

```
In [1]: # Importing the libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

In [2]:
# Importing the data
retail_df = pd.read_excel("Online Retail.xlsx")

In [3]:
retail_df.head(10)
```

InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country
0	536365	85123A WHITE HANGING HEART T-LIGHT HOLDER	6	2010-12-01 08:26:00	2.55	178500	United Kingdom
1	536365	71053 WHITE METAL LANTERN	6	2010-12-01 08:26:00	3.39	178500	United Kingdom
2	536365	844068 CREAM FLUID HEARTS COAT HANGER	8	2010-12-01 08:26:00	2.75	178500	United Kingdom
2	536365	840296 KNITTED UNION FLAG HOT WATER BOTTLE	6	2010-12-01 08:26:00	3.39	178500	United Kingdom
4	536365	840296 RED WOOLLY HOTTIE WHITE HEART.	6	2010-12-01 08:26:00	3.39	178500	United Kingdom
5	536365	22752 SET 7 BABUSHKA NESTING BOXES	2	2010-12-01 08:26:00	7.65	178500	United Kingdom
6	536365	21730 GLASS STAR FROSTED T-LIGHT HOLDER	6	2010-12-01 08:26:00	4.25	178500	United Kingdom
7	536366	22633 HAND WARMER UNION JACK	6	2010-12-01 08:28:00	1.85	178500	United Kingdom
8	536366	22632 HAND WARMER RED POLKA DOT	6	2010-12-01 08:28:00	1.85	178500	United Kingdom
9	536367	84879 ASSORTED COLOUR BIRD ORNAMENT	32	2010-12-01 08:34:00	1.69	130470	United Kingdom

1. Data Inspection and data cleaning

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

Out[6]:
(54199, 8)
```

```
In [7]: # data info
retail_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 54199 entries, 0 to 54198
Data columns (total 8 columns):
#   Column                Non-Null Count  Dtype
---  --
0   InvoiceNo              54199 non-null  object
1   StockCode             54199 non-null  object
2   Description            54045 non-null  object
3   Quantity              54199 non-null  int64
4   InvoiceDate            54199 non-null  datetime64[ns]
5   UnitPrice             54199 non-null  float64
6   CustomerID            40682 non-null  float64
7   Country               54199 non-null  object
dtypes: datetime64[ns](1), float64(2), int64(1), object(4)
memory usage: 33.1+ MB
```

```
In [8]: # Checking the null values for missing value treatment
retail_df.isnull().sum()
```

```
InvoiceNo      0
StockCode      0
Description    1454
Quantity        0
InvoiceDate    0
UnitPrice      0
CustomerID    135080
Country        0
dtype: int64
```

```
In [15]: # Checking the percentage of the customer ID
(retail_df["CustomerID"].isnull().sum()/len(retail_df)*100)
```

```
Out[15]:
24.9264943288598
```

The description column can be dropped as it has no use in the model

```
In [16]: retail_df=retail_df.drop("Description",axis=1)
```

```
In [17]: retail_df.head(5)
```

InvoiceNo	StockCode	Quantity	InvoiceDate	UnitPrice	CustomerID	Country
0	536365	85123A	6	2010-12-01 08:26:00	2.55	178500 United Kingdom
1	536365	71053	6	2010-12-01 08:26:00	3.39	178500 United Kingdom
2	536365	844068	8	2010-12-01 08:26:00	2.75	178500 United Kingdom
3	536365	840296	6	2010-12-01 08:26:00	3.39	178500 United Kingdom
4	536365	840296	6	2010-12-01 08:26:00	3.39	178500 United Kingdom

```
In [23]: # Missing value treatment for CustomerID
# The customer ID is an unique class and hence cannot be replaced with mean/mode/median.The null values needs to be dropped.

