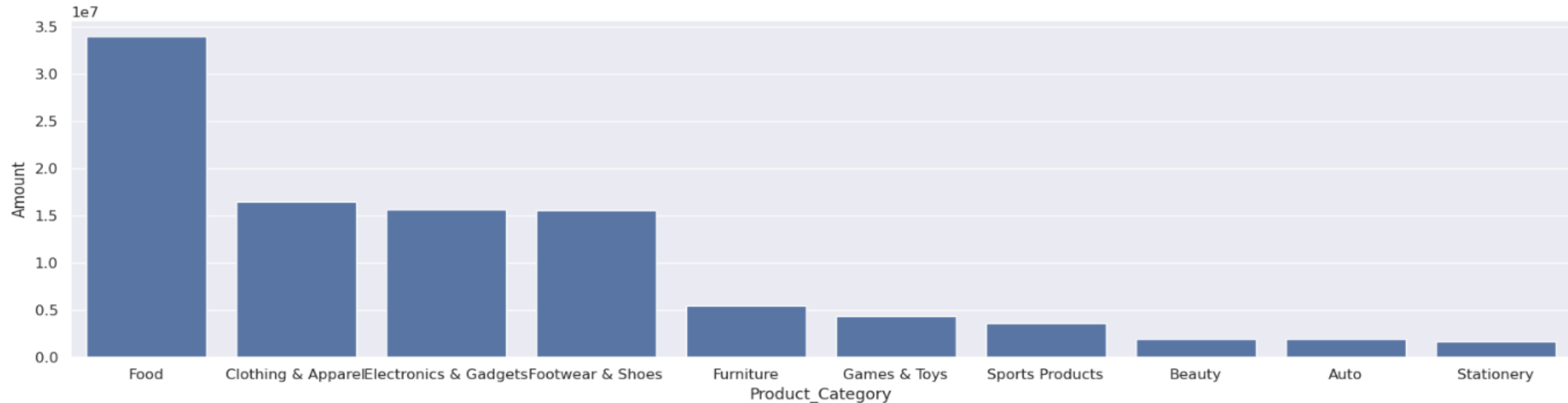
** From above graphs we can see that most of the buyers are working in IT, Healthcare and Aviation sector .*

Analysing how many orders came from each product category and then which product category generated highest revenue . This helps understand both customer preferences and high-earning product segments .

```
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')
```

<Axes: xlabel='Product_Category', ylabel='Amount'>



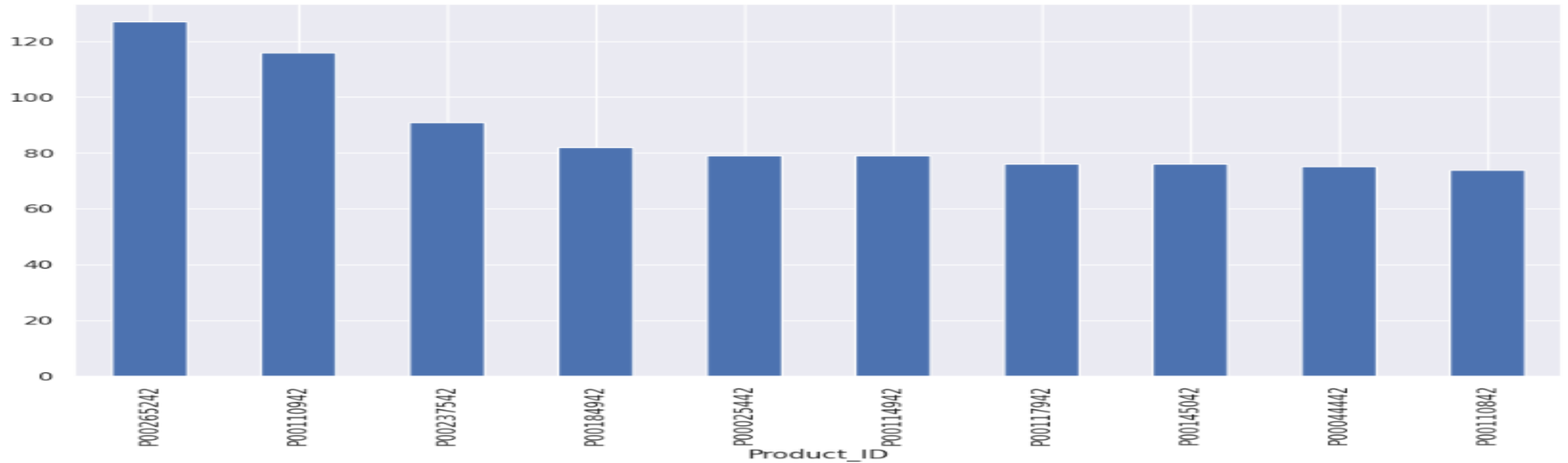
** From above graphs we can see that most of the sold products are from Food, Clothing and Electronics category*

Examining which items were most in demand during Diwali , this will help in planning stock and marketing.

```
# top 10 most sold products (same thing as above)
```

```
fig1, ax1 = plt.subplots(figsize=(12,7))  
df.groupby('Product_ID')['Orders'].sum().nlargest(10).sort_values(ascending=False).plot(kind='bar')
```

```
<Axes: xlabel='Product_ID'>
```



CONCLUSION

The highest sales came from married women aged 26–35, mainly from Uttar Pradesh, Maharashtra, and Karnataka. Most of them worked in IT, Healthcare, or Aviation and preferred buying products from the Food, Clothing, and Electronics categories.