

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.preprocessing import StandardScaler
from sklearn.cluster import KMeans
from sklearn.metrics import silhouette_score
from datetime import timedelta
from pandas import ExcelWriter
```

```
In [2]: df = pd.read_excel("Online Retail.xlsx")
df.head()
```

Out[2]:

	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	17850.0	United Kingdom
1	536365	71053	WHITE METAL LANTERN	6	2010-12-01 08:26:00	3.39	17850.0	United Kingdom
2	536365	84406B	CREAM CUPID HEARTS COAT HANGER	8	2010-12-01 08:26:00	2.75	17850.0	United Kingdom
3	536365	84029G	KNITTED UNION FLAG HOT WATER BOTTLE	6	2010-12-01 08:26:00	3.39	17850.0	United Kingdom
4	536365	84029E	RED WOOLLY HOTTIE WHITE HEART.	6	2010-12-01 08:26:00	3.39	17850.0	United Kingdom

```
In [3]: df.shape
```

Out[3]: (541909, 8)

```
In [4]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 541909 entries, 0 to 541908
Data columns (total 8 columns):
#   Column          Non-Null Count  Dtype
---  -
0   InvoiceNo        541909 non-null object
1   StockCode       541909 non-null object
2   Description      540455 non-null object
3   Quantity        541909 non-null int64
4   InvoiceDate      541909 non-null datetime64[ns]
5   UnitPrice       541909 non-null float64
6   CustomerID      406829 non-null float64
7   Country         541909 non-null object
dtypes: datetime64[ns](1), float64(2), int64(1), object(4)
memory usage: 33.1+ MB
```

```
In [5]: # Check missing values in dataset
df.isnull().sum()
```

```
Out[5]: InvoiceNo          0
StockCode          0
Description      1454
Quantity          0
InvoiceDate        0
UnitPrice          0
CustomerID      135080
Country           0
dtype: int64
```

```
In [6]: # Calculating the Missing Values % contribution in DF
df_null = round(df.isnull().sum()/len(df)*100,2)
df_null
```

```
Out[6]: InvoiceNo      0.00
        StockCode     0.00
        Description   0.27
        Quantity      0.00
        InvoiceDate    0.00
        UnitPrice      0.00
        CustomerID    24.93
        Country       0.00
        dtype: float64
```

```
In [7]: invoice_null_custid = set(df[df['CustomerID'].isnull()][ 'InvoiceNo'])
df[df['InvoiceNo'].isin(invoice_null_custid) & (~df['CustomerID'].isnull())]
```

```
Out[7]:
```

InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country
-----------	-----------	-------------	----------	-------------	-----------	------------	---------

```
In [8]: df = df.drop('Description', axis=1)
df = df.dropna()
df.shape
```

```
Out[8]: (406829, 7)
```

```
In [9]: df = df.drop_duplicates()
df.shape
```

```
Out[9]: (401602, 7)
```

```
In [10]: df['CustomerID'] = df['CustomerID'].astype(str)
```

```
In [11]: df.describe(include=['O'])
```

```
Out[11]:
```

	InvoiceNo	StockCode	CustomerID	Country
count	401602	401602	401602	401602
unique	22190	3684	4372	37
top	576339	85123A	17841.0	United Kingdom
freq	542	2065	7812	356726

## # Create month cohort of customers and analyze active customers in each cohort:

```
In [12]: # Convert to InvoiceDate to Year-Month format
df['month_year'] = df['InvoiceDate'].dt.to_period('M')
df['month_year'].nunique()
```

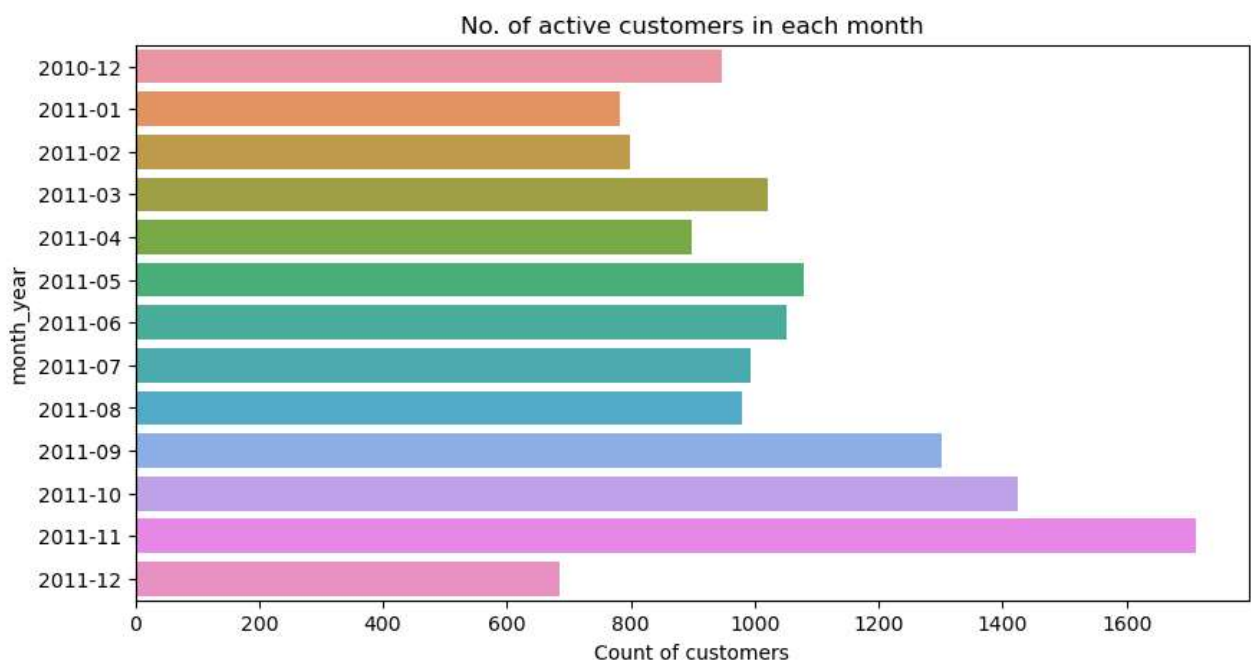
```
Out[12]: 13
```

```
In [13]: month_cohort = df.groupby('month_year')['CustomerID'].nunique()  
month_cohort
```

```
Out[13]: 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 [14]: 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[14]: Text(0.5, 1.0, 'No. of active customers in each month')
```



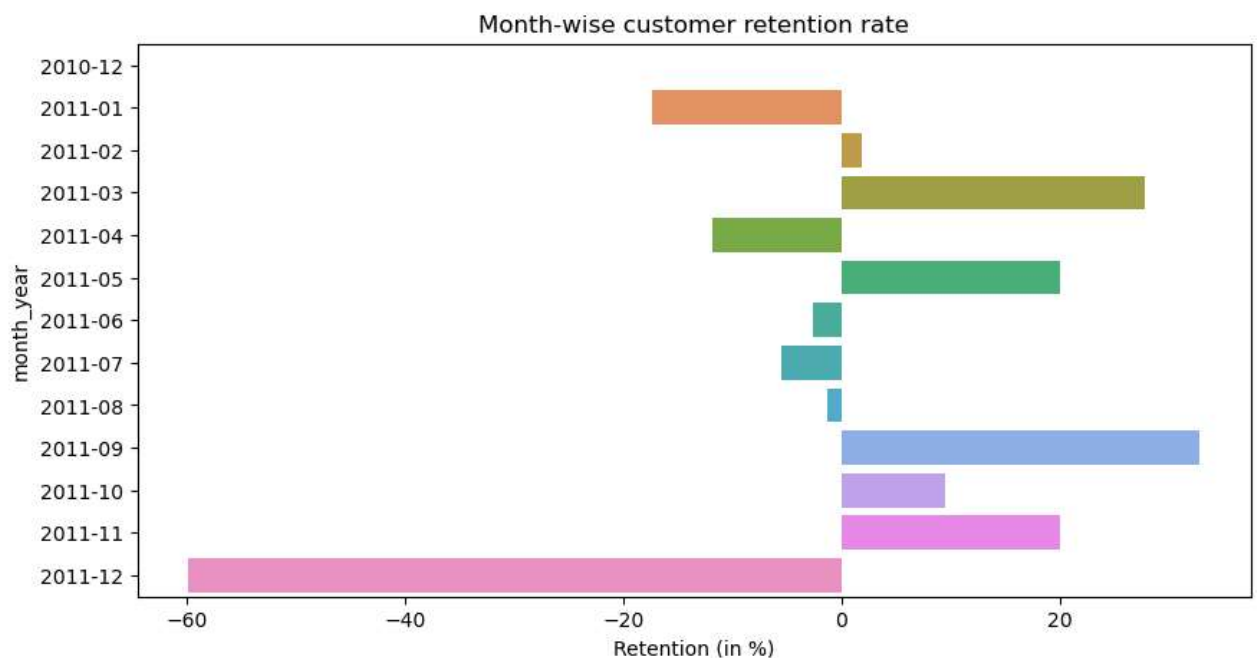
```
In [15]: month_cohort - month_cohort.shift(1)
```