In [28]: retail_df=retail_df.dropna()
```

```
In [29]: retail_df.isnull().sum()
```

```
InvoiceNo      0
StockCode      0
Quantity        0
InvoiceDate    0
UnitPrice      0
CustomerID    135080
Country        0
dtype: int64
```

```
In [30]: # Dropping the duplicate entries in the dataset
retail_df=retail_df.drop_duplicates()
retail_df.shape
```

```
Out[30]:
(401402, 7)
```

Performing Descriptive Analysis

```
In [32]: # Customer ID is of float type and this needs to be changed to string type
retail_df["CustomerID"]=retail_df["CustomerID"].astype(str)
```

```
In [34]: retail_df.describe(datetime_is_numeric=True)
```

	Quantity	InvoiceDate	UnitPrice
count	401602.000000	401602	401602.000000
mean	12.182579	2011-07-10 12:08:08.129839872	3.474064
min	-80995.000000	2010-12-01 08:26:00	0.000000
25%	2.000000	2011-04-06 15:02:00	1.250000
50%	5.000000	2011-07-29 15:40:00	1.950000
75%	12.000000	2011-10-20 11:58:00	3.750000
max	80995.000000	2011-12-09 12:50:00	38970.000000
std	250.283248	NaN	69.764209

Inferences

1. The minimum quantity is in negative which implies returns
2. The invoice date ranges from 2010-12-01 to 2010-12-01
3. Average price of each product is 3.47

```
In [39]: retail_df.describe(include="O")
```

	InvoiceNo	StockCode	CustomerID	Country
count	401602	401602	401602	401602
unique	22190	3684	4372	37
top	576339	85123A	178410	United Kingdom
freq	542	2065	7812	356726

Inferences

1. The count of invoice numbers are 401602 while unique transactions are 22190.
2. The invoice no with maximum frequency of 542 is "576339"
3. The stock code of 85123A is of max frequency of 2065.
4. The country with max transaction is UK and there are total 37 countries

Data Transformation

Cohort Analysis (a) Create month cohort of customers and analyze active customers in each cohort:

```
In [41]: retail_df["month_year"] = retail_df["InvoiceDate"].dt.to_period('M')
retail_df["month_year"].nunique()
```

```
Out[41]:
13
```

```
In [42]: month_cohort = retail_df.groupby('month_year')['CustomerID'].nunique()
month_cohort
```

```
Out[42]:
month_year
2010-12    948
2011-01    783
2011-02    798
2011-03   1020
2011-04    899
2011-05   1079
2011-06   1051
2011-07    993
2011-08    980
2011-09   1302
2011-10   1425
2011-11   1711
2011-12    686
Freq: M, Name: CustomerID, dtype: int64
```

```
In [43]: plt.figure(figsize=(10,5))
sns.barplot(y = month_cohort.index, x = month_cohort.values);
plt.xlabel("Count of customers")
plt.title("No. of active customers in each month")
```

```
Out[43]:
Text(0.5, 1.0, 'No. of active customers in each month')
```

(b) Analyze the retention rate of customers:

```
In [45]: month_cohort = month_cohort.shift()
```

```
Out[45]:
month_year
2010-12    NaN
2011-01   -165.0
2011-02    15.0
2011-03   222.0
2011-04   -121.0
2011-05    180.0
2011-06   -28.0
2011-07   -58.0
2011-08   -13.0
2011-09   322.0
2011-10   123.0
2011-11   286.0
2011-12  -102.0
Freq: M, Name: CustomerID, dtype: float64
```

```
In [46]: retention_rate = round(month_cohort.pct_change(periods=1)*100,2)
retention_rate
```

```
Out[46]:
month_year
2010-12    NaN
2011-01   -17.41
2011-02    1.92
2011-03   27.82
2011-04   -11.86
2011-05    20.02
2011-06   -2.59
2011-07   -5.52
2011-08   -1.31
2011-09   32.86
2011-10    9.45
2011-11   20.07
2011-12  -59.91
Freq: M, Name: CustomerID, dtype: float64
```

```
In [47]: plt.figure(figsize=(10,5))
sns.barplot(y = retention_rate.index, x = retention_rate.values);
plt.xlabel("Retention (in %)")
plt.title("Month-wise customer retention rate");
```

RFM ANALYSIS - RECENCY FREQUENCY MONETORY ANALYSIS

```
In [ ]: # MONETARY ANALYSIS
```

```
In [58]: retail_df["amount"] = retail_df["Quantity"]*retail_df["UnitPrice"]
retail_df.head()
```

InvoiceNo	StockCode	Quantity	InvoiceDate	UnitPrice	CustomerID	Country	month_year	amount	
0	536365	85123A	6	2010-12-01 08:26:00	2.55	178500	United Kingdom	2010-12	15.30
1	536365	71053	6	2010-12-01 08:26:00	3.39	178500	United Kingdom	2010-12	20.34
2	536365	844068	8	2010-12-01 08:26:00	2.75	178500	United Kingdom	2010-12	22.00
3	536365	840296	6	2010-12-01 08:26:00	3.39	178500	United Kingdom	2010-12	20.34
4	536365	840296	6	2010-12-01 08:26:00	3.39	178500	United Kingdom	2010-12	20.34

