

Black friday Dataset EDA and Feature Engineering

Cleaning and preparing the data for model training

Problem Statement

A retail company "ABC Private Limited" wants to understand the customer purchase behaviour (specifically, purchase amount) against various products of different categories. They have shared purchase summary of various customers for selected high volume products from last month. The data set also contains customer demographics (age, gender, marital status, city_type, stay_in_current_city), product details (product_id and product category) and Total purchase_amount from last month.

Now, they want to build a model to predict the purchase amount of customer against various products which will help them to create personalized offer for customers against different products.

Handling Categorical features

In [13]:

Handling Categorical feature Gender

In [14]:

df['Gender']=df['Gender'].replace({'M','F'},[1,0])

In [15]:

df.head()

Out [15]:

	Product_ID	Gender	Age	Occupation	City_Category	Stay_In_Current_City_Years	Marital_Status	Product_Category_1	Product_Category_2	Product_Category_3	Purchase
0	P00069042	0	0-17	10	A	2	0	3	NaN	NaN	8370.0
1	P00248942	0	0-17	10	A	2	0	1	6.0	14.0	15200.0
2	P00087842	0	0-17	10	A	2	0	12	NaN	NaN	1422.0
3	P00085442	0	0-17	10	A	2	0	12	14.0	NaN	1057.0
4	P00285442	1	55+	16	C	4+	0	8	NaN	NaN	7969.0

In [16]:

Handling Categorical feature Age

In [17]:

df.Age.unique()

Out [17]:

array(['0-17', '55+', '26-35', '46-50', '51-55', '36-45', '18-25'], dtype=object)

In [18]:

when training the model it will be able to understand
This is known as Target ordinal encoding or Target Guiding

In [19]:

ppd.get_dummies(df['Age'],drop_first=True)
df['Age']=df['Age'].map({'0-17':1,'18-25':2,'26-35':3,'36-45':4,'46-50':5,'51-55':6,'55+':7})

In [20]:

df.head()

Out [20]:

	Product_ID	Gender	Age	Occupation	City_Category	Stay_In_Current_City_Years	Marital_Status	Product_Category_1	Product_Category_2	Product_Category_3	Purchase
0	P00069042	0	1	10	A	2	0	3	NaN	NaN	8370.0
1	P00248942	0	1	10	A	2	0	1	6.0	14.0	15200.0
2	P00087842	0	1	10	A	2	0	12	NaN	NaN	1422.0
3	P00085442	0	1	10	A	2	0	12	14.0	NaN	1057.0
4	P00285442	1	7	16	C	4+	0	8	NaN	NaN	7969.0

In [21]:

Handling Categorical feature City_Category

In [22]:

drop_first=True
because if we have 3 categories then 2 categories are sufficient to represent all the 3 categories

In [23]:

df_city=pd.get_dummies(df['City_Category'],drop_first=True)
df_city

Out [23]:

	B	C
0	0	0
1	0	0
2	0	0
3	0	0
4	0	1

In [24]:

df=pd.concat([df,df_city],axis=1)
df.head()

Out [24]:

	Product_ID	Gender	Age	Occupation	City_Category	Stay_In_Current_City_Years	Marital_Status	Product_Category_1	Product_Category_2	Product_Category_3	Purchase	B	C
0	P00069042	0	1	10	A	2	0	3	NaN	NaN	8370.0	0	0
1	P00248942	0	1	10	A	2	0	1	6.0	14.0	15200.0	0	0
2	P00087842	0	1	10	A	2	0	12	NaN	NaN	1422.0	0	0
3	P00085442	0	1	10	A	2	0	12	14.0	NaN	1057.0	0	0
4	P00285442	1	7	16	C	4+	0	8	NaN	NaN	7969.0	0	1

Finding the Missing values and handling them

23980 1 0

783667 rows x 2 columns

In [24]:

df=pd.concat([df,df_city],axis=1)
df.head()

Out [24]:

	Product_ID	Gender	Age	Occupation	Stay_In_Current_City_Years	Marital_Status	Product_Category_1	Product_Category_2	Product_Category_3	Purchase	B	C	
0	P00069042	0	1	10	A	2	0	3	NaN	NaN	8370.0	0	0
1	P00248942	0	1	10	A	2	0	1	6.0	14.0	15200.0	0	0
2	P00087842	0	1	10	A	2	0	12	NaN	NaN	1422.0	0	0
3	P00085442	0	1	10	A	2	0	12	NaN	NaN	1057.0	0	0
4	P00285442	1	7	16	C	4+	0	8	NaN	NaN	7969.0	0	1

