

```
In [5]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
%matplotlib inline
import seaborn as sns
```

```
In [6]: df = pd.read_csv("Diwali Sales Dataset.csv")
df
```

```
Out[6]:
```

	User_ID	Cust_name	Product_ID	Gender	Age Group	Age	Marital_Status	State	Zone	Occupation	Product_Category	Orders	Amount	Status
0	1002903	Sanskriti	P00125942	F	26-35	28	0	Maharashtra	Western	Healthcare	Auto	1	23952.0	NaN
1	1000732	Karrik	P00110942	F	26-35	35	1	Andhra Pradesh	Southern	Govt	Auto	3	23934.0	NaN
2	1001990	Bindu	P00118542	F	26-35	35	1	Uttar Pradesh	Central	Automobile	Auto	3	23924.0	NaN
3	1001425	Sudevi	P00237842	M	0-17	16	0	Karnataka	Southern	Construction	Auto	2	23912.0	NaN
4	1000588	Joni	P00057942	M	26-35	28	1	Gujarat	Western	Food Processing	Auto	2	23877.0	NaN
...
11246	1000695	Manning	P00296942	M	18-25	19	1	Maharashtra	Western	Chemical	Office	4	370.0	NaN
11247	1004089	Reichenbach	P00171342	M	26-35	33	0	Haryana	Northern	Healthcare	Veterinary	3	367.0	NaN
11248	1001209	Oshin	P00201342	F	36-45	40	0	Madhya Pradesh	Central	Textile	Office	4	213.0	NaN
11249	1004023	Noonan	P00059442	M	36-45	37	0	Karnataka	Southern	Agriculture	Office	3	206.0	NaN
11250	1002744	Brunley	P00281742	F	18-25	19	0	Maharashtra	Western	Healthcare	Office	3	188.0	NaN

11251 rows × 15 columns

```
In [5]: 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 [6]: df.describe()
```

```
Out[6]:
```

	User_ID	Age	Marital_Status	Orders	Amount	Status	unnamed1
count	1.125100e+04	11251.000000	11251.000000	11251.000000	11239.000000	0.0	0.0
mean	1.003004e+06	35.421207	0.420318	2.489290	9453.610858	NaN	NaN
std	1.716125e+03	12.754122	0.493632	1.115047	5222.355869	NaN	NaN
min	1.000001e+06	12.000000	0.000000	1.000000	188.000000	NaN	NaN
25%	1.001492e+06	27.000000	0.000000	1.500000	5443.000000	NaN	NaN
50%	1.003065e+06	33.000000	0.000000	2.000000	8109.000000	NaN	NaN
75%	1.004430e+06	43.000000	1.000000	3.000000	12675.000000	NaN	NaN
max	1.006040e+06	92.000000	1.000000	4.000000	23952.000000	NaN	NaN

```
In [7]: df.isnull().sum()
```

```
Out[7]:
User_ID      0
Cust_name    0
Product_ID   0
Gender        0
Age Group    0
Marital_Status  0
State         0
Zone          0
Occupation    0
Product_Category  0
Orders        0
Amount       12
Status        0
unnamed1     11251
dtype: int64
```

Data Cleaning:

Step to DROP blank /unrelated columns

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

```
In [9]: df
```

```
Out[9]:
```

	User_ID	Cust_name	Product_ID	Gender	Age Group	Age	Marital_Status	State	Zone	Occupation	Product_Category	Orders	Amount
0	1002903	Sanskriti	P00125942	F	26-35	28	0	Maharashtra	Western	Healthcare	Auto	1	23952.0
1	1000732	Karrik	P00110942	F	26-35	35	1	Andhra Pradesh	Southern	Govt	Auto	3	23934.0
2	1001990	Bindu	P00118542	F	26-35	35	1	Uttar Pradesh	Central	Automobile	Auto	3	23924.0
3	1001425	Sudevi	P00237842	M	0-17	16	0	Karnataka	Southern	Construction	Auto	2	23912.0
4	1000588	Joni	P00057942	M	26-35	28	1	Gujarat	Western	Food Processing	Auto	2	23877.0
...
11246	1000695	Manning	P00296942	M	18-25	19	1	Maharashtra	Western	Chemical	Office	4	370.0
11247	1004089	Reichenbach	P00171342	M	26-35	33	0	Haryana	Northern	Healthcare	Veterinary	3	367.0
11248	1001209	Oshin	P00201342	F	36-45	40	0	Madhya Pradesh	Central	Textile	Office	4	213.0
11249	1004023	Noonan	P00059442	M	36-45	37	0	Karnataka	Southern	Agriculture	Office	3	206.0
11250	1002744	Brunley	P00281742	F	18-25	19	0	Maharashtra	Western	Healthcare	Office	3	188.0

11251 rows × 13 columns

```
In [11]: pd.isnull(df).sum()
```

```
Out[11]:
User_ID      0
Cust_name    0
Product_ID   0
Gender        0
Age Group    0
Marital_Status  0
State         0
Zone          0
Occupation    0
Product_Category  0
Orders        0
Amount       12
dtype: int64
```

```
In [13]: df.shape
```

```
Out[13]: (11251, 13)
```

DROP null values

```
In [14]: df.dropna(inplace = True)
```

```
In [15]: df.shape
```

```
Out[15]: (11239, 13)
```

```
In [16]: df.isnull().sum()
```

```
Out[16]:
User_ID      0
Cust_name    0
Product_ID   0
Gender        0
Age Group    0
Marital_Status  0
State         0
Zone          0
Occupation    0
Product_Category  0
Orders        0
Amount       0
dtype: int64
```

Exploratory Data Analysis

Gender

```
In [21]: ax = sns.countplot(x = 'Gender', data = df)
```



```
In [28]: df.groupby(['Gender'], as_index = False)['Amount'].sum().sort_values(by = 'Amount', ascending = False)
```

Gender	Amount
0	F 74335866.43
1	M 31913276.00

```
In [30]: 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)
```

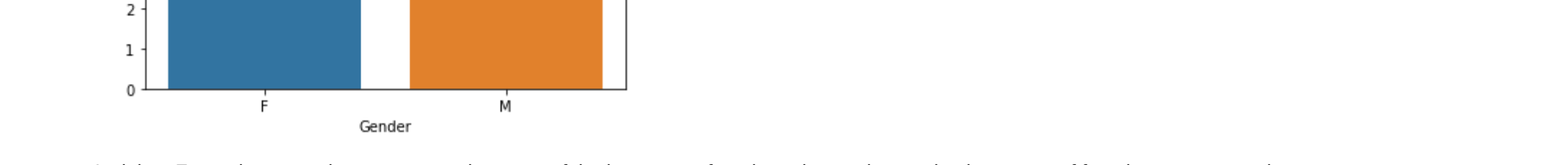
```
Out[30]: <AxesSubplot: xlabel='Gender', ylabel='Amount'>
```

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

Age

```
In [31]: sns.countplot(data = df, x = 'Age Group', hue = 'Gender')
```

```
Out[31]: <AxesSubplot: xlabel='Age Group', ylabel='count'>
```



Total Amount vs Age Group

```
In [25]: 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[25]: <AxesSubplot: xlabel='Age Group', ylabel='Amount'>
```



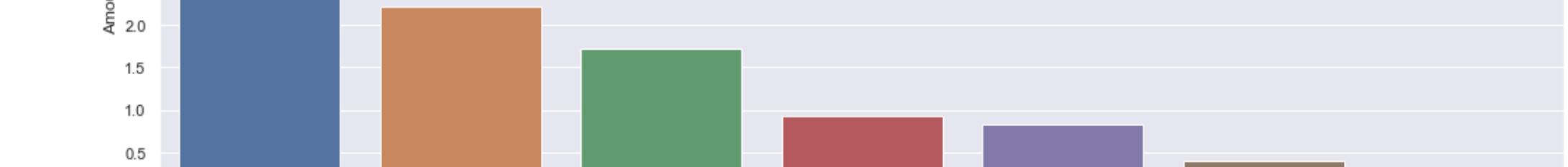
Insights: From the above graphs we can see that most of the buyers are of the age group between 26-35 yrs female

State

Total number of orders from Top 10 states

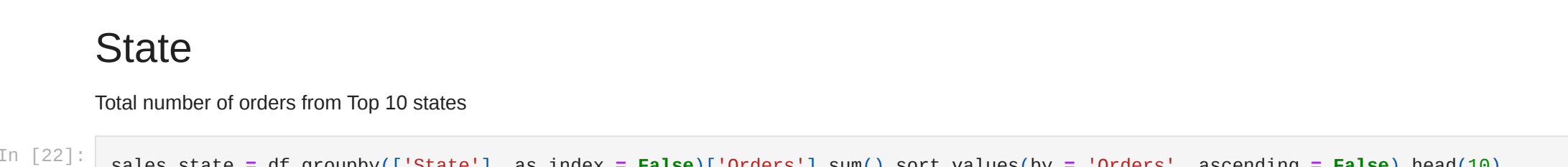
```
In [22]: sales_state = df.groupby(['State'], as_index = False)['Orders'].sum().sort_values(by = 'Orders', ascending = False).head(10)
sns.set(rc = {'figure.figsize' : (18,5)})
sns.barplot(data = sales_state, x = 'State', y = 'Orders')
```

```
Out[22]: <AxesSubplot: xlabel='State', ylabel='Orders'>
```



```
In [24]: #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' : (18,5)})
sns.barplot(data = sales_state, x = 'State', y = 'Amount')
```

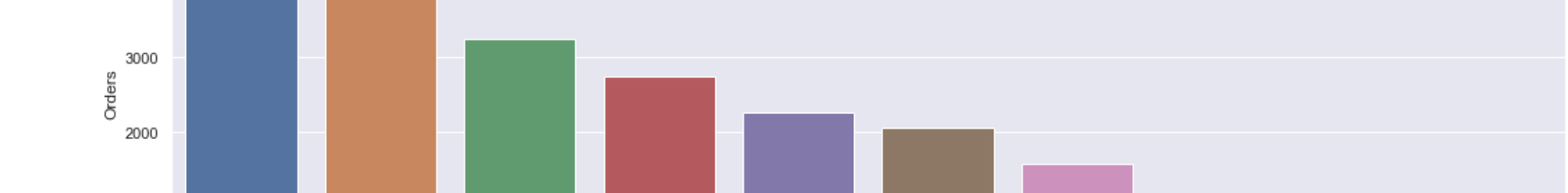
```
Out[24]: <AxesSubplot: xlabel='State', ylabel='Amount'>
```



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

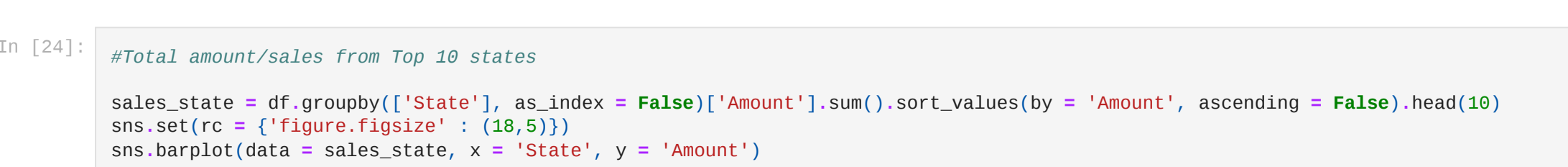
Marital Status

```
In [43]: ax = sns.countplot(data = df, x = 'Marital_Status')
sns.set(rc = {'figure.figsize' : (5,5)})
```



```
In [30]: sales_state = df.groupby(['Marital_Status', 'Gender'], as_index = False)['Amount'].sum().sort_values(by = 'Amount', ascending = False)
sns.set(rc = {'figure.figsize' : (18,5)})
sns.barplot(data = sales_state, x = 'Marital_Status', y = 'Amount', hue = 'Gender')
```

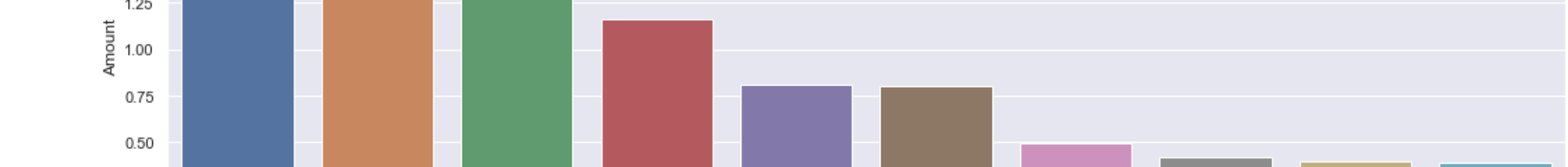
```
Out[30]: <AxesSubplot: xlabel='Marital_Status', ylabel='Amount'>
```



