```
In [1]:  # import python librray :-

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

```
In [2]:  # import csv file

         df= pd.read_csv("E:\Diwali Sales Data (1).csv",encoding='unicode_escape')
```

```
In [3]:  df
```

Out[3]:

| | 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 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 11246 | 1000695 | Manning | P00296942 | M | 18-25 | 19 | 1 | Maharashtra |
| 11247 | 1004089 | Reichenbach | P00171342 | M | 26-35 | 33 | 0 | Haryana |
| 11248 | 1001209 | Oshin | P00201342 | F | 36-45 | 40 | 0 | Madhya Pradesh |
| 11249 | 1004023 | Noonan | P00059442 | M | 36-45 | 37 | 0 | Karnataka |
| 11250 | 1002744 | Brumley | P00281742 | F | 18-25 | 19 | 0 | Maharashtra |

11251 rows × 15 columns

```
In [4]:  df.shape
```

Out[4]:  (11251, 15)

```
In [5]:  df.head()
```

Out[5]:

| | User_ID | Cust_name | Product_ID | Gender | Age Group | Age | Marital_Status | State | Zon |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1002903 | Sanskriti | P00125942 | F | 26-35 | 28 | 0 | Maharashtra | Wester |
| 1 | 1000732 | Kartik | P00110942 | F | 26-35 | 35 | 1 | Andhra Pradesh | Souther |
| 2 | 1001990 | Bindu | P00118542 | F | 26-35 | 35 | 1 | Uttar Pradesh | Centra |
| 3 | 1001425 | Sudevi | P00237842 | M | 0-17 | 16 | 0 | Karnataka | Souther |
| 4 | 1000588 | Joni | P00057942 | M | 26-35 | 28 | 1 | Gujarat | Wester |

In [6]: `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 [8]:
```python
# Drop unrelated/blank column

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

In [9]:
```python
# Check all null values
pd.isnull(df).sum()
```

Out[9]:
```
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 [10]:
```python
# Drop all null values
```

```
df.dropna(inplace=True)
```

In [11]:
```python
# Change datatype
df['Amount'] = df['Amount'].astype('int')
```

In [12]:
```python
df['Amount'].dtypes
```

Out[12]:
```
dtype('int32')
```

In [13]:
```python
df.columns
```

Out[13]:
```
Index(['User_ID', 'Cust_name', 'Product_ID', 'Gender', 'Age Group', 'Age',
       'Marital_Status', 'State', 'Zone', 'Occupation', 'Product_Category',
       'Orders', 'Amount'],
      dtype='object')
```

In [14]:
```python
# Rename columns
df.rename(columns = {'Marital_Status' : 'Shadi'})
```

Out[14]:

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

11239 rows × 13 columns

In [15]:
```python
# Describe () method returns description of the data in the dataframe( i.e count, r

df.describe()
```

Out[15]:

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

In [16]:
```python
# Use describe () for specific columns

df[['Age','Orders', 'Amount']].describe()
```

Out[16]:

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

# Exploratory Data Analysis

In [17]:
```python
# 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 [18]:
```python
# Plotting a bar chart for gender vs total amount
sales_gen = df.groupby (['Gender'], as_index=False)['Amount'].sum().sort_values(by
sns.barplot(x = 'Gender',y='Amount',data = sales_gen)
```