In [25]:

City_Category is not required anymore can be dropped

In [26]:

inplace= True is used to remove City_Category permanently

In [27]:

df.drop('City_Category',axis=1,inplace=True)

In [28]:

df.head()

Out [28]:

	Product_ID	Gender	Age	Occupation	Stay_In_Current_City_Years	Marital_Status	Product_Category_1	Product_Category_2	Product_Category_3	Purchase	B	C	
0	P00069042	0	1	10		2	0	3	NaN	NaN	8370.0	0	0
1	P00248942	0	1	10		2	0	1	6.0	14.0	15200.0	0	0
2	P00087842	0	1	10		2	0	12	NaN	NaN	1422.0	0	0
3	P00085442	0	1	10		2	0	12	14.0	NaN	1057.0	0	0
4	P00285442	1	7	16		4+	0	8	NaN	NaN	7969.0	0	1

Finding the Missing values and handling them

In [29]:

df.isnull().sum()

Out [29]:

Product_ID	0
Gender	0
Age	0
Occupation	0
Stay_In_Current_City_Years	0
Marital_Status	0
Product_Category_1	0
Product_Category_2	245982
Product_Category_3	545889
Purchase	235999
B	0
C	0
dtype: int64	

In [30]:

Purchase has the null data because of the test data

In [31]:

Focusing on replacing missing values
Product_Category_2 245982
Product_Category_3 545889

In [32]:

For 'Product_Category_2' missing values

In [33]:

The best way for replacing the missing values categorical or discrete variables is with the mode
.fillna() function is used to fill nan values

In [34]:

df['Product_Category_2']=df['Product_Category_2'].fillna(df['Product_Category_2'].mode()[0])

In [35]:

df['Product_Category_2'].isnull().sum()

Out [35]:

0

In [36]:

For 'Product_Category_3' missing values

In [37]:

df['Product_Category_3'].unique()

Out [37]:

array([nan, 14., 17., 5., 4., 16., 15., 8., 9., 13., 6., 12., 3.,
18., 11., 10.])

In [38]:

df['Product_Category_3'].mode()[0]

Out [38]:

16.0

In [39]:

df['Product_Category_3']=df['Product_Category_3'].fillna(df['Product_Category_3'].mode()[0])
df['Product_Category_3'].isnull().sum()

Out [39]:

0

In [40]:

df.isnull().sum()

Out [40]:

Product_ID	0
------------	---

Visualization

In [49]:	<pre>sns.barplot('Age','Purchase',data=df);</pre>
Out [49]:	<pre>F:\Anaconda\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variables as keyword args: x, y. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation. warnings.warn(</pre>
In [78]:	<pre>sns.barplot('Age','Purchase',hue='Gender',data=df,ci=None);</pre>
Out [78]:	<pre>F:\Anaconda\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variables as keyword args: x, y. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation. warnings.warn(</pre>
In [51]:	<pre># Purchasing of Men is higher than the Women</pre>
In [53]:	<pre># Visualization of Purchase with respective Occupations</pre>
In [79]:	<pre>sns.barplot('Occupation','Purchase',hue='Gender',data=df,ci=None);</pre>
Out [79]:	<pre>F:\Anaconda\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variables as keyword args: x, y. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation. warnings.warn(</pre>
In [53]:	<pre># Visualization of Purchase with respective 'Product_Category_1'</pre>
In [80]:	<pre>sns.barplot('Product_Category_1','Purchase',hue='Gender',data=df,ci=None);</pre>
Out [80]:	<pre>F:\Anaconda\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variables as keyword args: x, y. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation. warnings.warn(</pre>
In [55]:	<pre># Visualization of Purchase with respective 'Product_Category_2'</pre>
In [81]:	<pre>sns.barplot('Product_Category_2','Purchase',hue='Gender',data=df,ci=None);</pre>
Out [81]:	<pre>F:\Anaconda\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variables as keyword args: x, y. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation. warnings.warn(</pre>
In [57]:	<pre># Visualization of Purchase with respective 'Product_Category_3'</pre>
In [82]:	<pre>sns.barplot('Product_Category_3','Purchase',hue='Gender',data=df,ci=None);</pre>
Out [82]:	<pre>F:\Anaconda\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variables as keyword args: x, y. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation. warnings.warn(</pre>

Feature Scaling

In [69]:
