

```
In [2]: #import necessary libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

In [4]: #read the dataset
df = pd.read_csv('C:/Users/Nadim/Downloads/Diwali Sales Data.csv', encoding = 'unicode_escape')
df.shape

Out[4]: (11251, 15)
```

```
In [6]: 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 | N      |
| 1     | 1000732 | Kartik      | P00110942  | F      | 26-35     | 35  | 1              | Andhra Pradesh | Southern | Govt            | Auto             | 3      | 23934.0 | N      |
| 2     | 1001990 | Bindu       | P00118542  | F      | 26-35     | 35  | 1              | Uttar Pradesh  | Central  | Automobile      | Auto             | 3      | 23924.0 | N      |
| 3     | 1001425 | Sudevi      | P00237842  | M      | 0-17      | 16  | 0              | Karnataka      | Southern | Construction    | Auto             | 2      | 23912.0 | N      |
| 4     | 1000588 | Joni        | P00057942  | M      | 26-35     | 28  | 1              | Gujarat        | Western  | Food Processing | Auto             | 2      | 23877.0 | N      |
| ...   | ...     | ...         | ...        | ...    | ...       | ... | ...            | ...            | ...      | ...             | ...              | ...    | ...     | ...    |
| 11246 | 1000695 | Manning     | P00296942  | M      | 18-25     | 19  | 1              | Maharashtra    | Western  | Chemical        | Office           | 4      | 370.0   | N      |
| 11247 | 1004089 | Reichenbach | P00171342  | M      | 26-35     | 33  | 0              | Haryana        | Northern | Healthcare      | Veterinary       | 3      | 367.0   | N      |
| 11248 | 1001209 | Oshin       | P00201342  | F      | 36-45     | 40  | 0              | Madhya Pradesh | Central  | Textile         | Office           | 4      | 213.0   | N      |
| 11249 | 1004023 | Noonan      | P00059442  | M      | 36-45     | 37  | 0              | Karnataka      | Southern | Agriculture     | Office           | 3      | 206.0   | N      |
| 11250 | 1002744 | Brunley     | P00281742  | F      | 18-25     | 19  | 0              | Maharashtra    | Western  | Healthcare      | Office           | 3      | 188.0   | N      |

11251 rows × 15 columns

```
In [8]: #check data types
df.dtypes

Out[8]: User_ID          int64
Cust_name         object
Product_ID        object
Gender            object
Age_Group         object
Age              int64
Marital_Status    int64
State            object
Zone            object
Occupation        object
Product_Category  object
Orders           int64
Amount          float64
Status           float64
unnamed1         float64
dtype: object

In [10]: 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(13), int64(4), object(8)
memory usage: 1.3+ MB

In [12]: #see statistical measures
df.describe()