```
In [51]: df_monetary = retail_df.groupby('CustomerID').sum().reset_index()
df_monetary
```

CustomerID	amount
0	12346.0
1	12347.0
2	12348.0
3	12349.0
4	12350.0
...	...
4367	18280.0
4368	18281.0
4369	18282.0
4370	18283.0
4371	18287.0

4372 rows x 2 columns

```
In [ ]: # FREQUENCY ANALYSIS
```

```
In [53]: df_frequency = retail_df.groupby('CustomerID').nunique().reset_index()
df_frequency
```

CustomerID	InvoiceNo
0	12346.0
1	12347.0
2	12348.0
3	12349.0
4	12350.0
...	...
4367	18280.0
4368	18281.0
4369	18282.0
4370	18283.0
4371	18287.0

4372 rows x 2 columns

```
In [55]: # We will fix reference date for calculating recency as last transaction day in data + 1 day
retail_df["days_to_last_order"] = (retail_df["InvoiceDate"] + timedelta(days=1) - ref_day - retail_df["InvoiceDate"]).dt.days
retail_df.head()
```

InvoiceNo	StockCode	Quantity	InvoiceDate	UnitPrice	CustomerID	Country	month_year	days_to_last_order	
0	536365	85123A	6	2010-12-01 08:26:00	2.55	178500	United Kingdom	2010-12	15.30
1	536365	71053	6	2010-12-01 08:26:00	3.39	178500	United Kingdom	2010-12	20.34
2	536365	844068	8	2010-12-01 08:26:00	2.75	178500	United Kingdom	2010-12	22.00
3	536365	840296	6	2010-12-01 08:26:00	3.39	178500	United Kingdom	2010-12	20.34
4	536365	840296	6	2010-12-01 08:26:00	3.39	178500	United Kingdom	2010-12	20.34

```
In [61]: df_recency = retail_df.groupby('CustomerID')['days_to_last_order'].min().reset_index()
df_recency
```

CustomerID	days_to_last_order
0	12346.0
1	12347.0
2	12348.0
3	12349.0
4	12350.0
...	...
4367	18280.0
4368	18281.0
4369	18282.0
4370	18283.0
4371	18287.0

4372 rows x 2 columns

Calculating RFM Matrix

```
In [63]: df_rf = pd.merge(df_recency, df_frequency, on='CustomerID', how='inner')
df_rf = pd.merge(df_rf, df_monetary, on='CustomerID', how='inner')
df_rf.columns = ['CustomerID', 'Recency', 'Frequency', 'Monetary']
df_rf.head()
```

CustomerID	Recen	Frequency	Monetary
0	12346.0	326	2
1	12347.0	2	7
2	12348.0	75	4
3	12349.0	19	1
4	12350.0	310	1
...
4367	18280.0	180	6
4368	18281.0	80	82
4369	18282.0	176	60
4370	18283.0	2045	53
4371	18287.0	1837	28

4372 rows x 4 columns

```
In [64]: df_rf["recency_labels"] = pd.cut(df_rf["Recency"], bins=5, labels=['newest', 'newer', 'medium', 'older', 'oldest'])
df_rf["frequency_labels"] = pd.cut(df_rf["Frequency"], bins=5, labels=['lowest', 'lower', 'medium', 'higher', 'highest'])
df_rf["monetary_labels"] = pd.cut(df_rf["Monetary"], bins=5, labels=['smallest', 'smaller', 'medium', 'larger', 'largest'])
df_rf.head()
```

CustomerID	Recen	Frequency	Monetary
0	12346.0	326	2
1	12347.0	2	7
2	12348.0	75	4
3	12349.0	19	1
4	12350.0	310	1
...
4367	18280.0	180	6
4368	18281.0	80	82
4369	18282.0	176	60
4370	18283.0	2045	53
4371	18287.0	1837	28

4372 rows x 4 columns

RFM ANALYSIS - RECENCY FREQUENCY MONETORY ANALYSIS

