

Dataset Information

This dataset comprises of sales transactions captured at a retail store. It's a classic dataset to explore and expand your feature engineering skills and day to day understanding from multiple shopping experiences. This is a regression problem.

Attributes:

	Column ID	Column Name	Data type	Description	Masked
	0	User_ID	int64	Unique Id of customer	False
	1	Product_ID	object	Unique Id of product	False
	2	Gender	object	Sex of customer	False
	3	Age	object	Age of customer	False
	4	Occupation	int64	Occupation code of customer	True
	5	City_Category	object	City of customer	True
	6	Stay_In_Current_City_Years	object	Number of years of stay in city	False
	7	Marital_Status	int64	Marital status of customer	False
	8	Product_Category_1	int64	Category of product	True
	9	Product_Category_2	float64	Category of product	True
	10	Product_Category_3	float64	Category of product	True
	11	Purchase	int64	Purchase amount	False

Import modules

```
In [40]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import warnings
warnings.filterwarnings('ignore')
```

Loading the dataset

```
In [41]: df = pd.read_csv('train.csv')
df.head()
```

```
Out[41]:
```

User_ID	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	1000001	P0006962	F	0-17	10	A	2	0	3	NaN	NaN	8370
1	1000001	P00249842	F	0-17	10	A	2	0	1	6.0	14.0	15200
2	1000001	P0006962	F	0-17	10	A	2	0	12	NaN	NaN	1422
3	1000001	P0006962	F	0-17	10	A	2	0	12	14.0	NaN	1067
4	1000002	P00285442	M	55+	16	C	4+	0	8	NaN	NaN	7969

```
In [42]: # Statistical Info
df.describe()
```

```
Out[42]:
```

	User_ID	Info	Marital_Status	Product_Category_1	Product_Category_3	Product_Category_3	Purchase
count	5.000000e+05	550068.000000	550068.000000	550068.000000	376430.000000	196821.000000	550068.000000
mean	1.000000e+06	8.076707	0.496953	5.404270	9.842329	12.866243	8263.966713
std	1.127592e+03	1.620260	0.491770	3.936211	5.086990	4.125338	5023.065394
min	1.000001e+06	0.000000	0.000000	1.000000	2.000000	3.000000	12.000000
25%	1.003016e+06	2.000000	0.000000	1.000000	5.000000	9.000000	5823.000000
50%	1.003077e+06	6.000000	0.000000	5.000000	9.000000	14.000000	8047.000000
75%	1.004478e+06	14.000000	1.000000	8.000000	15.000000	16.000000	12054.000000
max	1.006040e+06	20.000000	1.000000	20.000000	18.000000	18.000000	23961.000000

```
In [43]: # Data type Info
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 550068 entries, 0 to 550067
Data columns (total 12 columns):
#   Column                Non-Null Count  Dtype
---  --
0   User_ID                558868 non-null    int64
1   Product_ID             558868 non-null    object
2   Gender                 558868 non-null    object
3   Age                    558868 non-null    object
4   Occupation              558868 non-null    int64
5   City_Category          558868 non-null    object
6   Stay_In_Current_City_Years  558868 non-null    object
7   Marital_Status         558868 non-null    int64
8   Product_Category_1     558868 non-null    int64
9   Product_Category_2     376430 non-null    float64
10  Product_Category_3     196821 non-null    float64
11  Purchase                558868 non-null    int64
dtypes: float64(2), int64(5), object(5)
memory usage: 80.4+ MB
```

```
In [44]: # Find unique values
df.apply(lambda x: len(x.unique()))
```

```
Out[44]:
```

User_ID	Product_ID	Gender	Age	Occupation	City_Category	Stay_In_Current_City_Years	Marital_Status	Product_Category_1	Product_Category_2	Product_Category_3	Purchase
5891	3831	2	7	21	5	2	20	18	18189		

Exploratory Data Analysis

```
In [45]: # distplot for purchase
plt.style.use('fivethirtyeight')
plt.figure(figsize=(15, 7))
sns.distplot(df['Purchase'], bins=25)
```

```
Out[45]:
```

<AxesSubplot: xlabel='Purchase', ylabel='Density'>

```
In [46]: # distribution of numeric variables
sns.countplot(df['Gender'])
```

```
Out[46]:
```

<AxesSubplot: xlabel='Gender', ylabel='count'>

```
In [47]: sns.countplot(df['Age'])
```

<AxesSubplot: xlabel='Age', ylabel='count'>

```
In [48]: sns.countplot(df['Marital_Status'])
```

<AxesSubplot: xlabel='Marital_Status', ylabel='count'>

```
In [49]: sns.countplot(df['Occupation'])
```

<AxesSubplot: xlabel='Occupation', ylabel='count'>

```
In [50]: sns.countplot(df['Product_Category_1'])
```

<AxesSubplot: xlabel='Product_Category_1', ylabel='count'>

```
In [51]: sns.countplot(df['Product_Category_2'])
```

<AxesSubplot: xlabel='Product_Category_2', ylabel='count'>

```
In [52]: sns.countplot(df['Product_Category_3'])
```

<AxesSubplot: xlabel='Product_Category_3', ylabel='count'>

```
In [53]: sns.countplot(df['City_Category'])
```

<AxesSubplot: xlabel='City_Category', ylabel='count'>

```
In [54]: sns.countplot(df['Stay_In_Current_City_Years'])
```