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

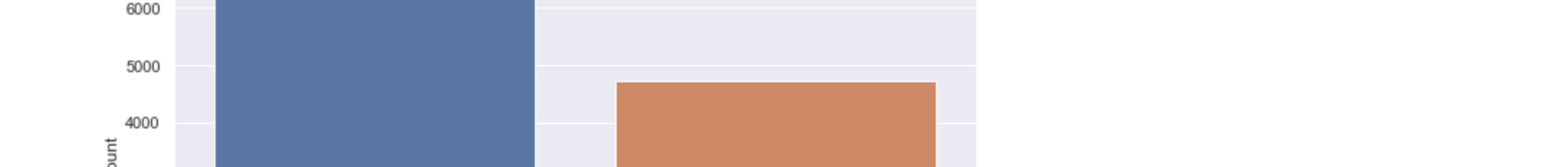
Occupation

```
In [37]: ax = sns.countplot(data = df, x = 'Occupation')
sns.set(rc = {'figure.figsize' : (25,5)})
```



```
In [44]: sales_state = df.groupby(['Occupation'], as_index = False)['Amount'].sum().sort_values(by = 'Amount', ascending = False).head(10)
sns.set(rc = {'figure.figsize' : (30,5)})
sns.barplot(data = sales_state, x = 'Occupation', y = 'Amount')
```

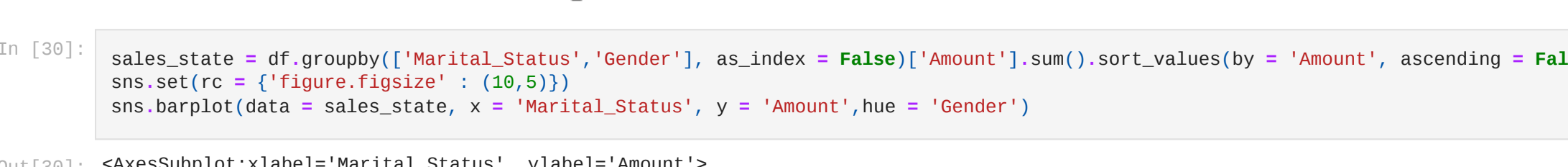
```
Out[44]: <AxesSubplot: xlabel='Occupation', ylabel='Amount'>
```



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

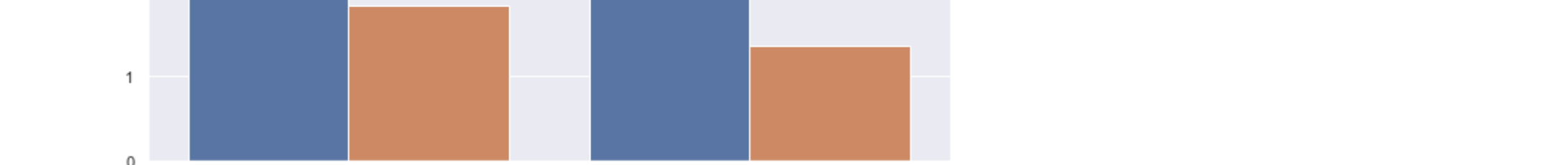
Product Category

```
In [60]: ax = sns.countplot(data = df, x = 'Product_Category')
sns.set(rc = {'figure.figsize' : (35,15)})
```



```
In [67]: sales_state = df.groupby(['Product_Category'], as_index = False)['Amount'].sum().sort_values(by = 'Amount', ascending = False).head(10)
sns.set(rc = {'figure.figsize' : (30,5)})
sns.barplot(data = sales_state, x = 'Product_Category', y = 'Amount')
```

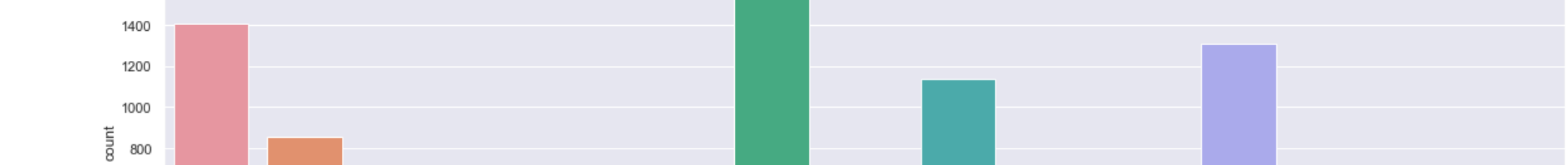
```
Out[67]: <AxesSubplot: xlabel='Product_Category', ylabel='Amount'>
```



Insights: From the above graphs we can see that most of the products are from Food, Clothing and Electrical Category

```
In [68]: sales_state = df.groupby(['Product_ID'], as_index = False)['Orders'].sum().sort_values(by = 'Orders', ascending = False).head(10)
sns.set(rc = {'figure.figsize' : (30,5)})
sns.barplot(data = sales_state, x = 'Product_ID', y = 'Orders')
```

```
Out[68]: <AxesSubplot: xlabel='Product_ID', ylabel='Orders'>
```



CONCLUSION

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

```
In [ ]:
```