

In [ ]:

In [12]: 

```
# import python libraries
```

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

In [15]: 

```
# import csv file
```

```
df = pd.read_csv(r'C:\Users\user\OneDrive\Desktop\Diwali Sales Data.csv', encoding=
```

In [16]: 

```
df.shape
```

Out[16]: (11251, 15)

In [18]: 

```
df.head(10)
```

Out[18]:

	User_ID	Cust_name	Product_ID	Gender	Age Group	Age	Marital_Status	State
0	1002903	Sanskriti	P00125942	F	26-35	28	0	Maharashtra
1	1000732	Kartik	P00110942	F	26-35	35	1	Andhra Pradesh
2	1001990	Bindu	P00118542	F	26-35	35	1	Uttar Pradesh
3	1001425	Sudevi	P00237842	M	0-17	16	0	Karnataka
4	1000588	Joni	P00057942	M	26-35	28	1	Gujarat
5	1000588	Joni	P00057942	M	26-35	28	1	Himachal Pradesh
6	1001132	Balk	P00018042	F	18-25	25	1	Uttar Pradesh
7	1002092	Shivangi	P00273442	F	55+	61	0	Maharashtra
8	1003224	Kushal	P00205642	M	26-35	35	0	Uttar Pradesh
9	1003650	Ginny	P00031142	F	26-35	26	1	Andhra Pradesh

In [20]: 

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

```
In [25]: 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 [26]: #check for null values
pd.isnull(df)
```

Out[26]:

	User_ID	Cust_name	Product_ID	Gender	Age Group	Age	Marital_Status	State	Zone
0	False	False	False	False	False	False	False	False	False
1	False	False	False	False	False	False	False	False	False
2	False	False	False	False	False	False	False	False	False
3	False	False	False	False	False	False	False	False	False
4	False	False	False	False	False	False	False	False	False
...	...	...	...	...	...	...	...	...	...
11246	False	False	False	False	False	False	False	False	False
11247	False	False	False	False	False	False	False	False	False
11248	False	False	False	False	False	False	False	False	False
11249	False	False	False	False	False	False	False	False	False
11250	False	False	False	False	False	False	False	False	False

11251 rows × 13 columns

In [27]: `#check for null values  
pd.isnull(df).sum()`

Out[27]:

In [28]: `df.shape`

Out[28]: (11251, 13)

In [29]: `df.dropna(inplace=True)`In [30]: `df.shape`

Out[30]: (11239, 13)

```
In [31]: pd.isnull(df).sum()
```

```
Out[31]: 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 [32]: df.info()
```

```
<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
```

```
In [33]: # drop null values
df.dropna(inplace=True)
```

```
In [34]: df.shape
```

```
Out[34]: (11239, 13)
```

```
In [35]: # change data type
df['Amount'] = df['Amount'].astype('int')
```

```
In [36]: df['Amount'].dtypes
```

```
Out[36]: dtype('int32')
```

In [37]: `df.columns`

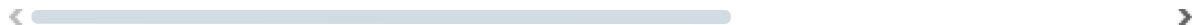
Out[37]: `Index(['User_ID', 'Cust_name', 'Product_ID', 'Gender', 'Age Group', 'Age', 'Marital_Status', 'State', 'Zone', 'Occupation', 'Product_Category', 'Orders', 'Amount'], dtype='object')`

In [38]: `#rename column  
df.rename(columns= {'Marital_Status': 'Shaadi'})`

Out[38]:

	User_ID	Cust_name	Product_ID	Gender	Age Group	Age	Shaadi	State
<b>0</b>	1002903	Sanskriti	P00125942	F	26-35	28	0	Maharashtra \
<b>1</b>	1000732	Kartik	P00110942	F	26-35	35	1	Andhra Pradesh \ S
<b>2</b>	1001990	Bindu	P00118542	F	26-35	35	1	Uttar Pradesh
<b>3</b>	1001425	Sudevi	P00237842	M	0-17	16	0	Karnataka \ S
<b>4</b>	1000588	Joni	P00057942	M	26-35	28	1	Gujarat \
...	...	...	...	...	...	...	...	...
<b>11246</b>	1000695	Manning	P00296942	M	18-25	19	1	Maharashtra \
<b>11247</b>	1004089	Reichenbach	P00171342	M	26-35	33	0	Haryana \ N
<b>11248</b>	1001209	Oshin	P00201342	F	36-45	40	0	Madhya Pradesh
<b>11249</b>	1004023	Noonan	P00059442	M	36-45	37	0	Karnataka \ S
<b>11250</b>	1002744	Brumley	P00281742	F	18-25	19	0	Maharashtra \

11239 rows × 13 columns



In [39]: `# describe() method returns description of the data in the DataFrame (i.e. count, mean etc.)  
df.describe()`

Out[39]:

	User_ID	Age	Marital_Status	Orders	Amount
<b>count</b>	1.123900e+04	11239.000000	11239.000000	11239.000000	11239.000000
<b>mean</b>	1.003004e+06	35.410357	0.420055	2.489634	9453.610553
<b>std</b>	1.716039e+03	12.753866	0.493589	1.114967	5222.355168
<b>min</b>	1.000001e+06	12.000000	0.000000	1.000000	188.000000
<b>25%</b>	1.001492e+06	27.000000	0.000000	2.000000	5443.000000
<b>50%</b>	1.003064e+06	33.000000	0.000000	2.000000	8109.000000
<b>75%</b>	1.004426e+06	43.000000	1.000000	3.000000	12675.000000
<b>max</b>	1.006040e+06	92.000000	1.000000	4.000000	23952.000000

In [40]:

```
# use describe() for specific columns
df[['Age', 'Orders', 'Amount']].describe()
```

Out[40]:

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

# Exploratory Data Analysis

## Gender

In [41]:

```
df.columns
```

Out[41]:

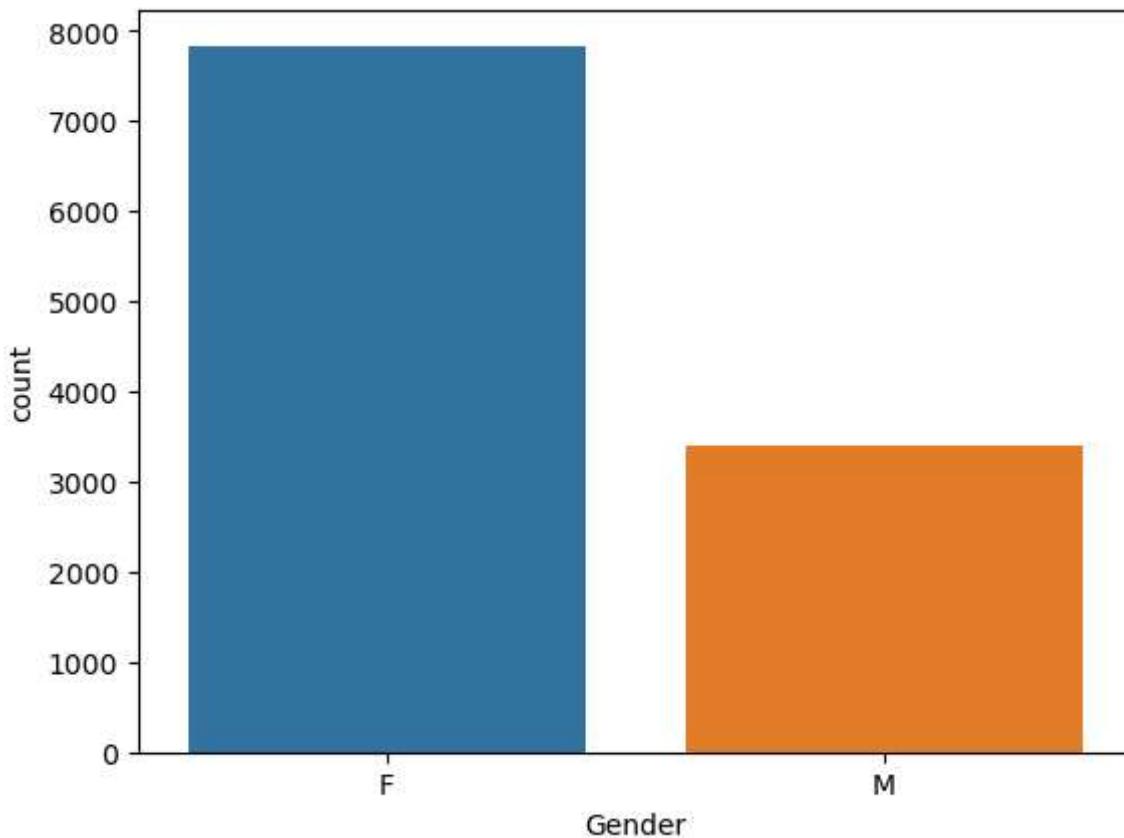
```
Index(['User_ID', 'Cust_name', 'Product_ID', 'Gender', 'Age Group', 'Age',
       'Marital_Status', 'State', 'Zone', 'Occupation', 'Product_Category',
       'Orders', 'Amount'],
      dtype='object')
```

In [42]:

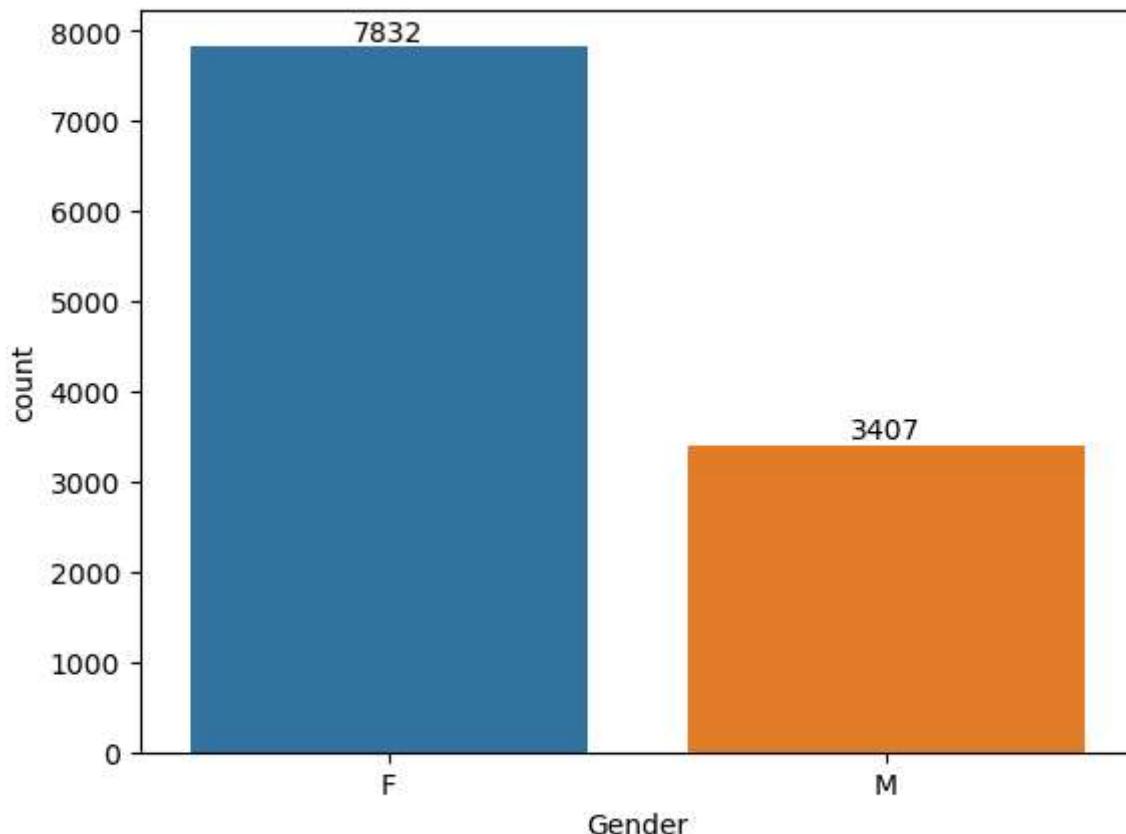
```
sns.countplot(x = 'Gender', data = df)
```

Out[42]:

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



```
In [43]: # 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)
```



```
In [44]: df.groupby(['Gender'], as_index=False)['Amount'].sum().sort_values(by='Amount', asc
```