```
Out[15]: 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   -1025.0
Freq: M, Name: CustomerID, dtype: float64
```

```
In [16]: retention_rate = round(month_cohort.pct_change( periods=1)*100,2)
         retention_rate
```

```
Out[16]: 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 [17]: 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");
```



```
In [18]: df['amount'] = df['Quantity']*df['UnitPrice']
df.head()
```

Out[18]:

	InvoiceNo	StockCode	Quantity	InvoiceDate	UnitPrice	CustomerID	Country	month_year	amount
0	536365	85123A	6	2010-12-01 08:26:00	2.55	17850.0	United Kingdom	2010-12	15.30
1	536365	71053	6	2010-12-01 08:26:00	3.39	17850.0	United Kingdom	2010-12	20.34
2	536365	84406B	8	2010-12-01 08:26:00	2.75	17850.0	United Kingdom	2010-12	22.00
3	536365	84029G	6	2010-12-01 08:26:00	3.39	17850.0	United Kingdom	2010-12	20.34
4	536365	84029E	6	2010-12-01 08:26:00	3.39	17850.0	United Kingdom	2010-12	20.34

```
In [19]: df_monetary = df.groupby('CustomerID').sum()['amount'].reset_index()
df_monetary
```

Out[19]:

	CustomerID	amount
0	12346.0	0.00
1	12347.0	4310.00
2	12348.0	1797.24
3	12349.0	1757.55
4	12350.0	334.40
...	...	...
4367	18280.0	180.60
4368	18281.0	80.82
4369	18282.0	176.60
4370	18283.0	2045.53
4371	18287.0	1837.28

4372 rows × 2 columns

```
In [20]: df_frequency = df.groupby('CustomerID').nunique()['InvoiceNo'].reset_index()
df_frequency
```

Out[20]:

	CustomerID	InvoiceNo
0	12346.0	2
1	12347.0	7
2	12348.0	4
3	12349.0	1
4	12350.0	1
...	...	...
4367	18280.0	1
4368	18281.0	1
4369	18282.0	3
4370	18283.0	16
4371	18287.0	3

4372 rows × 2 columns

```
In [21]: # calculating recently as last transaction day in data + 1 day
ref_day = max(df['InvoiceDate']) + timedelta(days=1)
df['days_to_last_order'] = (ref_day - df['InvoiceDate']).dt.days
df.head()
```

Out[21]:

	InvoiceNo	StockCode	Quantity	InvoiceDate	UnitPrice	CustomerID	Country	month_year	amount	days_to_last_order
0	536365	85123A	6	2010-12-01 08:26:00	2.55	17850.0	United Kingdom	2010-12	15.30	326
1	536365	71053	6	2010-12-01 08:26:00	3.39	17850.0	United Kingdom	2010-12	20.34	2
2	536365	84406B	8	2010-12-01 08:26:00	2.75	17850.0	United Kingdom	2010-12	22.00	75
3	536365	84029G	6	2010-12-01 08:26:00	3.39	17850.0	United Kingdom	2010-12	20.34	19
4	536365	84029E	6	2010-12-01 08:26:00	3.39	17850.0	United Kingdom	2010-12	20.34	43