Out[12]:
```

|       | 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 [14]: #check null values
df.isnull().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           12
Status           11251
unnamed1         11251
dtype: int64

In [16]: df.columns

Out[16]: Index(['User_ID', 'Cust_name', 'Product_ID', 'Gender', 'Age_Group', 'Age',
       'Marital_Status', 'State', 'Zone', 'Occupation', 'Product_Category',
       'Orders', 'Amount', 'Status', 'unnamed1'],
      dtype='object')

In [18]: df.drop(['Status', 'unnamed1'], axis=1, inplace=True)

In [20]: df.shape

Out[20]: (11251, 13)

In [22]: df.isnull().sum()

Out[22]: 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 [24]: df.dropna(inplace=True)

In [26]: 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 [28]: #change data type
df['Amount'] = df['Amount'].astype('int')

In [30]: df['Amount'].dtypes

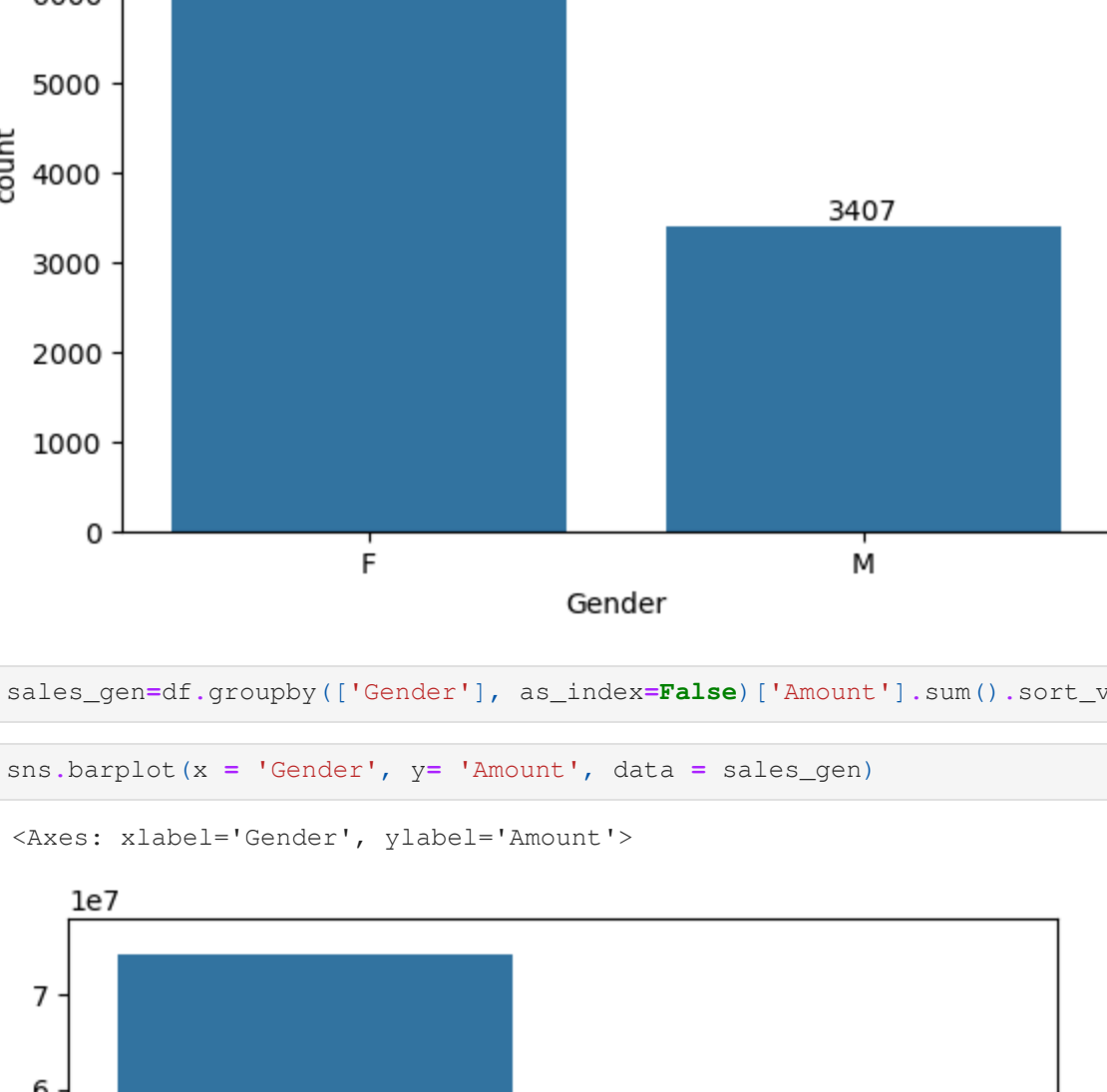
Out[30]: dtype('int32')
```

## Exploratory Data Analysis

### Gender

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

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



```
In [38]: sales_gen=df.groupby(['Gender'], as_index=False) ['Amount'].sum().sort_values(by='Amount', ascending=False)

In [40]: sns.barplot(x = 'Gender', y = 'Amount', data = sales_gen)

Out[40]: <Axes: xlabel='Gender', ylabel='Amount'>
```



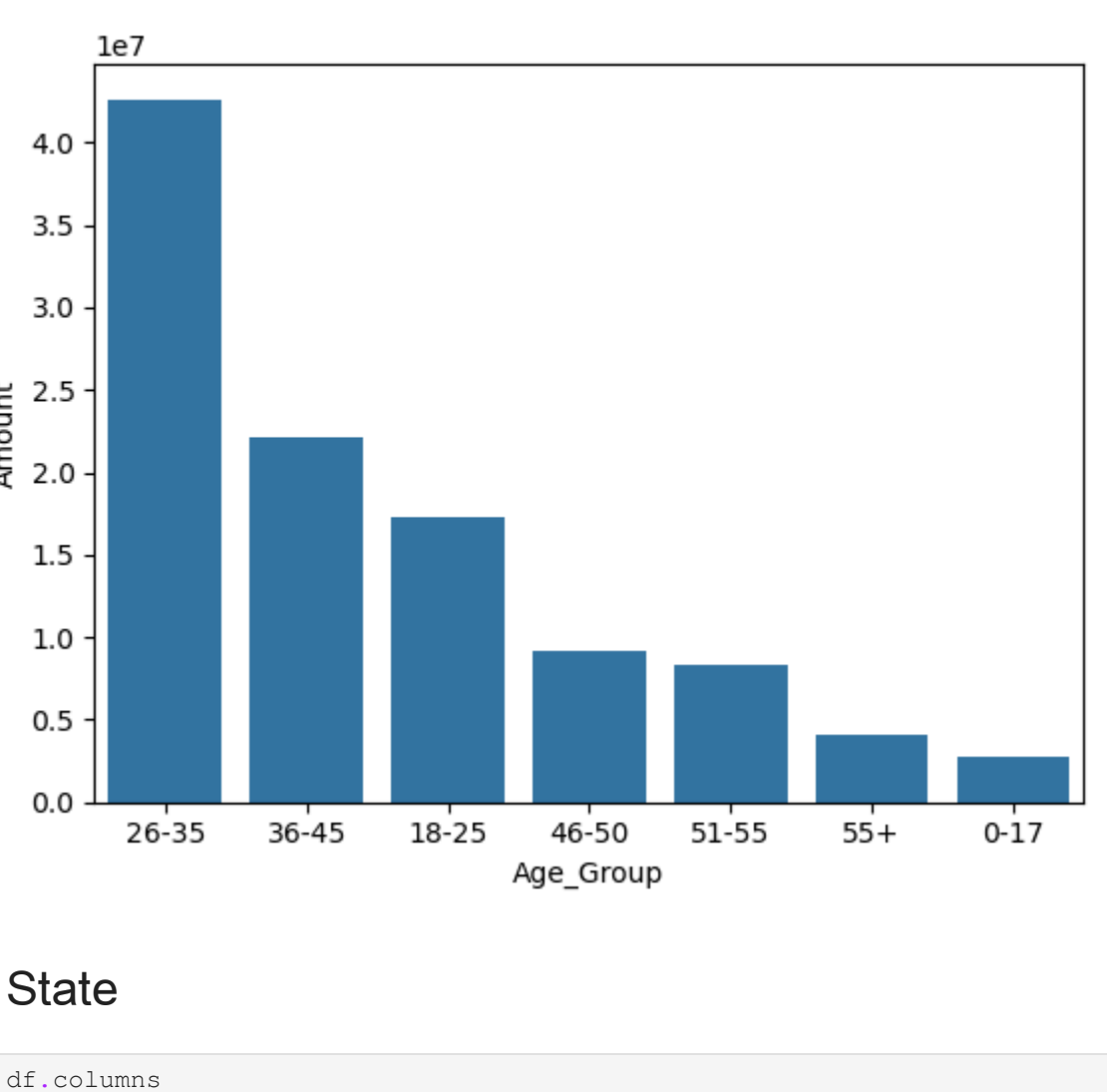
### Age

```
In [45]: df.columns

Out[45]: Index(['User_ID', 'Cust_name', 'Product_ID', 'Gender', 'Age_Group', 'Age',
       'Marital_Status', 'State', 'Zone', 'Occupation', 'Product_Category',
       'Orders', 'Amount'],
      dtype='object')

In [49]: ax = sns.countplot(data = df, x = 'Age_Group', hue = 'Gender')


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