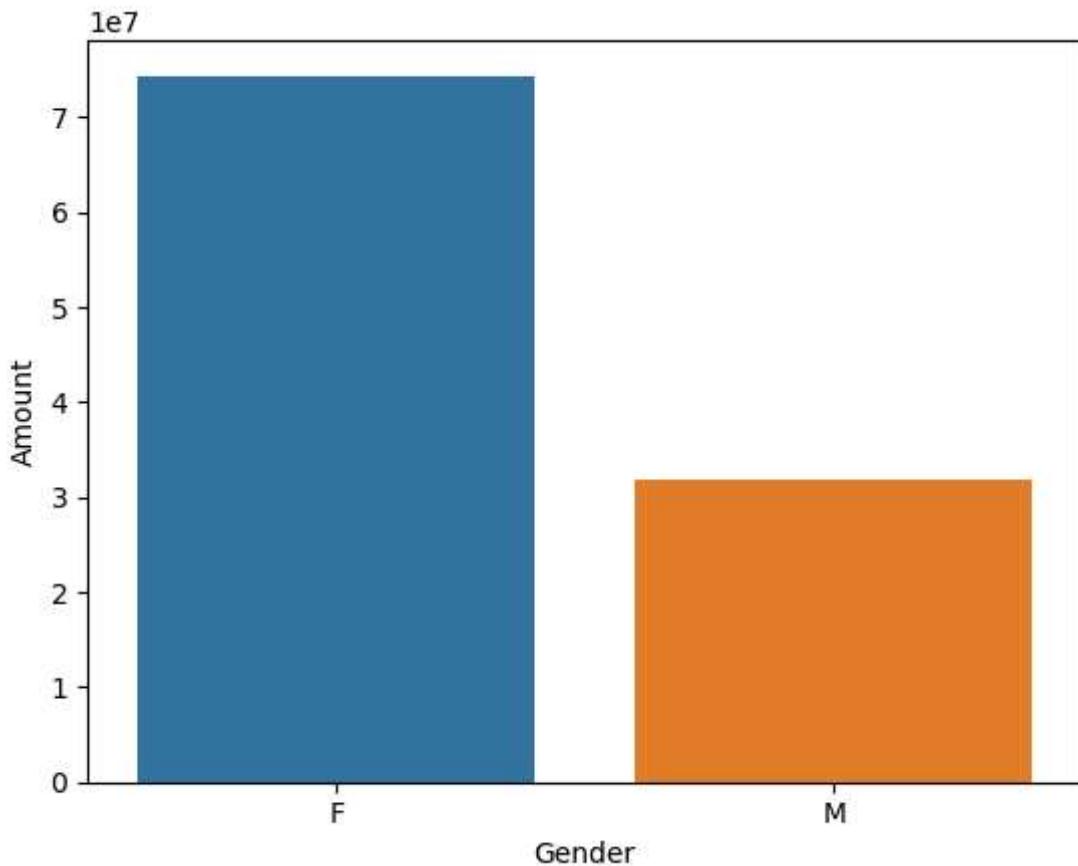
```
Out[44]:
```

	Gender	Amount
0	F	74335853
1	M	31913276

```
In [45]: # plotting a bar chart for gender vs total amount
```

```
sales_gen = df.groupby(['Gender'], as_index=False)['Amount'].sum().sort_values(by='Amount', asc)
sns.barplot(x = 'Gender', y= 'Amount' ,data = sales_gen)
```

```
Out[45]: <Axes: xlabel='Gender', ylabel='Amount'>
```

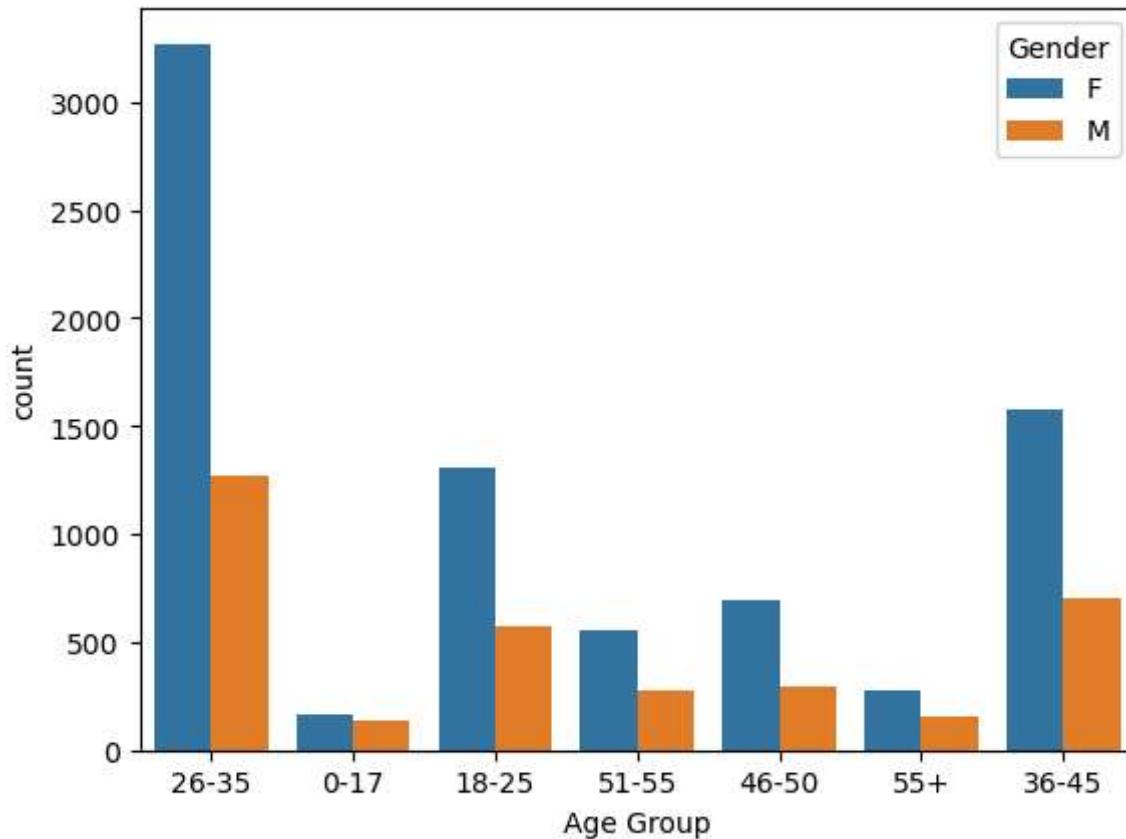


```
In [46]: df.columns
```

```
Out[46]: Index(['User_ID', 'Cust_name', 'Product_ID', 'Gender', 'Age Group', 'Age',  
       'Marital_Status', 'State', 'Zone', 'Occupation', 'Product_Category',  
       'Orders', 'Amount'],  
      dtype='object')
```

