

```
In [1]: import warnings
warnings.filterwarnings('ignore')
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
```

```
In [2]: df2=pd.ExcelFile('KPMG_Final.xlsx')
df2=pd.read_excel(df,'Transactions')
df2=pd.read_excel(df,'NewCustomerList')
df2=pd.read_excel(df,'CustomerDemographic')
df2=pd.read_excel(df,'CustomerAddress')
```

Data Quality issues in df1 - Transactions

```
In [3]: df1.to_csv('Transactions.csv')
```

```
In [4]: df1.head()
```

```
Out[4]:
```

	transaction_id	product_id	customer_id	transaction_date	online_order	order_status	brand	product_line	product_class	product_size	list_price	standard_cost	product_first_sold_date	
0	1	2	3	2020	2017-02-25	0.0	Approved	Solex	Standard	medium	medium	71.49	53.62	41245.0
1	2	3	3	3120	2017-05-21	1.0	Approved	Trek Bicycles	Standard	medium	large	1793.43	388.92	41701.0
2	3	37	37	402	2017-10-16	0.0	Approved	CHM Cycles	Standard	low	medium	2093.43	248.82	36361.0
3	4	88	3135	2017-08-31	0.0	Approved	Norco Bicycles	Standard	medium	medium	1198.46	381.10	36145.0	
4	5	78	787	2017-10-01	1.0	Approved	Giant Bicycles	Standard	medium	large	1765.30	709.48	42226.0	

```
In [5]: df1.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 20000 entries, 0 to 19999
Data columns (total 13 columns):
#   Column              Non-Null Count  Dtype
---  --
0   transaction_id       20000 non-null  int64
1   product_id          20000 non-null  int64
2   customer_id         20000 non-null  int64
3   transaction_date     20000 non-null  datetime64[ns]
4   online_order         19648 non-null  float64
5   order_status        20000 non-null  object
6   brand               19883 non-null  object
7   product_line        18893 non-null  object
8   product_class       19893 non-null  object
9   product_size        19893 non-null  object
10  list_price          20000 non-null  float64
11  standard_cost       19893 non-null  float64
12  product_first_sold_date 19893 non-null  float64
dtypes: datetime64[ns](1), float64(4), int64(3), object(5)
memory usage: 2.0+ MB
```

```
In [6]: #shape of the data
df1.shape
```

```
Out[6]: (20000, 13)
```

```
In [7]: #check for null values
df1.isnull().sum()
```

```
Out[7]:
```

transaction_id	0
product_id	0
customer_id	0
transaction_date	0
online_order	360
order_status	0
brand	197
product_line	197
product_class	197
product_size	197
list_price	0
standard_cost	197
product_first_sold_date	197
dtype:	int64

online_order,brand,product_line,product_class,product_size,standard_cost,product_first_sold_date are the columns with missing values.Lets do some analysis in those columns.

```
In [8]: #check if there are any duplicate values
df1.duplicated().sum()
```

```
Out[8]: 0
```

There are no duplicate values in the transactions dataframe which means that the dataset is unique.

```
In [9]: #uniqueness of the data
df1.nunique()
```

```
Out[9]:
```

transaction_id	20000
product_id	481
customer_id	3494
transaction_date	364
online_order	2
order_status	2
brand	6
product_line	4
product_class	3
product_size	3
list_price	296
standard_cost	183
product_first_sold_date	180
dtype:	int64

```
In [10]: #exploring and understanding the columns
df1.columns
```

```
Out[10]: Index(['transaction_id', 'product_id', 'customer_id', 'transaction_date',
              'online_order', 'order_status', 'brand', 'product_line',
              'product_class', 'product_size', 'list_price', 'standard_cost',
              'product_first_sold_date'],
              dtype='object')
```

```
In [11]: #number of approved and cancelled orders
df1['order_status'].value_counts()
```

```
Out[11]:
```

Approved	19821
Cancelled	179
Name: order_status, dtype: int64	

```
In [12]: #different brands
df1['brand'].value_counts()
```

```
Out[12]:
```

Solex	4253
Giant Bicycles	3312
Warrior	3295
CHM Cycles	3643
Trek Bicycles	2980
Norco Bicycles	2910
Name: brand, dtype: int64	

```
In [13]: #different product lines
df1['product_line'].value_counts()
```

```
Out[13]:
```