```
In [22]: df_recency = df.groupby('CustomerID')['days_to_last_order'].min().reset_index()
df_recency
```

Out[22]:

	CustomerID	days_to_last_order
0	12346.0	326
1	12347.0	2
2	12348.0	75
3	12349.0	19
4	12350.0	310
...	...	...
4367	18280.0	278
4368	18281.0	181
4369	18282.0	8
4370	18283.0	4
4371	18287.0	43

4372 rows × 2 columns

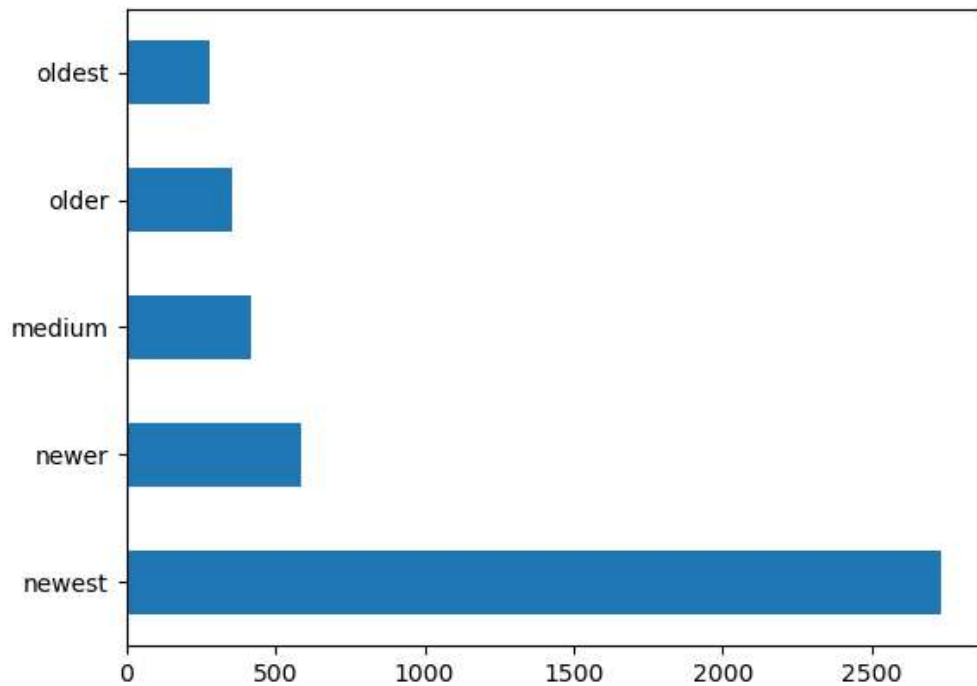
```
In [23]: df_rf = pd.merge(df_recency, df_frequency, on='CustomerID', how='inner')
df_rfm = pd.merge(df_rf, df_monetary, on='CustomerID', how='inner')
df_rfm.columns = ['CustomerID', 'Recency', 'Frequency', 'Monetary']
df_rfm.head()
```

Out[23]:

	CustomerID	Recency	Frequency	Monetary
0	12346.0	326	2	0.00
1	12347.0	2	7	4310.00
2	12348.0	75	4	1797.24
3	12349.0	19	1	1757.55
4	12350.0	310	1	334.40

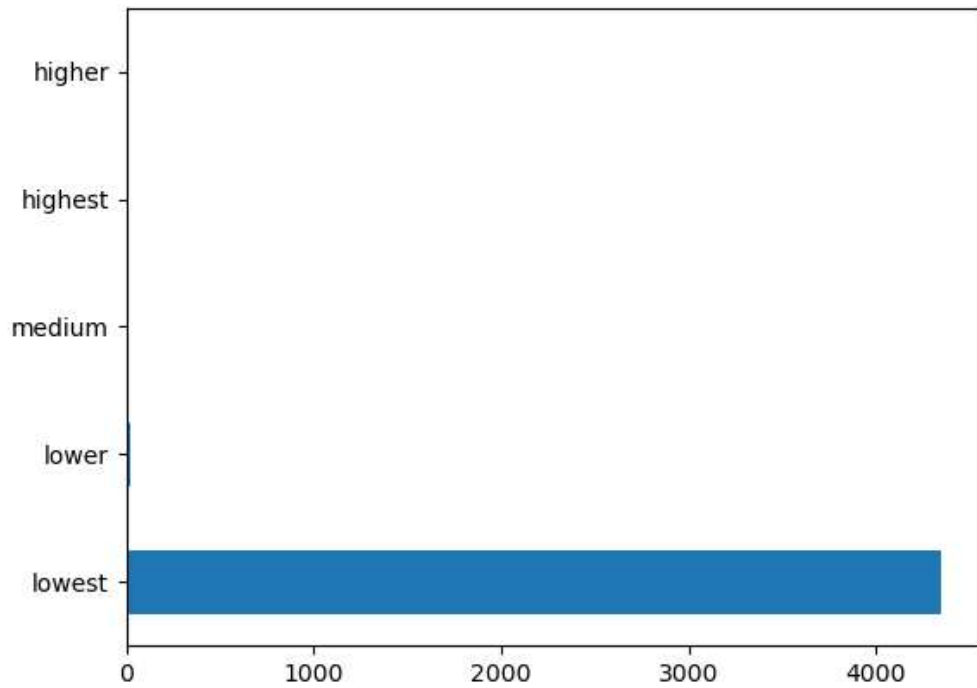
```
In [25]: df_rfm['recency_labels'] = pd.cut(df_rfm['Recency'], bins=5, labels=['newest', 'newer', 'medium',  
df_rfm['recency_labels'].value_counts().plot(kind='barh');  
df_rfm['recency_labels'].value_counts())
```