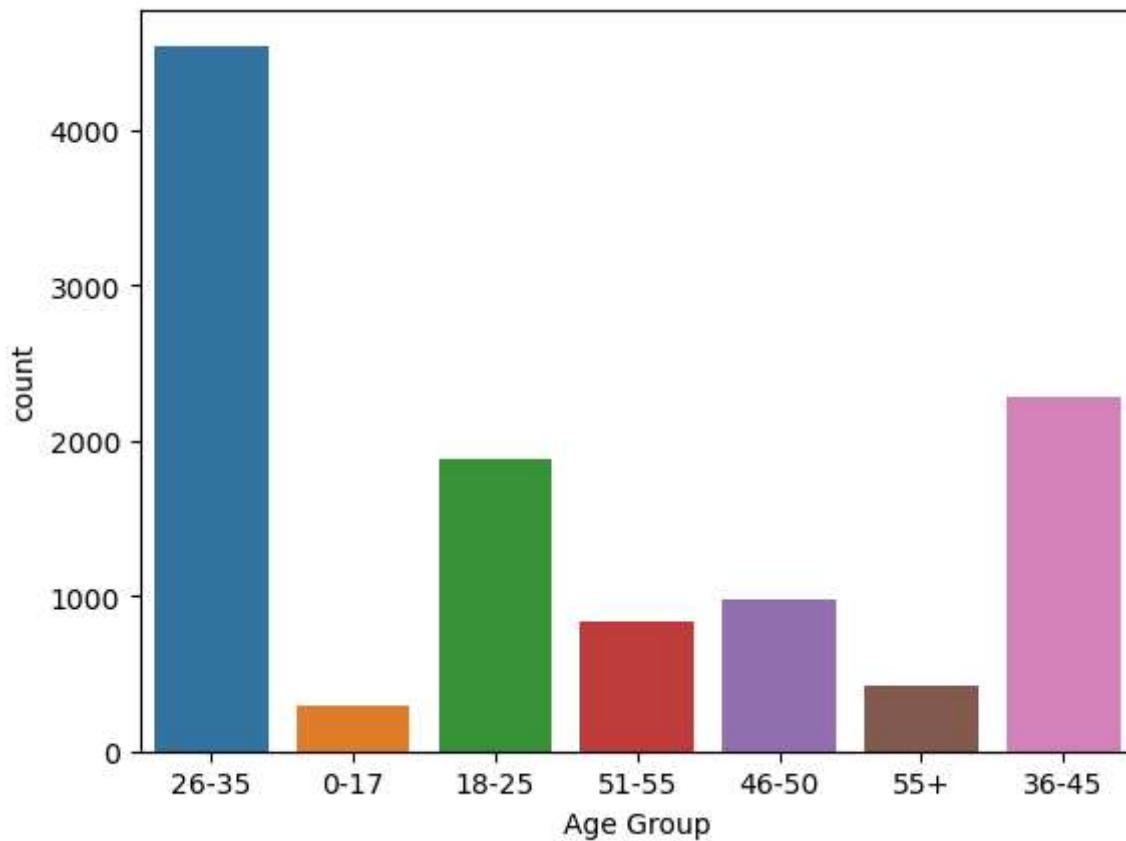
```
In [47]: sns.countplot(data = df, x = 'Age Group', hue = 'Gender')
```

```
Out[47]: <Axes: xlabel='Age Group', ylabel='count'>
```



```
In [48]: sns.countplot(data = df, x = 'Age Group')
```

```
Out[48]: <Axes: xlabel='Age Group', ylabel='count'>
```

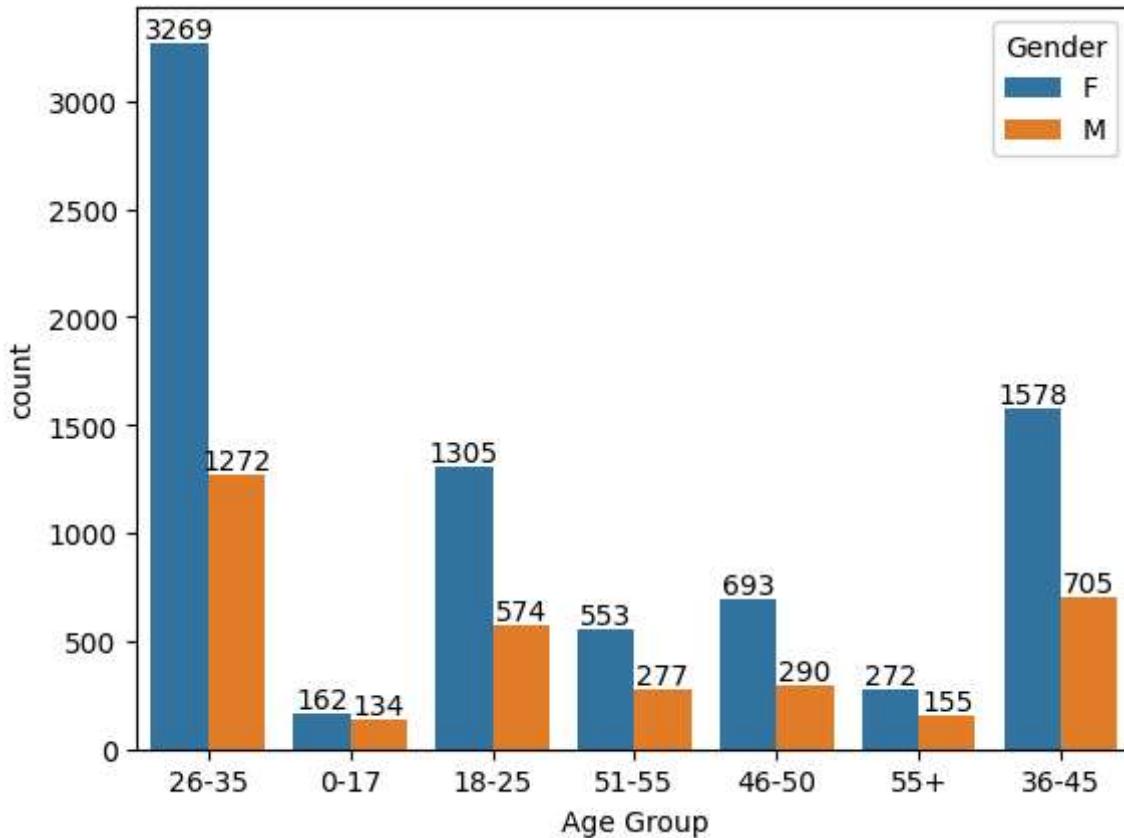


From above graphs we can see that most of the buyers are females and even the purchasing power of females are greater than men