Standard	14376
Road	3970
Touring	1234
Mountain	423
Name: product_line, dtype: int64	

```
In [14]: #different product classes
df1['product_class'].value_counts()
```

```
Out[14]:
```

medium	13826
high	3613
low	2964
Name: product_class, dtype: int64	

```
In [15]: #different product sizes
df1['product_size'].value_counts()
```

```
Out[15]:
```

medium	12999
large	3976
small	2837
Name: product_size, dtype: int64	

```
In [16]: #accuracy check
df1['product_first_sold_date'].value_counts()
```

```
Out[16]:
```

33879.0	234
41064.0	229
37823.0	227
39889.0	222
38216.0	220
...	...
41848.0	169
42844.0	168
41922.0	166
37659.0	163
34586.0	162
Name: product_first_sold_date, Length: 180, dtype: int64	

```
In [17]: #product_first_sold_date is in the form of numbers.Let's convert it into date.
df1['product_first_sold_date'] = pd.to_datetime(df1['product_first_sold_date'],unit='s')
```

```
In [18]: df1['product_first_sold_date'].head()
```

```
Out[18]:
```

0	1970-01-01 11:27:25
1	1970-01-01 11:35:01
2	1970-01-01 18:08:01
3	1970-01-01 10:06:01
4	1970-01-01 11:43:46
Name: product_first_sold_date, dtype: datetime64[ns]	

```
In [19]: df1['product_first_sold_date'].head(25)
```

```
Out[19]:
```

0	1970-01-01 11:27:25
1	1970-01-01 11:35:01
2	1970-01-01 18:08:01
3	1970-01-01 10:06:25
4	1970-01-01 11:43:46
5	1970-01-01 18:50:31
6	1970-01-01 09:29:26
7	1970-01-01 11:05:15
8	1970-01-01 09:17:35
9	1970-01-01 18:36:56
10	1970-01-01 11:19:44
11	1970-01-01 11:42:52
12	1970-01-01 09:25:27
13	1970-01-01 09:36:26
14	1970-01-01 18:36:33
15	1970-01-01 18:31:13
16	1970-01-01 18:30:46
17	1970-01-01 09:24:48
18	1970-01-01 11:05:15
19	1970-01-01 18:22:17
20	1970-01-01 18:05:34
21	1970-01-01 11:42:52
22	1970-01-01 11:42:25
23	1970-01-01 11:40:44
24	1970-01-01 09:27:59
Name: product_first_sold_date, dtype: datetime64[ns]	

The column 'product_first_sold_date' seems to be wrong since it shows the same date with different timestamps.This looks inaccurate.

Data Quality issues in df2- NewCustomerList

```
In [20]: df2.to_csv('NewCustomerList.csv')
```

```
In [21]: df2.head()
```

```
Out[21]:
```

	first_name	last_name	gender	past_3_years_bike_related_purchases	DOB	job_title	job_industry_category	wealth_segment	deceased_indicator	owns_car	...	state	country	property_valuation	Unnam	
0	Chickie	Brisler	Male	86	1957-07-22	General Manager	Manufacturing	Mass Customer		N	Yes	...	QLD	Australia	6	0
1	Morly	Genery	Male	69	1970-03-22	Structural Engineer	Property	Mass Customer		N	No	...	NSW	Australia	11	0
2	Adellis	Fomerster	Female	10	1974-08-28	Senior Cost Accountant	Financial Services	Affluent Customer		N	No	...	VIC	Australia	5	0
3	Lucine	Stutt	Female	64	1979-01-28	Account Representative III	Manufacturing	Affluent Customer		N	Yes	...	QLD	Australia	1	0
4	Melinda	Hadlee	Female	34	1965-09-21	Financial Analyst	Financial Services	Affluent Customer		N	No	...	NSW	Australia	9	0

