Experiment 1

Aim -

Data preparation using NumPy and Pandas.

Todo -

- 1. Load data in Pandas.
- 2. Description of the dataset.
- 3. Drop columns that aren't useful.
- 4. Drop rows with maximum missing values.
- 5. Take care of missing data.
- 6. Create dummy variables.
- 7. Find out outliers (manually).
- 8. Standardization and normalization of columns.

Dataset -

The dataset contains data about customers' purchases during the Black Friday sale. This dataset was taken from Kaggle. The dataset has 550k rows and 12 columns. The various columns of the dataset are age, marital status, gender, total purchase amount, and many other features.

Theory -

pandas is a software library written for the Python programming language for data manipulation and analysis. In particular, it offers data structures and operations for manipulating numerical tables and time series. It is a fast, powerful, flexible, and easy-to-use open-source data analysis and manipulation tool.

NumPy is a library for the Python programming language, adds support for large, multi-dimensional arrays and matrices, along with a large collection of high-level mathematical functions to operate on these arrays. It offers comprehensive mathematical functions, random number generators, linear algebra routines, Fourier transforms, and more.

Results -

1. Load dataset.

```
In [1]: import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
```

read csv():

Read a comma-separated values (CSV) file into DataFrame.

head():

This function returns the first n rows for the object based on position. It is useful for quickly testing if your object has the right type of data in it.

		# Read the data into datafrome df = pd.read_csv("train.csv")														
	<pre># Display the first `n` rows. df.head()</pre>															
Out[2]:		User_ID	Product_ID	Gender	Age	Occupation	City_Category	Stay_In_Current_City_Years	Marital_Status	Product_Category_1	Product_Category_2	Product				
	0	1000001	P00069042	F	0- 17	10	А	2	0	3	NaN					
	1	1000001	P00248942	F	0- 17	10	А	2	0	1	6.0					
	2	1000001	P00087842	F	0- 17	10	А	2	0	12	NaN					
	3	1000001	P00085442	F	0- 17	10	А	2	0	12	14.0					
	4	1000002	P00285442	М	55+	16	С	4+	0	8	NaN					
	4)				

info():

This method prints information about a DataFrame including the index dtype and columns, non-null values and memory usage.

```
In [3]: # Print a concise summary of a DataFrame.
          df.info()
          <class 'pandas.core.frame.DataFrame'>
          RangeIndex: 550068 entries, 0 to 550067
          Data columns (total 12 columns):
          # Column
                                                Non-Null Count
                                                                     Dtype
          --- -----
                                                 -----
                                                 550068 non-null int64
           0 User ID
           1 Product ID
                                                550068 non-null object
          2 Gender 550068 non-null object
3 Age 550068 non-null object
4 Occupation 550068 non-null int64
5 City_Category 550068 non-null object
           6 Stay_In_Current_City_Years 550068 non-null object
           7 Marital_Status 550068 non-null int64
          8 Product_Category_1 550068 non-null int64
9 Product_Category_2 376430 non-null float64
10 Product_Category_3 166821 non-null float64
           11 Purchase
                                                 550068 non-null int64
          dtypes: float64(2), int64(5), object(5)
          memory usage: 50.4+ MB
```

2. Describe dataset.

describe():

Descriptive statistics include those that summarize the central tendency, dispersion and shape of a dataset's distribution, excluding NaN values.

	# Generate descriptive statistics. df.describe()														
	User_ID	Occupation	Marital_Status	Product_Category_1	Product_Category_2	Product_Category_3	Purchas								
count	5.500680e+05	550068.000000	550068.000000	550068.000000	376430.000000	166821.000000	550068.00000								
mean	1.003029e+06	8.076707	0.409653	5.404270	9.842329	12.668243	9263.9687								
std	1.727592e+03	6.522660	0.491770	3.936211	5.086590	4.125338	5023.06539								
min	1.000001e+06	0.000000	0.000000	1.000000	2.000000	3.000000	12.00000								
25%	1.001516e+06	2.000000	0.000000	1.000000	5.000000	9.000000	5823.0000								
50%	1.003077e+06	7.000000	0.000000	5.000000	9.000000	14.000000	8047.00000								
75%	1.004478e+06	14.000000	1.000000	8.000000	15.000000	16.000000	12054.00000								
max	1.006040e+06	20.000000	1.000000	20.000000	18.000000	18.000000	23961.00000								

3. Drop columns that aren't useful.

drop():

Remove rows or columns by specifying label names and corresponding axis, or by specifying direct index or column names. When using a multi-index, labels on different levels can be removed by specifying the level.



isnull():

Detect missing values.

Return a boolean same-sized object indicating if the values are NA. NA values, such as None or numpy. NaN gets mapped to True values. Everything else gets mapped to False values. Characters such as empty strings " or numpy.inf are not considered NA values (unless you set pandas.options.mode.use inf as na = True).

sum():

Returns a data frame or Series of the same size containing the cumulative sum.

```
In [6]: # Count of the number of missing values in each column.
        df.isnull().sum()
Out[6]: User_ID
                                       0
        Product_ID
                                       0
        Gender
                                       0
        Age
                                       0
        Occupation
                                       0
        City_Category
                                       0
        Stay_In_Current_City_Years
        Marital Status
        Product_Category_1
                                       0
        Purchase
                                       0
        dtype: int64
```

4. Create dummy variables.

value_counts():

Return a Series containing counts of unique rows in the DataFrame.

```
In [7]: # Printing count of unique values in each columns.
        cols = list(df.columns)
        delete_cols = ["User_ID", "Product_ID", "Product_Category_1", "Purchase"]
        for i in delete cols:
           cols.remove(i)
        for col in cols:
           print(df[col].value_counts())
        М
             414259
             135809
        Name: Gender, dtype: int64
        26-35
                219587
        36-45
                 110013
        18-25
                 99660
                  45701
        46-50
        51-55
                  38501
        55+
                  21504
        0-17
                 15102
        Name: Age, dtype: int64
              72308
        4
        0
              69638
              59133
        7
        1
              47426
        17
              40043
        20
              33562
        12
              31179
              27309
        2
              26588
        16
              25371
              20355
        6
        3
              17650
        10
              12930
              12177
        15
             12165
        11
              11586
               8461
        19
               7728
        13
               6622
```

concat():

Concatenate pandas objects along a particular axis.

Allows optional set logic along the other axes.

Can also add a layer of hierarchical indexing on the concatenation axis, which may be useful if the labels are the same (or overlapping) on the passed axis number.

get_dummies():

Convert categorical variable into dummy/indicator variables.

In [8]:			dummies.	us! !c	ondon'	leity ca	+ogony'1									
	dumn	i = [Ma nies = col in		us , G	enaer	, City_Ca	tegory]									
			s.append(pd)									
			pd.concat(d ncat((df,du													
	df.h	nead()														
Out[8]:		User_ID	Product_ID	Gender	Age (Occupation	City_Category	Stay_In_Current_City_Year	Marital_Statu	s Product_Ca	egory_1	Purchase	0	1 F	М	А В
	0	1000001	P00069042	F	0- 17	10	А		2	0	3	8370	1	0 1	0	1 0
	1	1000001	P00248942	F	0- 17	10	А		2	0	1	15200	1	0 1	0	1 0
	2	1000001	P00087842	F	0- 17	10	А		2	0	12	1422	1	0 1	0	1 0
	3	1000001	P00085442	F	0- 17	10	А		2	0	12	1057	1	0 1	0	1 0
	4 ′	1000002	P00285442	М	55+	16	С	4		0	8	7969	1	0 0	1	0 0
	4															r
In [9]			ng the colum	nns												
	df	.columr	'Marit	:al_Sta	tus',	'Product_	Category_1',	ge', 'Occupation', 'C 'Purchase', 'Married	ity_Category ', 'Not_Marr	', 'Stay_In_ ied', 'Femal	Current_ e', 'Mal	City_Yea e',	rs'	,		
	df	.head()		_Catego	ry_A',	'City_Ca	itegory_B', '	City_Category_C']								
Out[9]				Ganda	r Age	Occupation	City Category	Stay_In_Current_City_Yea	e Marital Stat	us Product Ca	tegory 1	Durchase	Ma	rried	Not	Marr
	0	100000		F	0	10			2	0	3	8370		1		
	1	100000	1 P00248942	F	0	10) Д		2	0	1	15200		1		
	2	100000	1 P00087842	F	0- 17	10	Α		2	0	12	1422		1		
	3	100000	1 P00085442	F	0- 17	10) Д		2					1		
	4	100000	D0000E440						2	0	12	1057				
			2 P00285442	N	1 55+	16	C	;	+	0	12	1057 7969		1		
	[2 P00285442	N	1 55+	16	S C									>
In [10]	: # df		ng redundant 'Gender','M	: colum	ns.			xis=1,inplace=True)								•
In [10]	: # df df	.drop([.head()	ng redundant 'Gender','M	: colum Marital	ns. _Statu	s','City_	Category'],a	xis=1,inplace=True)	+	0	8	7969	ty (1	ory A	A Ci
	: # df df	.drop([.head()	ng redundant 'Gender','M	: colum Marital	ns. _Statu	s','City_	Category'],a		+	0	8	7969	ty_C	1		A Ci
	: # df df	.drop([.head()	ng redundant 'Gender','M	colum Marital	ns. _Statu Occupati	s','City_ on Stay_In	Category'],a	xis=1,inplace=True) ears Product_Category_1	+ Purchase Mari	0 ied Not_Marrie	g d Female	7969 Male Ci	ty_C	1	1	
	: # df df :	.drop([.head() .head() .drop([.head() .head() .drop([.head() .head() .head()	og redundant 'Gender','M D Product_ID 1 P00069042	Age C	ns. _Statu Occupati	s','City_ on Stay_In	Category'],a	xis=1,inplace=True) lears Product_Category_1 2 3	+ Purchase Mari 8370	0 lied Not_Marrie	d Female	7969 Male Ci 0 0	ty_C	1	1	1
	: # df df :	.drop([.head() .head() .drop([.head() .head() .drop([.head() .head() .head()	ng redundant 'Gender',' Product_ID 1 P00069042 1 P00248942 1 P00087842	Age CO-17 O-17 O-	ns. _Statu Occupati	s','City_ on Stay_In 10	Category'],a	xis=1,inplace=True) Years Product_Category_1 2 3 2 1	Purchase Mari 8370	ied Not_Marrie	d Female 0 1 0 1	7969 Male Ci 0 0 0	ty_C	1	1	1
	: # dff dff:::	.drop([.head()	ng redundant 'Gender',' Product_ID 1 P00069042 1 P00248942 1 P00087842	Age C 0-17 0-17 0-17	ns. _Statu Occupati	s','City_ on Stay_In 10 10	Category'],a	xis=1,inplace=True) ears Product_Category_1 2	Purchase Mari 8370 15200	0 ied Not_Marrie 1 1 1	8 d Female 0 1 0 1	7969 Ci	ty_C	1	1 1 1	1 1 1

```
In [11]: df['Purchase'].describe()
Out[11]: count
                   550068,000000
          mean
                     9263.968713
          std
                     5023.065394
          min
                       12.000000
          25%
                     5823.000000
          50%
                     8047.000000
          75%
                    12054.000000
                    23961.000000
          max
          Name: Purchase, dtype: float64
```