```
Out[25]: newest      2734  
newer      588  
medium     416  
older      353  
oldest     281  
Name: recency_labels, dtype: int64
```



```
In [26]: df_rfm['frequency_labels'] = pd.cut(df_rfm['Frequency'], bins=5, labels=['lowest', 'lower', 'medium', 'highest', 'higher'],  
df_rfm['frequency_labels'].value_counts().plot(kind='barh');  
df_rfm['frequency_labels'].value_counts()
```

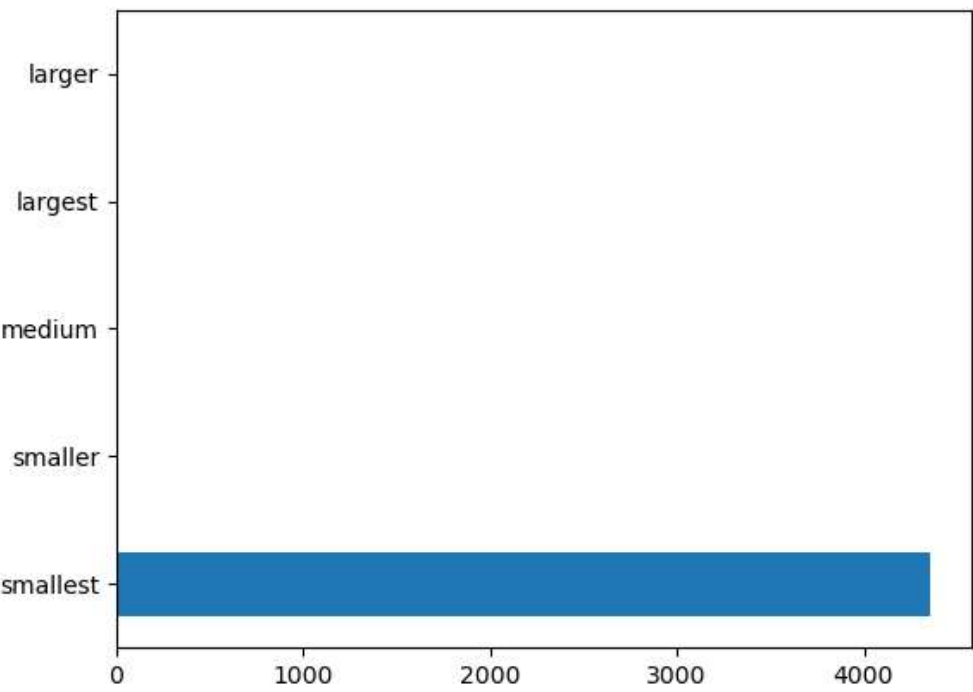
```
Out[26]: lowest      4348  
lower         18  
medium         3  
highest        2  
higher         1  
Name: frequency_labels, dtype: int64
```





```
In [27]: df_rfm['monetary_labels'] = pd.cut(df_rfm['Monetary'], bins=5, labels=['smallest', 'smaller', 'medium', 'largest', 'larger'],
df_rfm['monetary_labels'].value_counts().plot(kind='barh');
df_rfm['monetary_labels'].value_counts()
```

Out[27]: smallest 4357  
smaller 9  
medium 3  
largest 2  
larger 1  
Name: monetary\_labels, dtype: int64