```
In [53]: df_frequency = retail_df.groupby('CustomerID').nunique().reset_index()
df_frequency
```

CustomerID	InvoiceNo
0	12346.0
1	12347.0
2	12348.0
3	12349.0
4	12350.0
...	...
4367	18280.0
4368	18281.0
4369	18282.0
4370	18283.0
4371	18287.0

4372 rows x 2 columns

```
In [55]: # We will fix reference date for calculating recency as last transaction day in data + 1 day
retail_df["days_to_last_order"] = (retail_df["InvoiceDate"] + timedelta(days=1) - ref_day - retail_df["InvoiceDate"]).dt.days
retail_df.head()
```

InvoiceNo	StockCode	Quantity	InvoiceDate	UnitPrice	CustomerID	Country	month_year	days_to_last_order	
0	536365	85123A	6	2010-12-01 08:26:00	2.55	178500	United Kingdom	2010-12	15.30
1	536365	71053	6	2010-12-01 08:26:00	3.39	178500	United Kingdom	2010-12	20.34
2	536365	844068	8	2010-12-01 08:26:00	2.75	178500	United Kingdom	2010-12	22.00
3	536365	840296	6	2010-12-01 08:26:00	3.39	178500	United Kingdom	2010-12	20.34
4	536365	840296	6	2010-12-01 08:26:00	3.39	178500	United Kingdom	2010-12	20.34

```
In [61]: df_recency = retail_df.groupby('CustomerID')['days_to_last_order'].min().reset_index()
df_recency
```

CustomerID	days_to_last_order
0	12346.0
1	12347.0
2	12348.0
3	12349.0
4	12350.0
...	...
4367	18280.0
4368	18281.0
4369	18282.0
4370	18283.0
4371	18287.0

4372 rows x 2 columns

Calculating RFM Matrix

```
In [63]: df_rf = pd.merge(df_recency, df_frequency, on='CustomerID', how='inner')
df_rf = pd.merge(df_rf, df_monetary, on='CustomerID', how='inner')
df_rf.columns = ['CustomerID', 'Recency', 'Frequency', 'Monetary']
df_rf.head()
```

CustomerID	Recen	Frequency	Monetary
0	12346.0	326	2
1	12347.0	2	7
2	12348.0	75	4
3	12349.0	19	1
4	12350.0	310	1
...
4367	18280.0	180	6
4368	18281.0	80	82
4369	18282.0	176	60
4370	18283.0	2045	53
4371	18287.0	1837	28

4372 rows x 4 columns

```
In [64]: df_rf["recency_labels"] = pd.cut(df_rf["Recency"], bins=5, labels=['newest', 'newer', 'medium', 'older', 'oldest'])
df_rf["frequency_labels"] = pd.cut(df_rf["Frequency"], bins=5, labels=['lowest', 'lower', 'medium', 'higher', 'highest'])
df_rf["monetary_labels"] = pd.cut(df_rf["Monetary"], bins=5, labels=['smallest', 'smaller', 'medium', 'larger', 'largest'])
df_rf.head()
```

CustomerID	Recen	Frequency	Monetary
0	12346.0	326	2
1	12347.0	2	7
2	12348.0	75	4
3	12349.0	19	1
4	12350.0	310	1
...
4367	18280.0	180	6
4368	18281.0	80	82
4369	18282.0	176	60
4370	18283.0	2045	53
4371	18287.0	1837	28

4372 rows x 4 columns

RFM ANALYSIS - RECENCY FREQUENCY MONETORY ANALYSIS

```
In [53]: df_frequency = retail_df.groupby('CustomerID').nunique().reset_index()
df_frequency
```

CustomerID	InvoiceNo
0	12346.0
1	12347.0
2	12348.0
3	12349.0
4	12350.0
...	...
4367	18280.0
4368	18281.0
4369	18282.0
4370	18283.0
4371	18287.0

4372 rows x 2 columns