5. Find out outliers (manually).

```
In [12]: # Calculating the Outliers
   interQuartileRange = (df['Purchase'].describe()[6] - df['Purchase'].describe()[4])
   upper = df['Purchase'].describe()[6] + (1.5 * interQuartileRange)
   lower = df['Purchase'].describe()[4] - (1.5 * interQuartileRange)

   print("IQR (Q3-Q1): " + str(interQuartileRange))
   print("Lower: " + str(lower))
   print("Upper: " + str(upper))

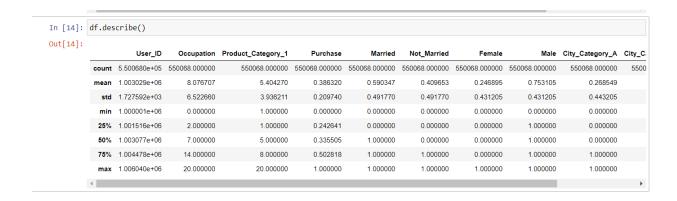
outliers = []
   for i in range(df['Purchase'].shape[0]):
        if(df['Purchase'].iloc[i])upper):
            outliers.append([i,df['Purchase'].iloc[i]])

print("Total number of outliers: " + str(len(outliers)))

IQR (Q3-Q1): 6231.0
   Lower: -3523.5
   Upper: 21400.5
   Total number of outliers: 2677
```

6. Standardization and normalization of columns.

```
In [13]: # Normalization on Purchase column.
          # normal_value = (value - min)/(max - min)
          df.Purchase = (df.Purchase - df.Purchase.min()) / (df.Purchase.max() - df.Purchase.min())
         df.head()
Out[13]:
             User_ID Product_ID Age Occupation Stay_In_Current_City_Years Product_Category_1 Purchase Married Not_Married Female Male City_Category_A Ci
          0 1000001 P00069042
                                            10
                                                                                      3 0.348992
                                                                                                                              0
           1 1000001 P00248942
                                            10
                                                                    2
                                                                                      1 0.634181
                                                                                                                              0
          2 1000001 P00087842
                                                                                     12 0.058875
                                            10
           3 1000001 P00085442
                                            10
                                                                    2
                                                                                     12 0.043634
          4 1000002 P00285442 55+
                                                                                      8 0.332248
```



Conclusion-

We have successfully preprocessed the data. We cleaned the dataset by removing missing values and redundant columns. We also used dummy variables and then normalized the values.