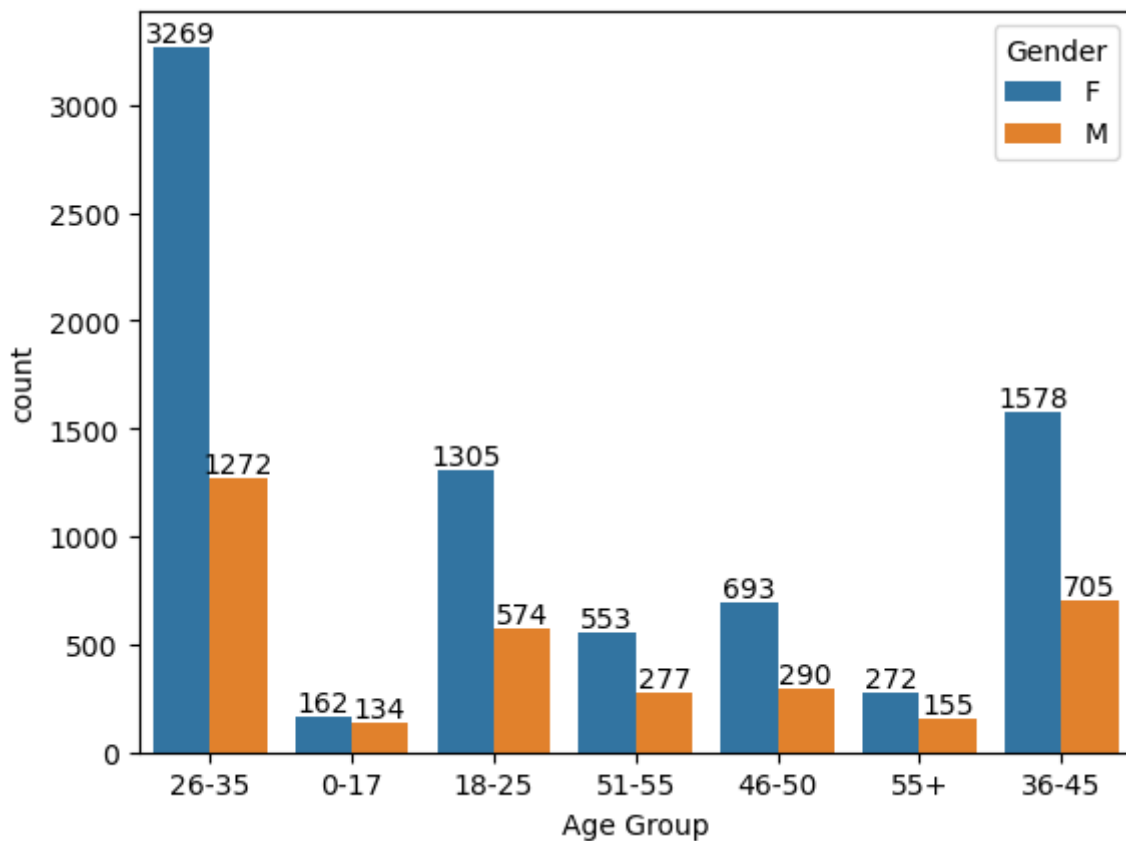Out[18]: <AxesSubplot:xlabel='Gender', ylabel='Amount'>

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

## Age

```
In [20]: ax = sns.countplot(data = df, x='Age Group', hue ='Gender')

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



```
In [21]: # Total amount Vs Total Age Group

         sales_age = df.groupby(['Age Group'],as_index=False)['Amount'].sum().sort_values(by

         sns.barplot(x= 'Age Group', y = 'Amount', data = sales_age)
```

```
Out[21]: <AxesSubplot:xlabel='Age Group', ylabel='Amount'>
```

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

## STATE

In [25]:
```python
# Total number of orders from top 10 states

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

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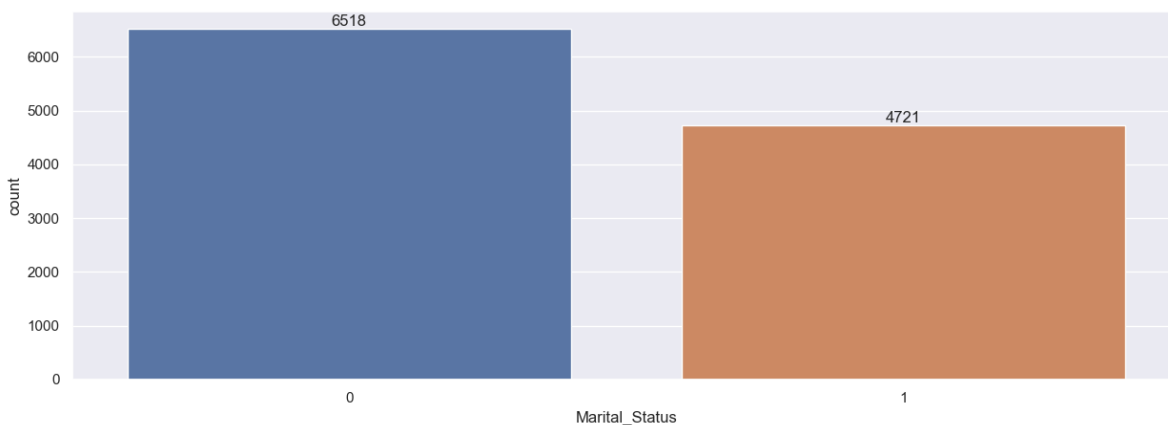```