```
In [55]: # We will fix reference date for calculating recency as last transaction day in data + 1 day
retail_df["days_to_last_order"] = (retail_df["InvoiceDate"] + timedelta(days=1) - ref_day - retail_df["InvoiceDate"]).dt.days
retail_df.head()
```

InvoiceNo	StockCode	Quantity	InvoiceDate	UnitPrice	CustomerID	Country	month_year	days_to_last_order	
0	536365	85123A	6	2010-12-01 08:26:00	2.55	178500	United Kingdom	2010-12	15.30
1	536365	71053	6	2010-12-01 08:26:00	3.39	178500	United Kingdom	2010-12	20.34
2	536365	844068	8	2010-12-01 08:26:00	2.75	178500	United Kingdom	2010-12	22.00
3	536365	840296	6	2010-12-01 08:26:00	3.39	178500	United Kingdom	2010-12	20.34
4	536365	840296	6	2010-12-01 08:26:00	3.39	178500	United Kingdom	2010-12	20.34

```
In [61]: df_rec
```


Out[92]:	clusters		inertia
	0	2	7113.097396
	1	3	5343.153295
	2	4	4480.998412
	3	5	3730.755397
	4	6	3044.935139
	5	7	2598.297835
	6	8	2299.260286
	7	9	2045.630814
	8	10	1852.933027
	9	11	1701.069737
	10	12	1574.969301

```
In [93]: # Finding the Optimal Number of Clusters with the help of Silhouette Analysis
range_n_clusters = [2, 3, 4, 5, 6, 7, 8, 9, 10]

for num_clusters in range_n_clusters:
    kmeans = KMeans(n_clusters=num_clusters, max_iter=50)
    kmeans.fit(df_rfm_scaled)

    cluster_labels = kmeans.labels_

    silhouette_avg = silhouette_score(df_rfm_scaled, cluster_labels)
    print("For n_clusters={0}, the silhouette score is {1}".format(num_clusters, silhouette_avg))

For n_clusters=2, the silhouette score is 0.44132753537785846
For n_clusters=3, the silhouette score is 0.3803925903795938
For n_clusters=4, the silhouette score is 0.3623606424972478
For n_clusters=5, the silhouette score is 0.3648879010841238
For n_clusters=6, the silhouette score is 0.3443420328414096
For n_clusters=7, the silhouette score is 0.34297910328413955
For n_clusters=8, the silhouette score is 0.3356899635770755
For n_clusters=9, the silhouette score is 0.3471822122915281
For n_clusters=10, the silhouette score is 0.35608282484485043
```