```
In [28]: df_rfm['rfm_segment'] = df_rfm[['recency_labels', 'frequency_labels', 'monetary_labels']].agg('-'.join, axis=1)
df_rfm.head()
```

Out[28]:

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

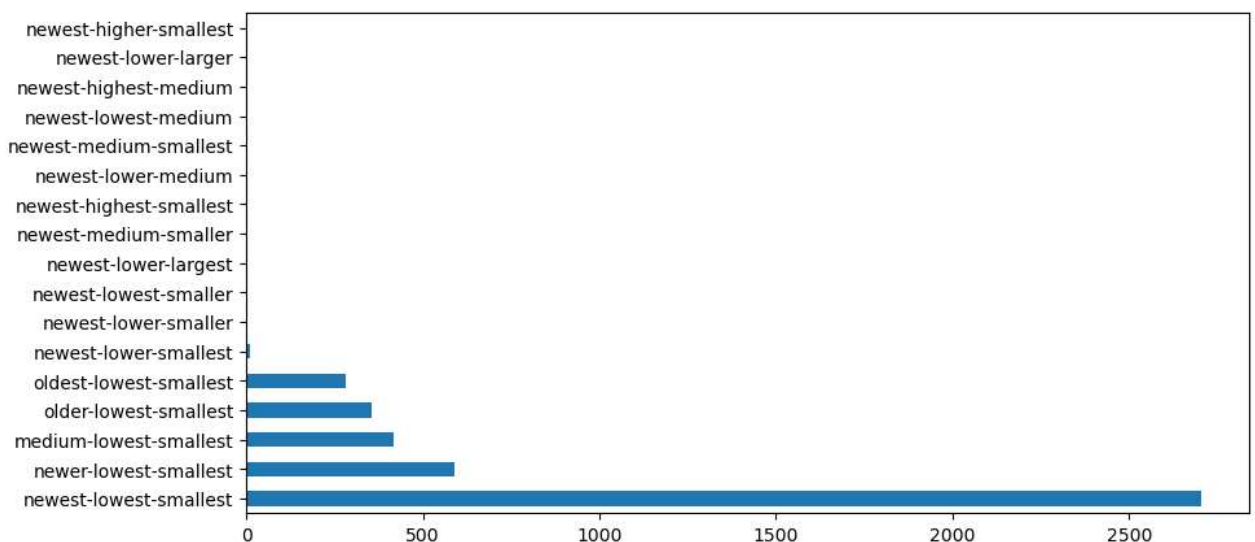
```
In [29]: recency_dict = {'newest': 5, 'newer':4, 'medium': 3, 'older':2, 'oldest':1}
frequency_dict = {'lowest':1, 'lower':2, 'medium': 3, 'higher':4, 'highest':5}
monetary_dict = {'smallest':1, 'smaller':2, 'medium': 3, 'larger':4, 'largest':5}

df_rfm['rfm_score'] = df_rfm['recency_labels'].map(recency_dict).astype(int)+ df_rfm['frequency_
df_rfm.head(10)
```

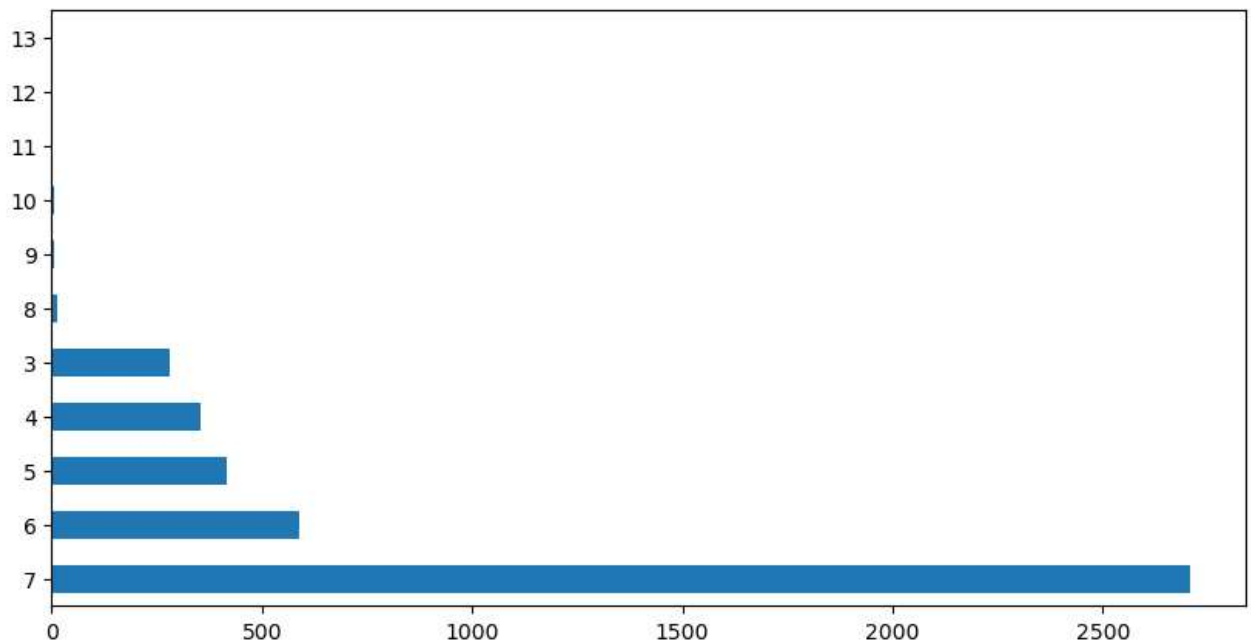
Out[29]:

	CustomerID	Recency	Frequency	Monetary	recency_labels	frequency_labels	monetary_labels	rfm_segment	rfm_sc
0	12346.0	326	2	0.00	oldest	lowest	smallest	oldest-lowest-smallest	
1	12347.0	2	7	4310.00	newest	lowest	smallest	newest-lowest-smallest	
2	12348.0	75	4	1797.24	newest	lowest	smallest	newest-lowest-smallest	
3	12349.0	19	1	1757.55	newest	lowest	smallest	newest-lowest-smallest	
4	12350.0	310	1	334.40	oldest	lowest	smallest	oldest-lowest-smallest	
5	12352.0	36	11	1545.41	newest	lowest	smallest	newest-lowest-smallest	
6	12353.0	204	1	89.00	medium	lowest	smallest	medium-lowest-smallest	
7	12354.0	232	1	1079.40	older	lowest	smallest	older-lowest-smallest	
8	12355.0	214	1	459.40	medium	lowest	smallest	medium-lowest-smallest	
9	12356.0	23	3	2811.43	newest	lowest	smallest	newest-lowest-smallest	