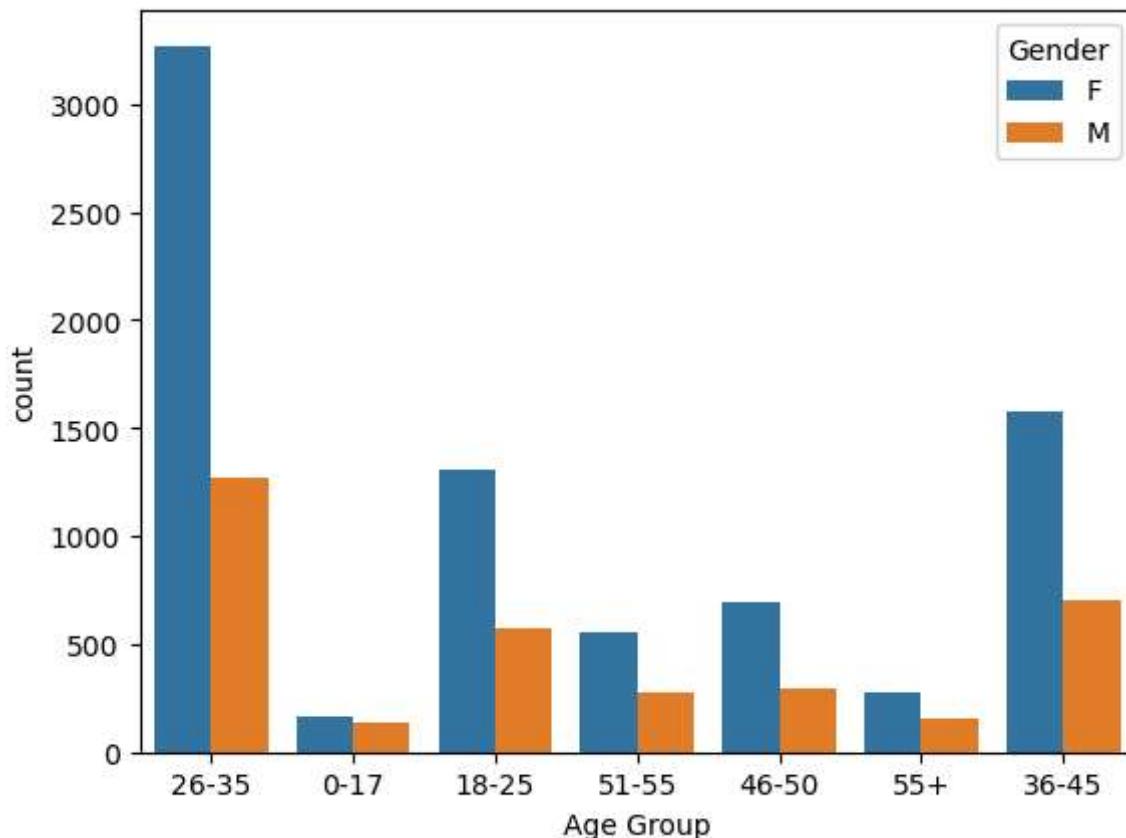
## Age

```
In [17]: ax = sns.countplot(data = df, x = 'Age Group', hue = 'Gender')

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

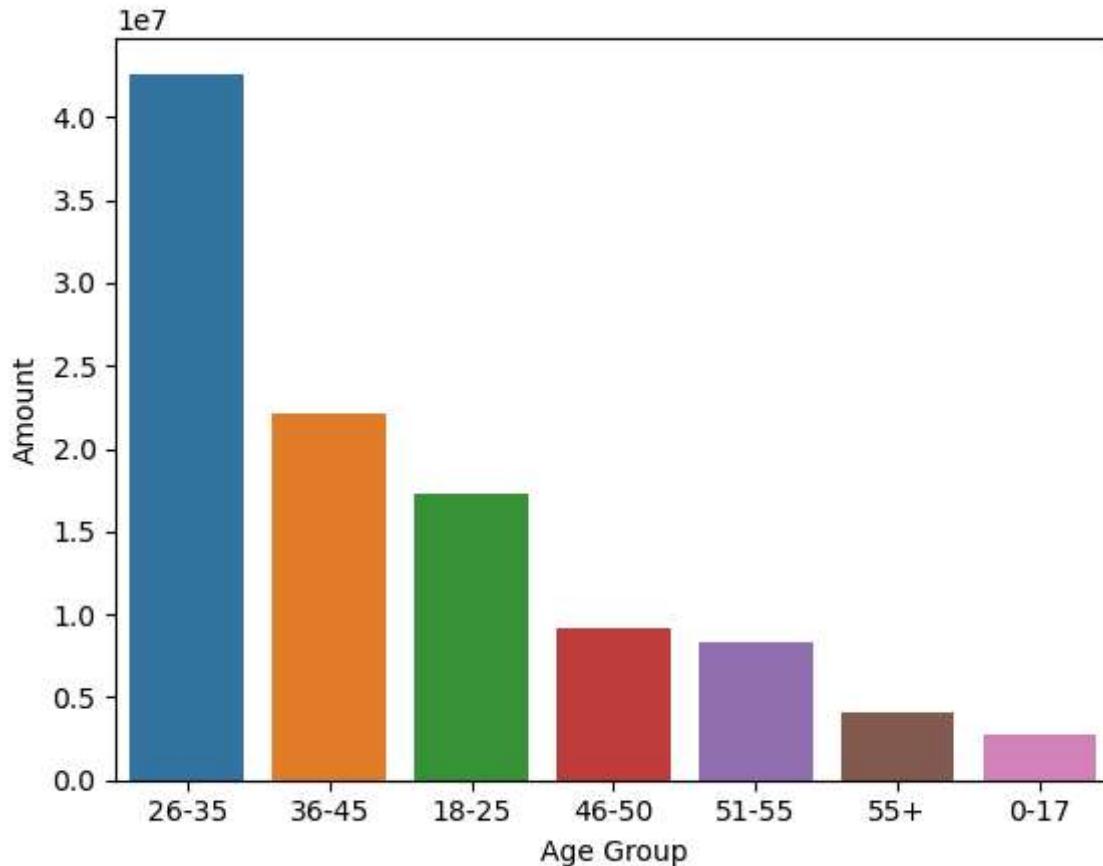


```
In [36]: ax = sns.countplot(data = df, x = 'Age Group', hue = 'Gender')
```



```
In [49]: # 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[49]: <Axes: xlabel='Age Group', ylabel='Amount'>
```



In [50]: `df.columns`

Out[50]: `Index(['User_ID', 'Cust_name', 'Product_ID', 'Gender', 'Age Group', 'Age', 'Marital_Status', 'State', 'Zone', 'Occupation', 'Product_Category', 'Orders', 'Amount'],  
dtype='object')`

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

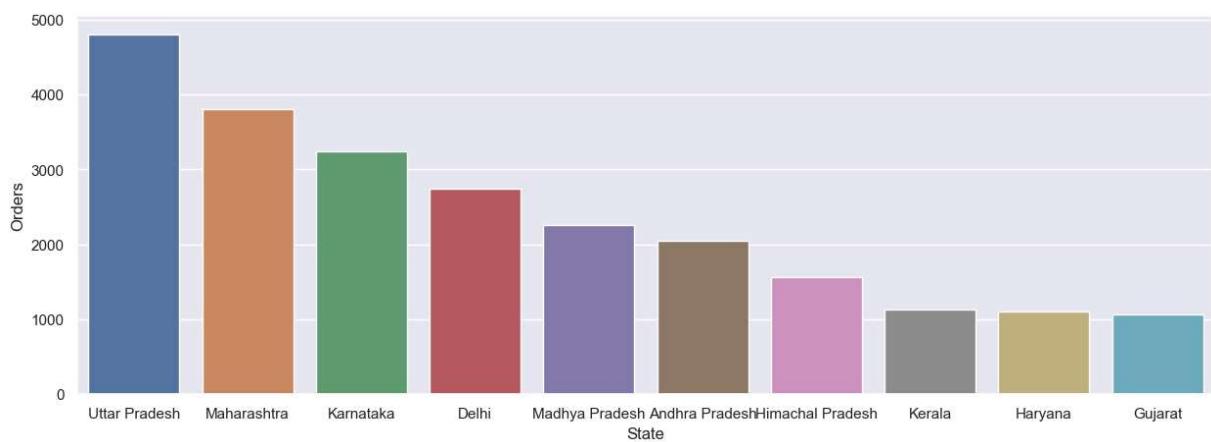
## State

In [51]: `# 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[51]: <Axes: xlabel='State', ylabel='Orders'>

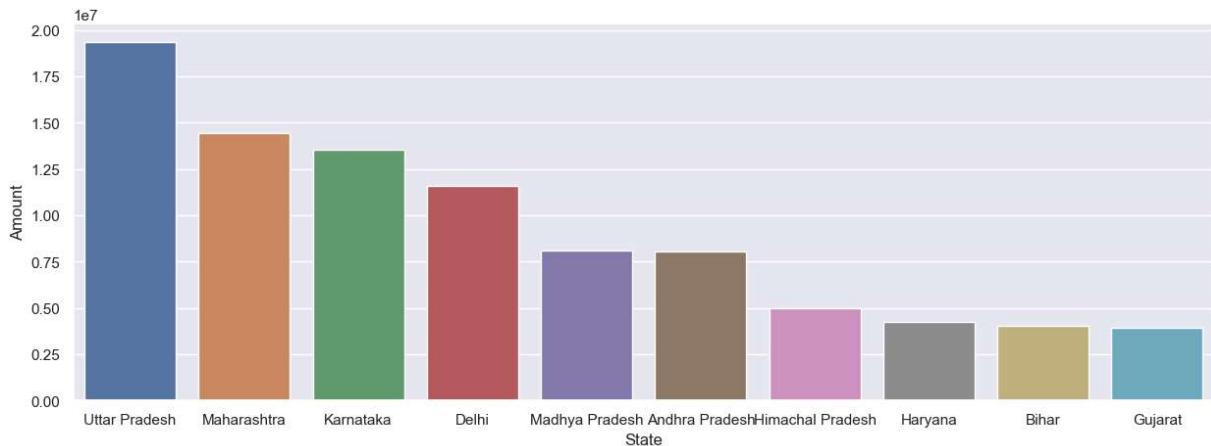