```
In [94]: # Final model with k=3
kmeans = KMeans(n_clusters=3, max_iter=50)
kmeans.fit(df_rfm_scaled)
```

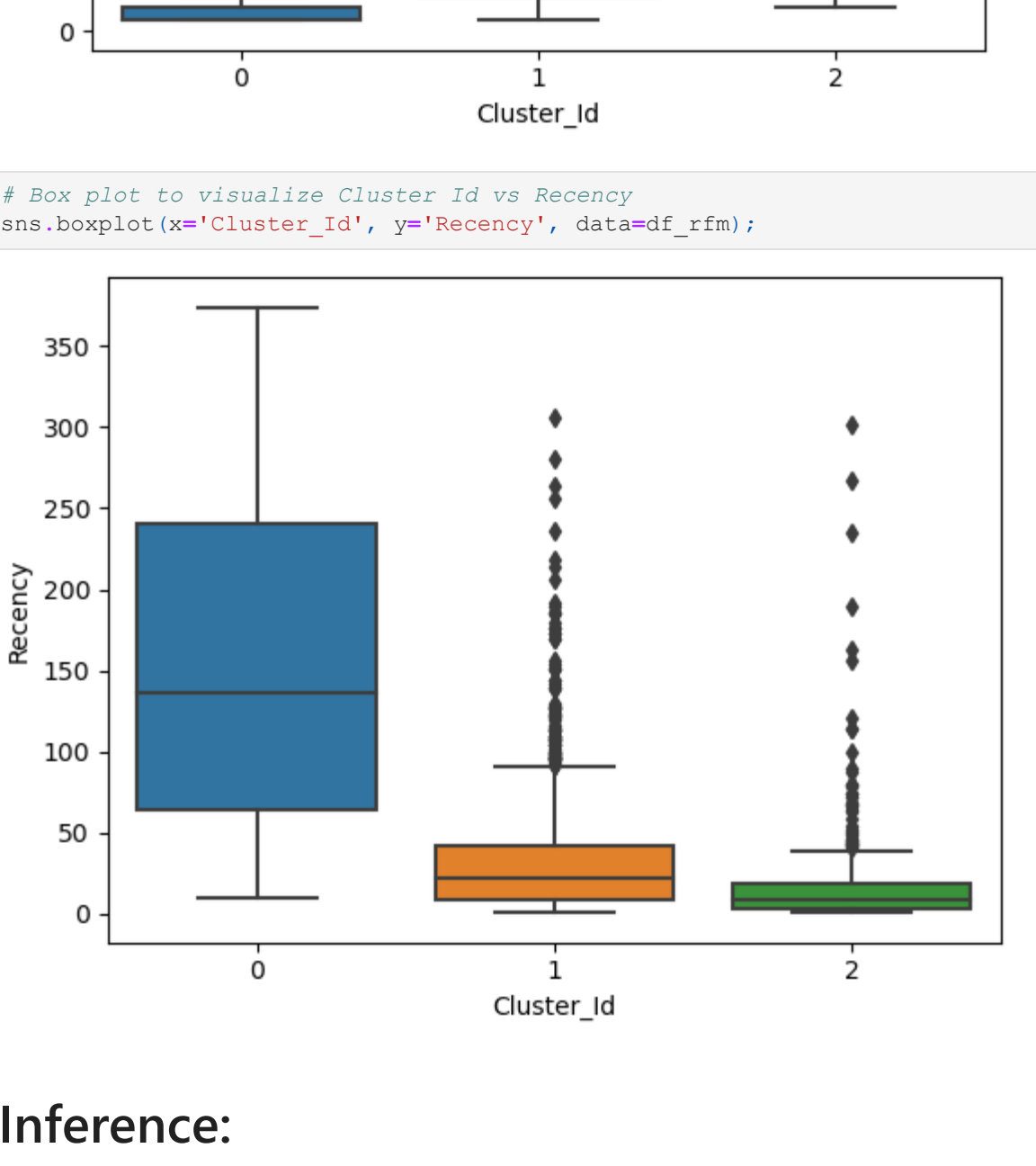
Out[94]: KMeans(max_iter=50, n_clusters=3)

c. Analyze these clusters and comment on the results.

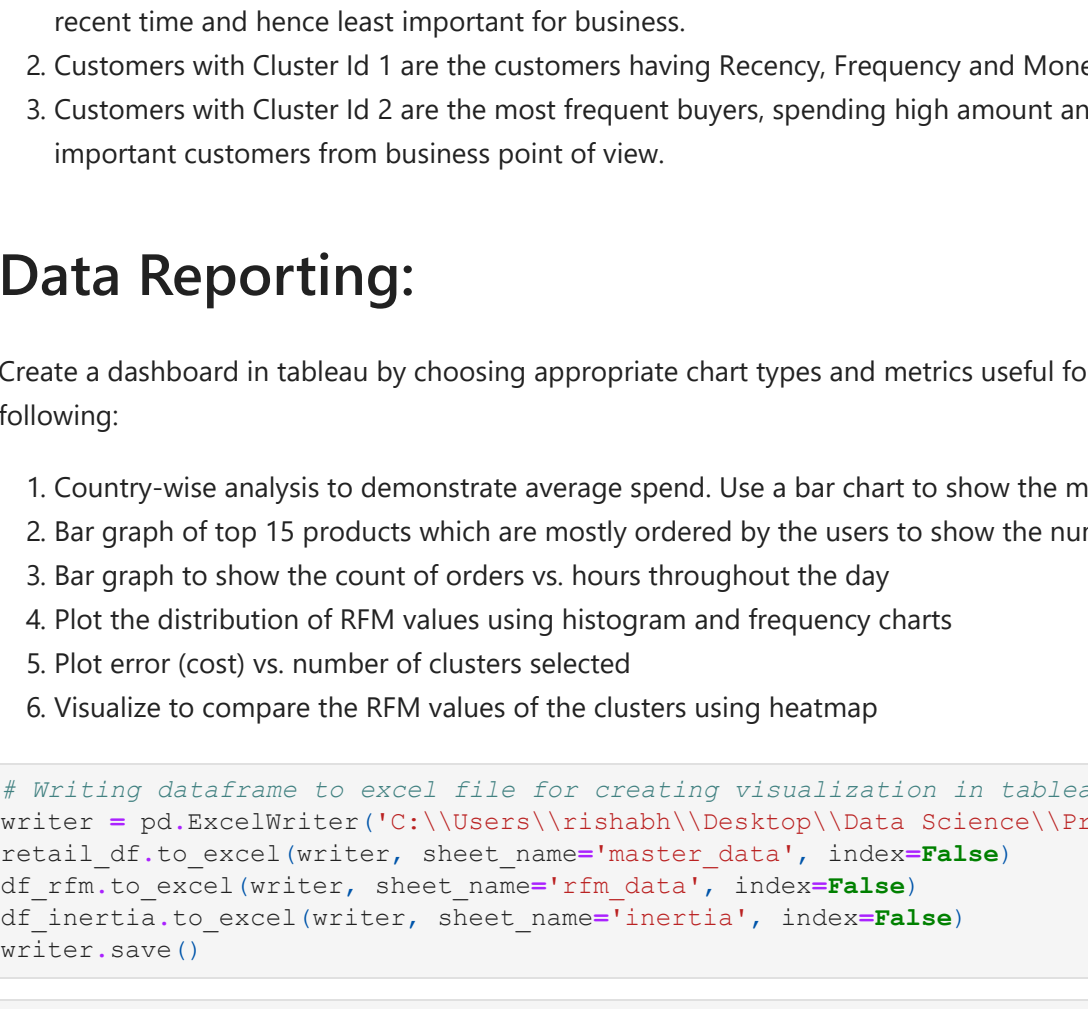
```
In [95]: # assign the label
df_rfm['Cluster_Id'] = kmeans.labels_
df_rfm.head()
```

	CustomerID	Recency	Frequency	Monetary	recency_labels	frequency_labels	monetary_labels	rfm_segment	rfm_score	Cluster_Id
0	12346.0	326	2	0.00	oldest	lowest	smallest	oldest-lowest-smallest	3	0
1	12347.0	2	7	4310.00	newest	lowest	smallest	newest-lowest-smallest	7	2
2	12348.0	75	4	1797.24	newest	lowest	smallest	newest-lowest-smallest	7	1
3	12349.0	19	1	1757.55	newest	lowest	smallest	newest-lowest-smallest	7	0
4	12350.0	310	1	334.40	oldest	lowest	smallest	oldest-lowest-smallest	3	0