```
In [30]: df_rfm['rfm_segment'].value_counts().plot(kind='barh', figsize=(10, 5));
```



```
In [31]: df_rfm['rfm_score'].value_counts().plot(kind='barh', figsize=(10, 5));
```



```
In [32]: print(df_rfm.shape)
df_rfm.head()
```

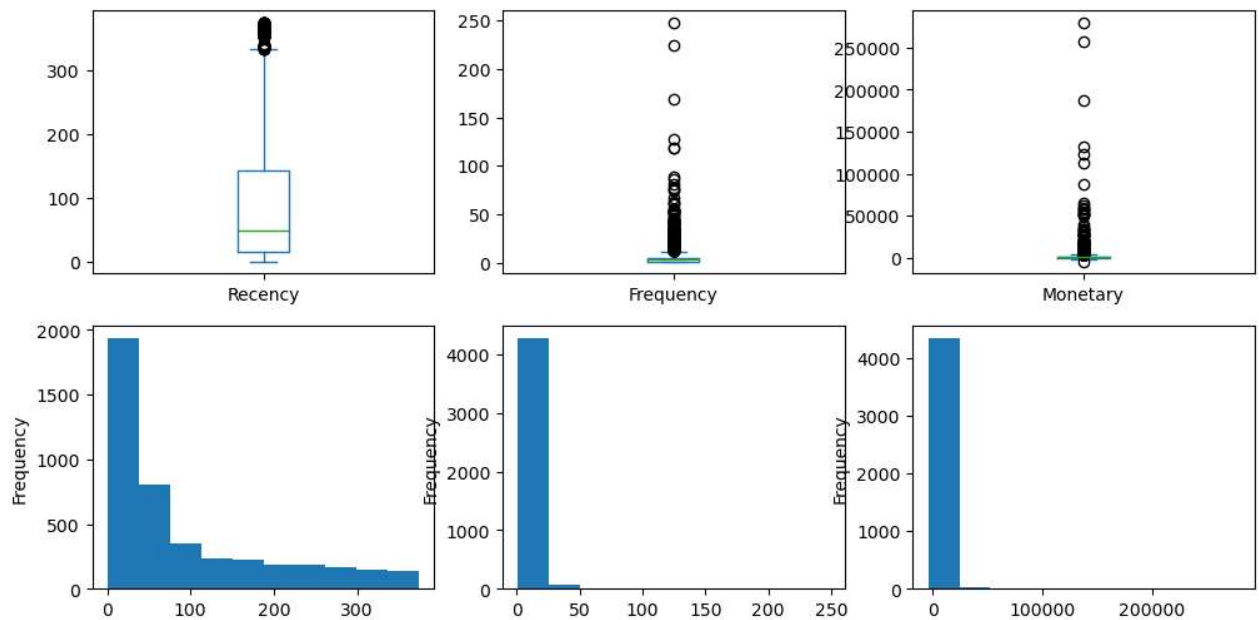
(4372, 9)

Out[32]:

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