<AxesSubplot: xlabel='Stay_In_Current_City_Years', ylabel='count'>

```
In [55]: # bivariate analysis
occupation_plot = df.pivot_table(index='Occupation', values='Purchase', aggfunc=np.mean)
age_plot = df.pivot_table(index='Age', values='Purchase', aggfunc=np.mean)
plt.figure(figsize=(15, 7))
plt.xlabel('Age')
plt.ylabel('Purchase')
plt.title('Occupation and Purchase Analysis')
plt.xticks(rotation=0)
plt.show()
```

Occupation and Purchase Analysis

```
In [56]: age_plot = df.pivot_table(index='Age', values='Purchase', aggfunc=np.mean)
age_plot.plot(kind='bar', figsize=(15, 7))
plt.xlabel('Age')
plt.ylabel('Purchase')
plt.title('Age and Purchase Analysis')
plt.xticks(rotation=0)
plt.show()
```

Age and Purchase Analysis

```
In [57]: gender_plot = df.pivot_table(index='Gender', values='Purchase', aggfunc=np.mean)
gender_plot.plot(kind='bar', figsize=(15, 7))
plt.xlabel('Gender')
plt.ylabel('Purchase')
plt.title('Gender and Purchase Analysis')
plt.xticks(rotation=0)
plt.show()
```

Gender and Purchase Analysis

Preprocessing the dataset

```
In [58]: # check for null values
df.isnull().sum()
```

```
Out[58]:
```

User_ID	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	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	0
5	0	0	0	0	0	0	0	0	0	0	0
6	0	0	0	0	0	0	0	0	0	0	0
7	0	0	0	0	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0	0	0	0
9	0	0	0	0	0	0	0	0	0	0	0
10	0	0	0	0	0	0	0	0	0	0	0
11	0	0	0	0	0	0	0	0	0	0	0
12	0	0	0	0	0	0	0	0	0	0	0
13	0	0	0	0	0	0	0	0	0	0	0
14	0	0	0	0	0	0	0	0	0	0	0
15	0	0	0	0	0	0	0	0	0	0	0
16	0	0	0	0	0	0	0	0	0	0	0
17	0	0	0	0	0	0	0	0	0	0	0
18	0	0	0	0	0	0	0	0	0	0	0
19	0	0	0	0	0	0	0	0	0	0	0
20	0	0	0	0	0	0	0	0	0	0	0

```
In [59]: df['Product_Category_2'] = df['Product_Category_2'].fillna(-2.0).astype('float32')
df['Product_Category_3'] = df['Product_Category_3'].fillna(-2.0).astype('float32')
```

```
df.isnull().sum()
```

```
Out[60]:
```

User_ID	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	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	0
5	0	0	0	0	0	0	0	0	0	0	0
6	0	0	0	0	0	0	0	0	0	0	0
7	0	0	0	0	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0	0	0	0
9	0	0	0	0	0	0	0	0	0	0	0
10	0	0	0	0	0	0	0	0	0	0	0
11	0	0	0	0	0	0	0	0	0	0	0
12	0	0	0	0	0	0	0	0	0	0	0
13	0	0	0	0	0	0	0	0	0	0	0
14	0	0	0	0	0	0	0	0	0	0	0
15	0	0	0	0	0	0	0	0	0	0	0
16	0	0	0	0	0	0	0	0	0	0	0
17	0	0	0	0	0	0	0	0	0	0	0
18	0	0	0	0	0	0	0	0	0	0	0
19	0	0	0	0	0	0	0	0	0	0	0
20	0	0	0	0	0	0	0	0	0	0	0

```
In [61]: # encoding values using dict
gender_dict = {'F':0, 'M':1}
df['Gender'] = df['Gender'].apply(lambda x: gender_dict[x])
df.head()
```

```
Out[61]:
```

User_ID	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	1000001	P0006962	0	0-17	10	A	2	0	3	-2.0	-2.0	8370
1	1000001	P00249842	0	0-17	10	A	2	0	1	6.0	14.0	15200
2	1000001	P0006962	0	0-17	10	A	2	0	12	-2.0	-2.0	1422
3	1000001	P0006962	0	0-17	10	A	2	0	12	14.0	-2.0	1067
4	1000002	P00285442	1	55+	16	C	4+	0	8	-2.0	-2.0	7969

```
In [62]: # to improve the metric use one hot encoding
# Label encoding
cols = ['Age', 'City_Category', 'Stay_In_Current_City_Years']
from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
df[cols] = le.fit_transform(df[cols])
df.head()
```

```
Out[62]:
```

User_ID	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	1000001	P0006962	0	0	10	0	2	0	3	-2.0	-2.0	8370
1	1000001	P00249842	0	0	10	0	2	0	1	6.0	14.0	15200
2	1000001	P0006962	0	0	10	0	2	0	12	-2.0	-2.0	1422
3	1000001	P0006962	0	0	10	0	2	0	12	14.0	-2.0	1067
4	1000002	P00285442	1	16	16	2	4	0	8	-2.0	-2.0	7969

Correlation Matrix

```
In [63]: corr = df.corr()
plt.figure(figsize=(14, 7))
sns.heatmap(corr, annot=True, cmap='coolwarm')
```

```
Out[63]:
```

<AxesSubplot: >

Input Split

```
In [64]: df.head()
```

```
Out[64]:
```

User_ID	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	1000001	P0006962	0	0	10	0	2	0	3	-2.0	-2.0	8370
1	1000001	P00249842	0	0	10	0	2	0	1	6.0	14.0	15200
2	1000001	P0006962	0	0	10	0	2	0	12	-2.0	-2.0	1422
3	1000001	P0006962	0	0	10	0	2	0	12	14.0	-2.0	1067
4	1000002	P00285442	1	16	16	2	4	0	8	-2.0	-2.0	