```
In [55]: 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[55]: <Axes: xlabel='Age_Group', ylabel='Amount'>
```



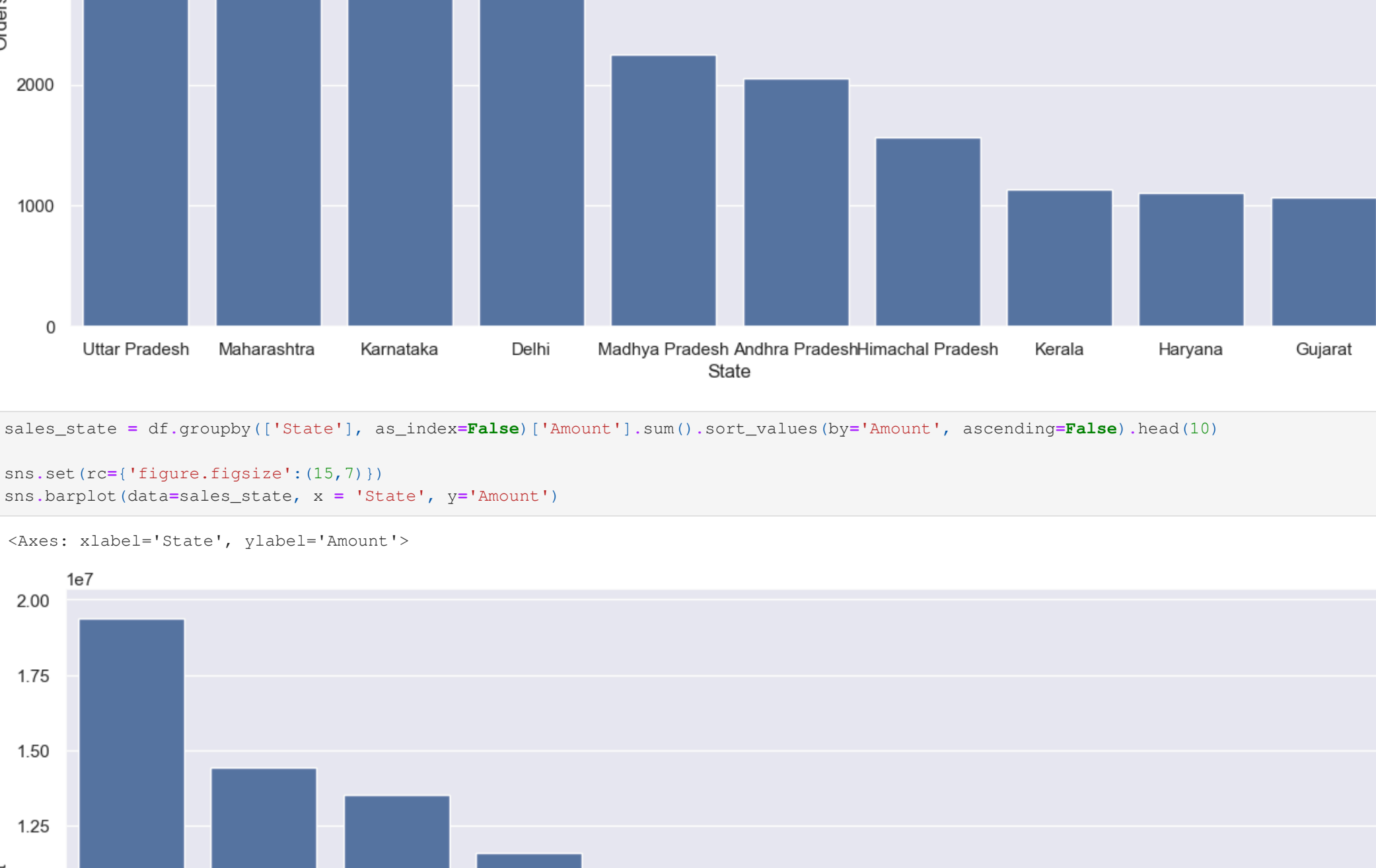
### State

```
In [58]: df.columns

Out[58]: Index(['User_ID', 'Cust_name', 'Product_ID', 'Gender', 'Age_Group', 'Age',
       'Marital_Status', 'State', 'Zone', 'Occupation', 'Product_Category',
       'Orders', 'Amount'],
      dtype='object')


In [60]: sales_state = df.groupby(['State'], as_index=False) ['Orders'].sum().sort_values(by='Orders', ascending=False).head(10)

sns.set(rc={'figure.figsize':(15,7)})
sns.barplot(data=sales_state, x = 'State', y='Orders')
```



```
In [62]: sales_state = df.groupby(['State'], as_index=False) ['Amount'].sum().sort_values(by='Amount', ascending=False).head(10)

sns.set(rc={'figure.figsize':(15,7)})
sns.barplot(data=sales_state, x = 'State', y='Amount')
```



### Marital Status


```
In [65]: ax = sns.countplot(data = df, x = 'Marital_Status')

sns.set(rc={'figure.figsize':(15,5)})
for bars in ax.containers:
    ax.bar_label(bars)
```



```
In [67]: 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')
```



### Occupation


```
In [70]: sns.set(rc={'figure.figsize':(20,5)})
ax = sns.countplot(data = df, x = 'Occupation')

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



```
In [72]: 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')
```



### Product Category

```
In [75]: 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 [77]: sales_state = df.groupby(['Product_Category'], as_index=False) ['Amount'].sum().sort_values(by='Amount', ascending=False).head(10)

sns.barplot(data=sales_state, x = 'Product_Category', y='Amount')
```



### Conclusion:

Age group of 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.

```
In [ ]:
```