5 rows × 23 columns

```
In [22]: df2.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1900 entries, 0 to 999
Data columns (total 23 columns):
#   Column              Non-Null Count  Dtype
---  --
0   first_name          1900 non-null  object
1   last_name           1900 non-null  object
2   gender              1900 non-null  object
3   past_3_years_bike_related_purchases 1900 non-null  int64
4   DOB                 1900 non-null  object
5   job_title           894 non-null   object
6   job_industry_category 835 non-null   object
7   wealth_segment      1900 non-null  object
8   deceased_indicator   1900 non-null  object
9   owns_car            1900 non-null  object
10  tenure              1900 non-null  int64
11  address             1900 non-null  object
12  postcode            1900 non-null  int64
13  state               1900 non-null  object
14  country             1900 non-null  object
15  property_valuation   1900 non-null  int64
16  Unnamed: 16         1900 non-null  float64
17  Unnamed: 17         1900 non-null  float64
18  Unnamed: 18         1900 non-null  float64
19  Unnamed: 19         1900 non-null  float64
20  Unnamed: 20         1900 non-null  int64
21  Rank                1900 non-null  int64
22  Value               1900 non-null  float64
dtypes: datetime64[ns](1), float64(5), int64(8), object(11)
memory usage: 179.8+ KB
```

```
In [23]: #lets remove the unwanted columns in the dataset
df2.drop(['Unnamed: 16', 'Unnamed: 17', 'Unnamed: 18', 'Unnamed: 19', 'Unnamed: 20'],axis=1,inplace=True)
```

```
Out[23]: df2.shape
```

```
Out[24]: (1900, 18)
```

```
In [25]: df2.isnull().sum()
```

```
Out[25]:
```

first_name	0
last_name	29
gender	0
past_3_years_bike_related_purchases	0
DOB	17
job_title	106
job_industry_category	165
wealth_segment	0
deceased_indicator	0
owns_car	0
tenure	0
address	1900
postcode	522
state	1
country	0
property_valuation	12
Rank	324
Value	0
dtype:	int64

last_name,DOB,job_title,job_industry_category are the columns that has missing values. Lets perform analysis on these columns of the dataset.

```
In [26]: df2.duplicated().sum()
```

```
Out[26]: 0
```

The NewCustomerList doesnt have any duplicated values which means that the dataset is unique.

```
In [27]: df2.nunique()
```

```
Out[27]:
```

first_name	940
last_name	3
gender	3
past_3_years_bike_related_purchases	190
DOB	95
job_title	184
job_industry_category	1
wealth_segment	3
deceased_indicator	2
owns_car	2
tenure	25
address	1900
postcode	522
state	1
country	0
property_valuation	12
Rank	324
Value	324
dtype:	int64

```
In [28]: #columns of NewCustomerList
df2.columns
```

```
Out[28]: Index(['first_name', 'last_name', 'gender',
              'past_3_years_bike_related_purchases', 'DOB', 'job_title',
              'job_industry_category', 'wealth_segment', 'deceased_indicator',
              'owns_car', 'tenure', 'address', 'postcode', 'state', 'country',
              'property_valuation', 'Rank', 'Value'],
              dtype='object')
```

```
In [29]: #check the column gender since it says that the unique values are 3.
df2['gender'].value_counts()
```

```
Out[29]:
```

Female	513
Male	478
U	17
Name: gender, dtype: int64	

```
In [30]: #lets remove the values that has 'U' in it since it doesnt apply to any gender and is irrelevant.
df2 = df2[df2['gender'].str.contains('U') == False]
```

```
Out[30]:
```

Female	513
Male	478
Name: gender, dtype: int64	

```
In [31]: #check if the column values of U are removed
df2['gender'].value_counts()
```

```
Out[31]:
```

Female	513
Male	478
Name: gender, dtype: int64	

```
In [32]: #different job industry categories
df2['job_industry_category'].value_counts()
```

```
Out[32]:
```

Financial Services	202
Manufacturing	199
Health	152
Retail	78
Property	64
IT	38
Entertainment	36
Agriculture	30
Telecommunications	25
Name: job_industry_category, dtype: int64	

```
In [33]: #different wealth segments
df2['wealth_segment'].value_counts()
```

```
Out[33]:
```

Mass Customer	499
High Net Worth	249
Affluent Customer	235
Name: wealth_segment, dtype: int64	

```
In [34]: #different states
df2['state'].value_counts()
```

```
Out[34]:
```