```
In [52]: # total amount/sales from top 10 states
```

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

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

```
Out[52]: <Axes: xlabel='State', ylabel='Amount'>
```

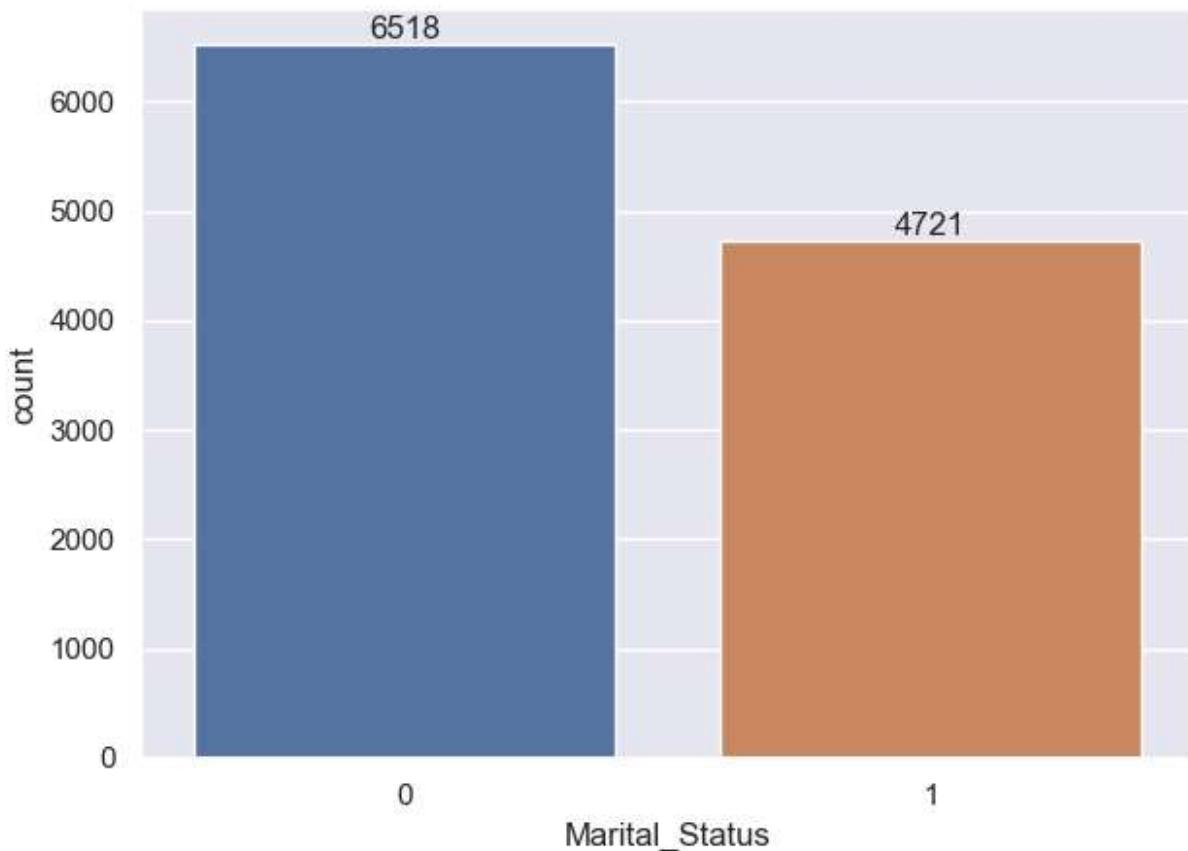


From above graphs we can see that most of the orders & total sales/amount are from Uttar Pradesh, Maharashtra and Karnataka respectively

## Marital Status

```
In [54]: ax = sns.countplot(data = df, x = 'Marital_Status')