Out[25]:  <AxesSubplot:xlabel='State', ylabel='Orders'>



In [26]:
```python
# Total amount/sales from top 10 states

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

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

Out[26]:  <AxesSubplot:xlabel='State', ylabel='Amount'>



## From the above graph we can see that most of the orders & total sales / amount are from uttar pradesh , Maharashtra and Karnataka respectively.

## Marital Status

In [28]:
```
ax = sns.countplot(data = df, x = 'Marital_Status')

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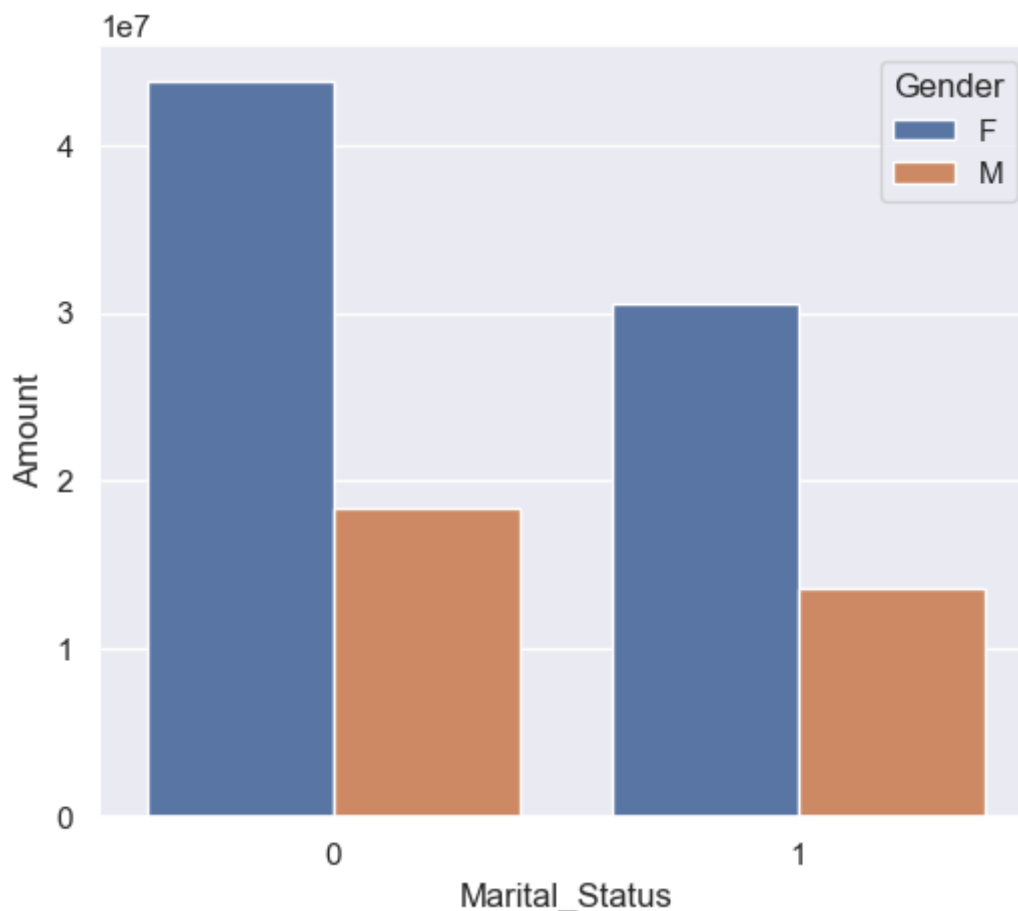


In [30]:
```
sales_state = df.groupby(['Marital_Status', 'Gender'], as_index=False)['Amount'].su

sns.set(rc={'figure.figsize':(6,5)})
sns.barplot(data = sales_state,x ='Marital_Status',y = 'Amount',hue = 'Gender')
```
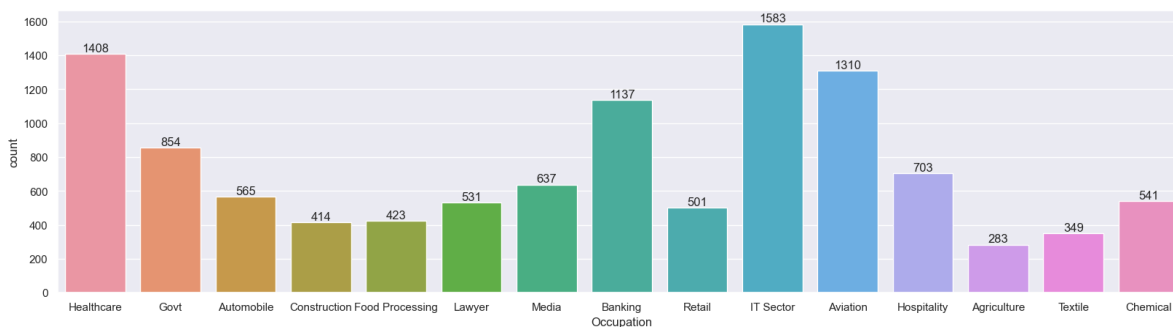
Out[30]:  <AxesSubplot:xlabel='Marital_Status', ylabel='Amount'>

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

## Occupation

In [31]:
```python
sns.set(rc={'figure.figsize':(20,5)})
ax = sns.countplot(data = df, x = 'Occupation')

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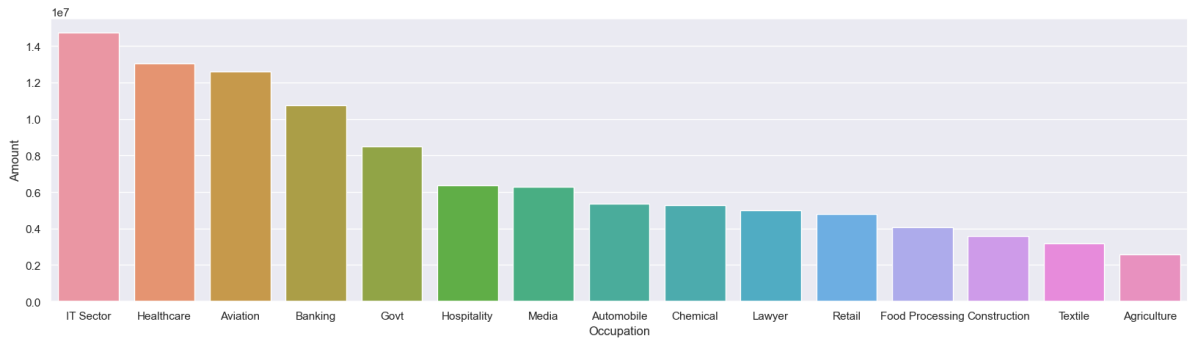


In [32]:
```python
sales_state = df.groupby(['Occupation'],as_index=False)['Amount'].sum().sort_value