```
In [96]: # Box plot to visualize Cluster_Id vs Monetary
sns.boxplot(x='Cluster_Id', y='Monetary', data=df_rfm);
```



```
In [97]: # Box plot to visualize Cluster_Id vs Frequency
sns.boxplot(x='Cluster_Id', y='Frequency', data=df_rfm);
```



```
In [98]: # Box plot to visualize Cluster_Id vs Recency
sns.boxplot(x='Cluster_Id', y='Recency', data=df_rfm);
```



Inference:

As we can observe from above boxplots that our model has nicely created 3 segments of customer with the interpretation as below:

- Customers with Cluster Id 0 are less frequent buyers with low monetary expenditure and also they have not purchased anything in recent time and hence least important for business.
- Customers with Cluster Id 1 are the customers having Recency, Frequency and Monetary score in the medium range.
- Customers with Cluster Id 2 are the most frequent buyers, spending high amount and recently placing orders so they are the most important customers from business point of view.

Data Reporting:

Create a dashboard in tableau by choosing appropriate chart types and metrics useful for the business. The dashboard must entail the following:

- Country-wise analysis to demonstrate average spend. Use a bar chart to show the monthly figures
- Bar graph of top 15 products which are mostly ordered by the users to show the number of products sold
- Bar graph to show the count of orders vs. hours throughout the day
- Plot the distribution of RFM values using histogram and frequency charts
- Plot error (cost) vs. number of clusters selected
- Visualize to compare the RFM values of the clusters using heatmap

```
In [185]: # Writing dataframe to excel file for creating visualization in tableau
writer = pd.ExcelWriter('C:\Users\irishabb\Desktop\Data Science\Project\Capstone Project\Online Retail1\RFM_data.xlsx')
df_rfm.to_excel(writer, sheet_name='RFM_data', index=False)
df_inertia.to_excel(writer, sheet_name='Inertia', index=False)
writer.save()
```

```
In [ ]:
```

```
In [ ]:
```