NSW	495
VIC	258
QLD	226
Name: state, dtype: int64	

Data Quality issues in df3-CustomerDemographic

```
In [35]: df3.to_csv('CustomerDemographic.csv')
```

```
In [36]: df3.head()
```

```
Out[36]:
```

	customer_id	first_name	last_name	gender	past_3_years_bike_related_purchases	DOB	job_title	job_industry_category	wealth_segment	deceased_indicator	default	owns_car	tenure	
0	1	Laraine	Medendorp	F		93	1953-10-12	Executive Secretary	Health	Mass Customer	N	-	Yes	11.0
1	2	Eli	Bockman	Male		81	1980-12-16	Administrative Officer	Financial Services	Mass Customer	N	<script>alert('X')</script>	Yes	16.0
2	3	Arlin	Dearle	Male		61	1954-01-20	Recruiting Manager	Property	Mass Customer	N	2018-02-01 00:00:00	Yes	15.0
3	4	Talbot	NaN	Male		33	1961-10-03	NaN	IT	Mass Customer	N	01-...>,{00{ touch /mp/bhns.shellsh...	No	7.0
4	5	Sheila-kathryn	Calton	Female		56	1977-05-13	Senior Editor	NaN	Affluent Customer	N	NaN	Yes	8.0

```
In [37]: df3.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4000 entries, 0 to 3999
Data columns (total 13 columns):
#   Column              Non-Null Count  Dtype
---  --
0   customer_id         4000 non-null  int64
1   first_name          4000 non-null  object
2   last_name           3999 non-null  object
3   gender              4000 non-null  object
4   past_3_years_bike_related_purchases 4000 non-null  int64
5   DOB                 3913 non-null  object
6   job_title           3494 non-null  object
7   job_industry_category 3044 non-null  object
8   wealth_segment      4000 non-null  object
9   deceased_indicator   4000 non-null  object
10  default             3999 non-null  object
11  owns_car            4000 non-null  object
12  tenure              3913 non-null  float64
dtypes: datetime64[ns](1), float64(1), int64(2), object(9)
memory usage: 406.4+ KB
```

```
In [38]: df3.shape
```

```
Out[38]: (4000, 13)
```

```
In [39]: df3.isnull().sum()
```

```
Out[39]:
```

customer_id	0
first_name	0
last_name	125
gender	0
past_3_years_bike_related_purchases	0
DOB	87
job_title	508
job_industry_category	656
wealth_segment	0
deceased_indicator	0
default	382
owns_car	0
tenure	87
dtype:	int64

last_name,DOB,job_title,job_industry_category,default and tenure are the columns with null values. lets Perform some analysis on these columns.

```
In [40]: df3.duplicated().sum()
```

```
Out[40]: 0
```

df3 doesnt have any duplicate value so its unique dataset.

```
In [41]: df3.nunique()
```

```
Out[41]:
```

customer_id	4000
first_name	3139
last_name	3725
gender	6
past_3_years_bike_related_purchases	190
DOB	3445
job_title	195
job_industry_category	3
wealth_segment	3
deceased_indicator	0
owns_car	99
owns_car	2
tenure	22
dtype:	int64

```
In [42]: #columns of CustomerDemographic
df3.columns
```

```
Out[42]: Index(['customer_id', 'first_name', 'last_name', 'gender',
              'past_3_years_bike_related_purchases', 'DOB', 'job_title',
              'job_industry_category', 'wealth_segment', 'deceased_indicator',
              'default', 'owns_car', 'tenure'],
              dtype='object')
```

```
In [43]: #df3 has 6 different genders mentioned.Lets analyse and drop the unwanted values
df3.gender.value_counts()
```

```
Out[43]:
```

Female	2037
Male	1872
U	88
F	1
M	1
Name: gender, dtype: int64	

```
In [44]: #U is not a gender.Lets drop the rows with U values.
df3 = df3[df3['gender'].str.contains('U') == False]
```

```
Out[44]:
```

Female	2037
Male	1873
Name: gender, dtype: int64	

```
In [45]: #F and Male refers to female so lets tag them into Female and M to Male
df3['gender'] = df3['gender'].replace('F', 'Female').replace('M', 'Male')
```

```
Out[45]:
```

Female	2037
Male	1873
Name: gender, dtype: int64	

```
In [46]: df3.gender.value_counts()
```

```
Out[46]:
```

Female	2039
Male	1873
Name: gender, dtype: int64	

```
In [47]: df3.wealth_segment.value_counts()
```

```
Out[47]:
```

Mass Customer