sns.set(rc={'figure.figsize':(20,5)})
sns.barplot(data = sales_state,x= 'Occupation',y = 'Amount')
```
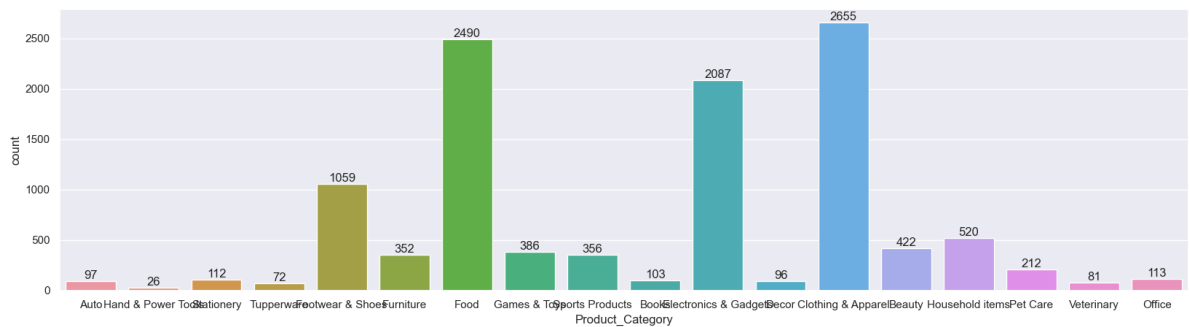
Out[32]:
```
<AxesSubplot:xlabel='Occupation', ylabel='Amount'>
```

## From the above graph we can see thar most of the buyers are working in IT , Healthcare and Aviation sector

# Product Category

In [33]:
```python
sns.set(rc={'figure.figsize':(20,5)})
ax = sns.countplot(data=df, x= 'Product_Category')

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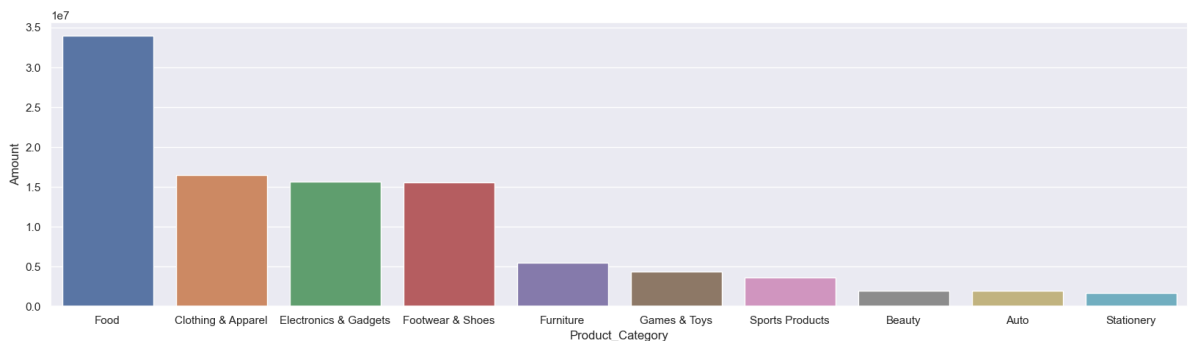


In [34]:
```python
sales_state = df.groupby(['Product_Category'], as_index=False)['Amount'].sum().sor
```
```python
sns.set(rc={'figure.figsize' :(20,5)})
sns.barplot(data = sales_state, x = 'Product_Category', y = 'Amount')
```

Out[34]:
```
<AxesSubplot:xlabel='Product_Category', ylabel='Amount'>
```
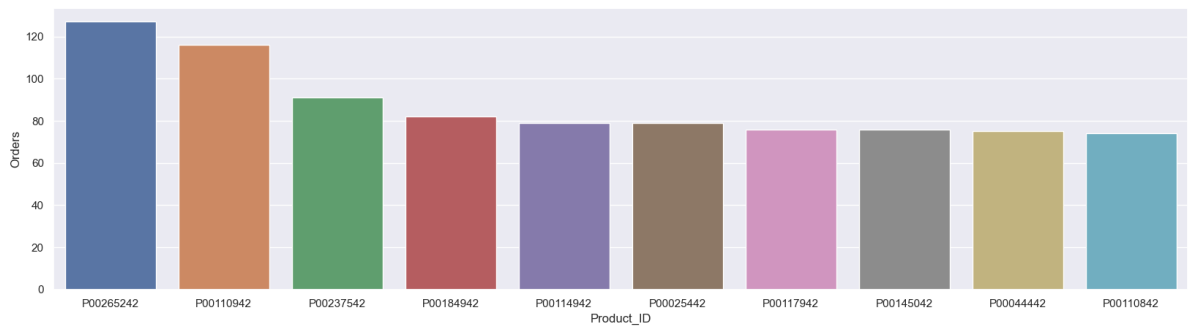