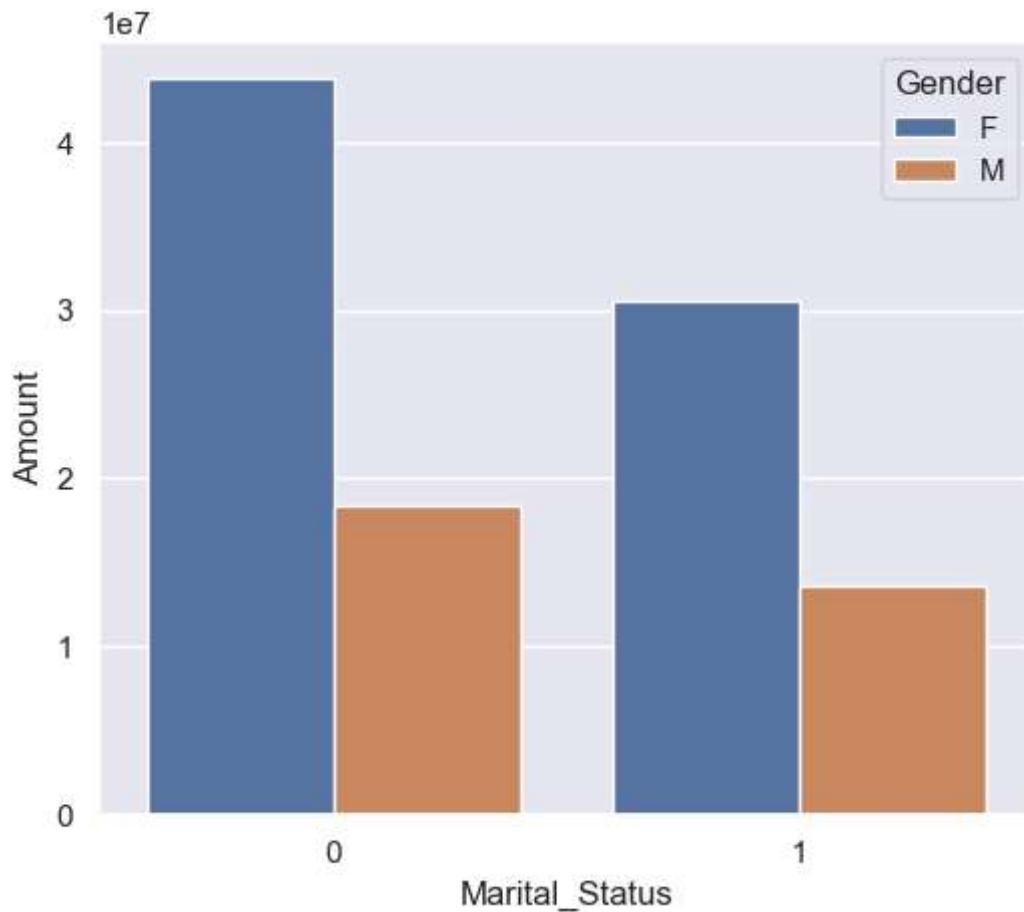
sns.set(rc={'figure.figsize':(7,5)})
for bars in ax.containers:
    ax.bar_label(bars)
```



```
In [55]: sales_state = df.groupby(['Marital_Status', 'Gender'], as_index=False)[['Amount']].sum()

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

```
Out[55]: <Axes: xlabel='Marital_Status', ylabel='Amount'>
```

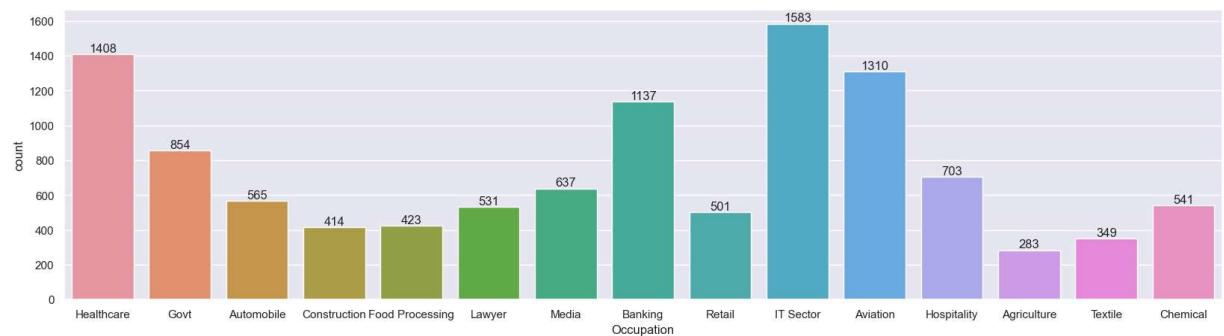


From above graphs we can see that most of the buyers are married (women) and they have high purchasing power

## Occupation

```
In [56]: sns.set(rc={'figure.figsize':(20,5)})
ax = sns.countplot(data = df, x = 'Occupation')

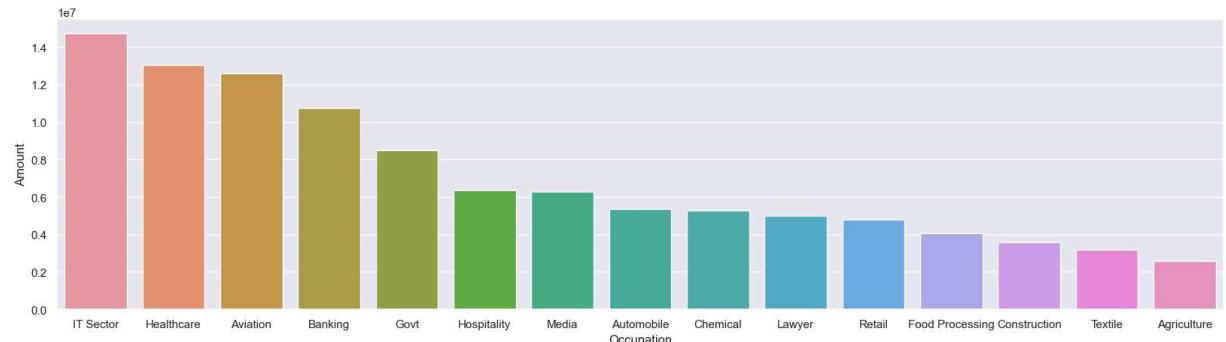
for bars in ax.containers:
    ax.bar_label(bars)
```



```
In [57]: sales_state = df.groupby(['Occupation'], as_index=False)[['Amount']].sum().sort_values

sns.set(rc={'figure.figsize':(20,5)})
sns.barplot(data = sales_state, x = 'Occupation',y= 'Amount')
```

Out[57]: <Axes: xlabel='Occupation', ylabel='Amount'>

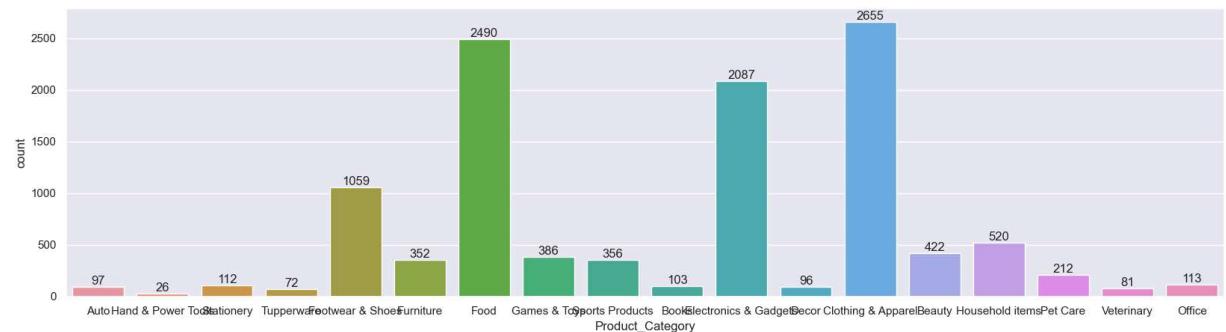


From above graphs we can see that most of the buyers are working in IT, Healthcare and Aviation sector

## Product Category

```
In [58]: sns.set(rc={'figure.figsize':(20,5)})
ax = sns.countplot(data = df, x = 'Product_Category')

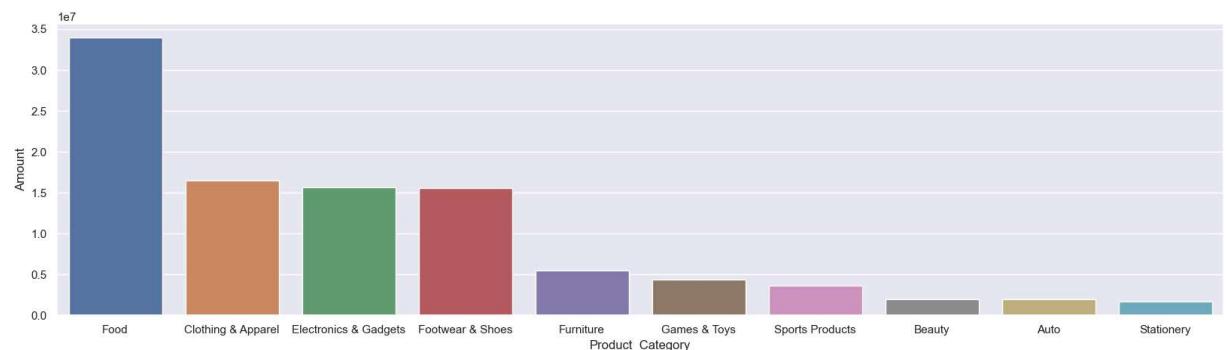
for bars in ax.containers:
    ax.bar_label(bars)
```



```
In [26]: sales_state = df.groupby(['Product_Category'], as_index=False)[['Amount']].sum().sort

sns.set(rc={'figure.figsize':(20,5)})
sns.barplot(data = sales_state, x = 'Product_Category',y= 'Amount')
```

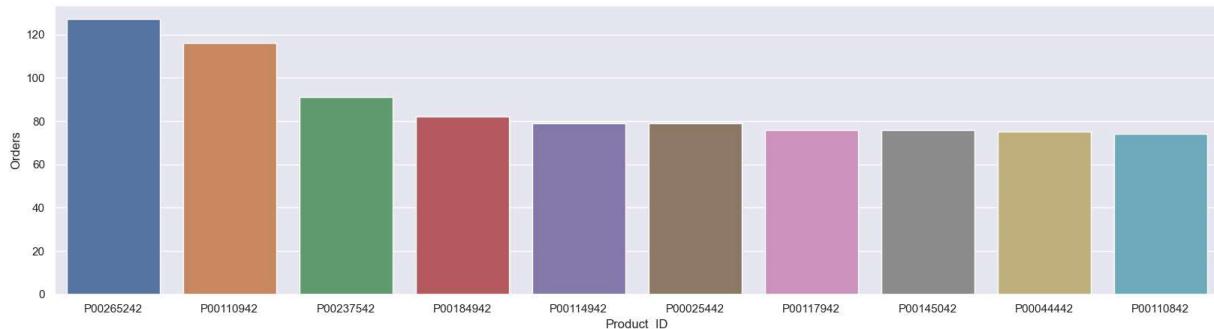
Out[26]: <Axes: xlabel='Product\_Category', ylabel='Amount'>



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

```
In [60]: sales_state = df.groupby(['Product_ID'], as_index=False)['Orders'].sum().sort_values(sns.set(rc={'figure.figsize':(20,5)})sns.barplot(data = sales_state, x = 'Product_ID',y= 'Orders')
```

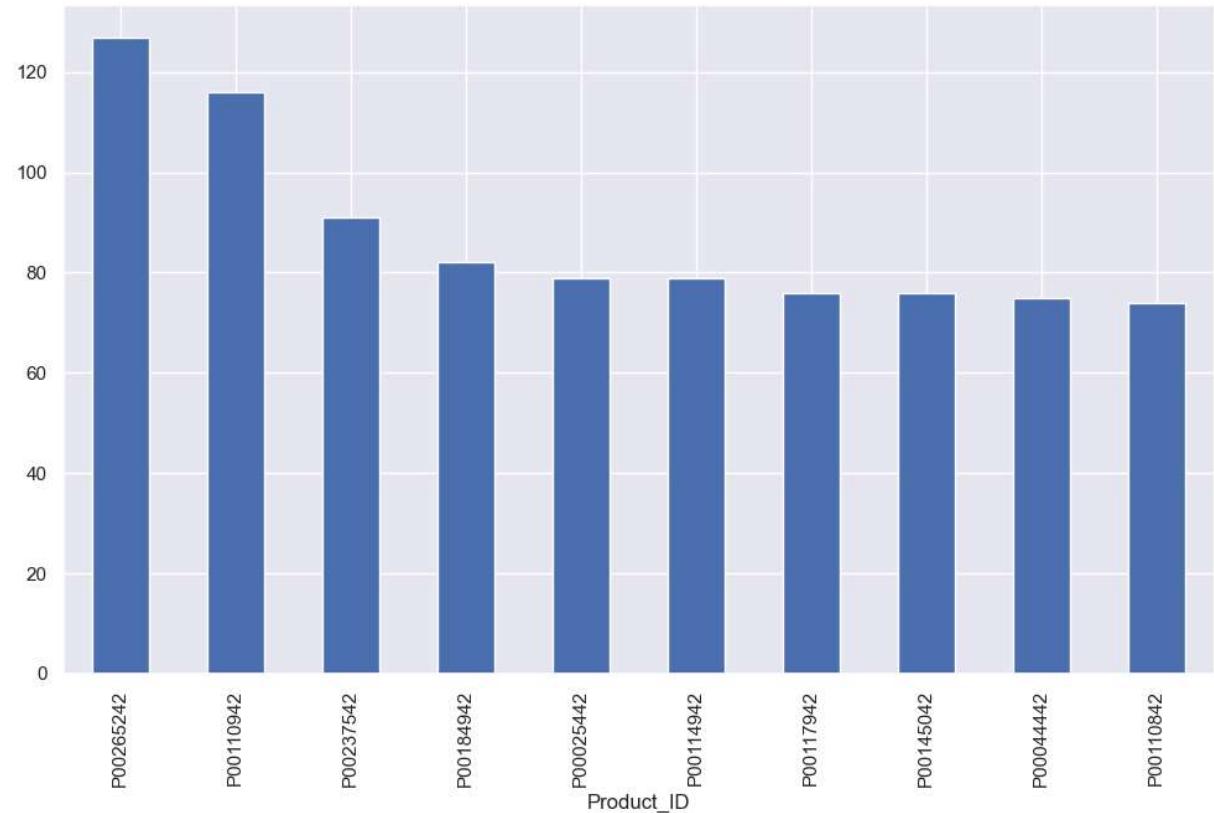
```
Out[60]: <Axes: xlabel='Product_ID', ylabel='Orders'>
```



```
In [61]: # 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).
```

```
Out[61]: <Axes: xlabel='Product_ID'>
```



## Conclusion:

*Married women age group 26-35 yrs from UP, Maharashtra and Karnataka working in IT, Healthcare and Aviation are more likely to buy products from Food, Clothing and Electronics category*

Thank you!