## From the above graph we can see that most of the sold products are from Food, Clothing,and Electonics category

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

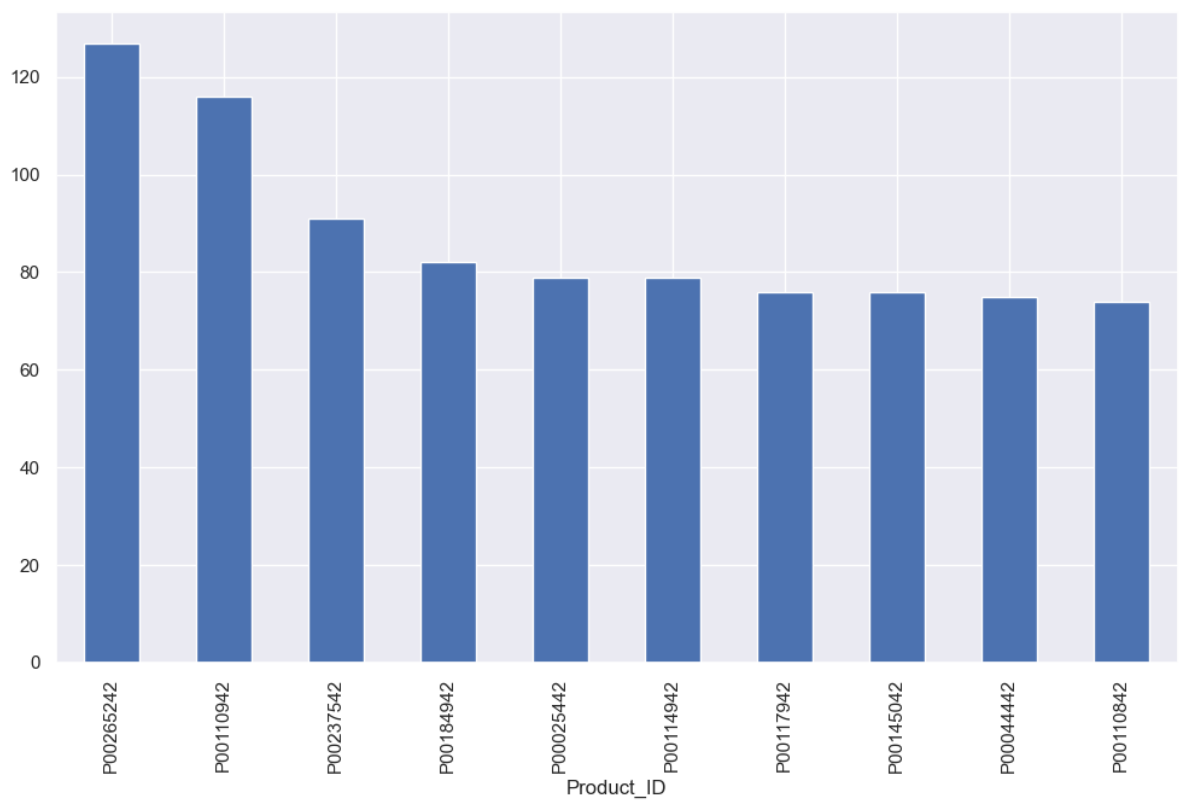Out[37]:     `<AxesSubplot:xlabel='Product_ID', ylabel='Orders'>`



In [38]:
```python
# 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[38]:     `<AxesSubplot:xlabel='Product_ID'>`



## Conclusion :-

Married women age group between 26-35 yrs from Up, Maharashtra and Karnataka working in IT, Healrhcare and Aviation are more likely to but products from Food,Clothing and Electronics category