```
In [33]: plt.figure(figsize=(12,6))

for i, feature in enumerate(['Recency', 'Frequency', 'Monetary']):
    plt.subplot(2,3,i+1)
    df_rfm[feature].plot(kind='box')
    plt.subplot(2,3,i+1+3)
    df_rfm[feature].plot(kind='hist')
```

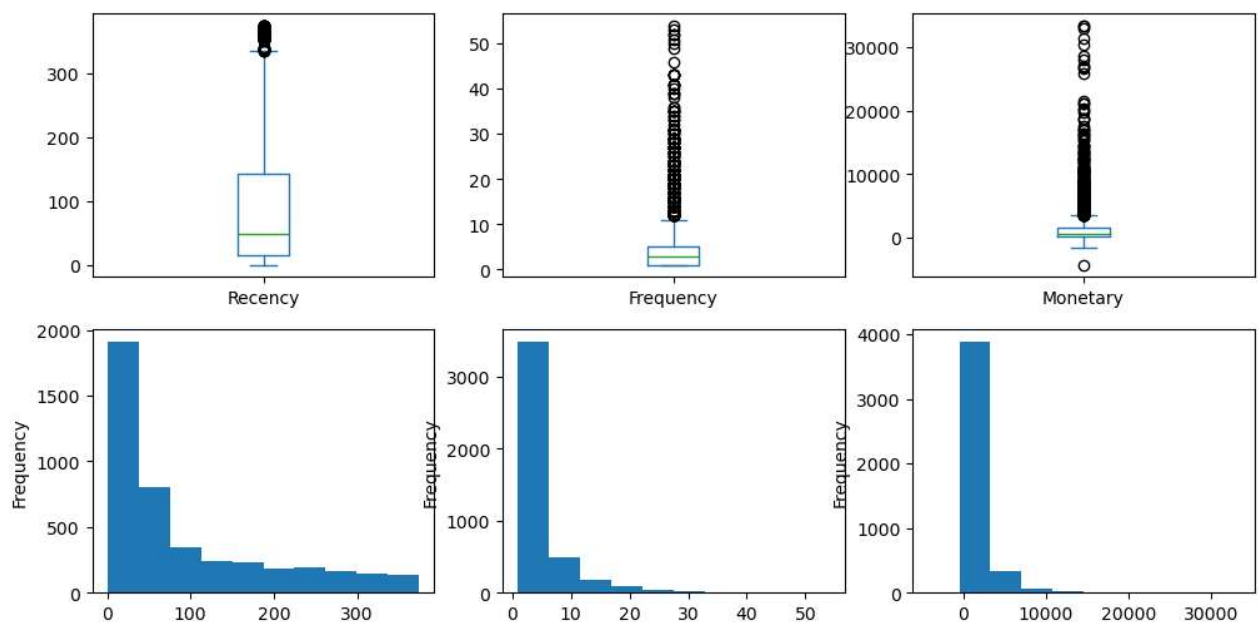


```
In [34]: df_rfm = df_rfm[(df_rfm['Frequency'] < 60) & (df_rfm['Monetary'] < 40000)]
df_rfm.shape
```

Out[34]: (4346, 9)

```
In [35]: plt.figure(figsize=(12,6))

for i, feature in enumerate(['Recency', 'Frequency', 'Monetary']):
    plt.subplot(2,3,i+1)
    df_rfm[feature].plot(kind='box')
    plt.subplot(2,3,i+1+3)
    df_rfm[feature].plot(kind='hist')
```



In [ ]: *#Log Transformation:*

```
In [36]: df_rfm_log_trans = pd.DataFrame()
df_rfm_log_trans['Recency'] = np.log(df_rfm['Recency'])
df_rfm_log_trans['Frequency'] = np.log(df_rfm['Frequency'])
df_rfm_log_trans['Monetary'] = np.log(df_rfm['Monetary']-df_rfm['Monetary'].min()+1)
```

In [ ]: *#Standard Scalar Transformation:*

```
In [37]: scaler = StandardScaler()

df_rfm_scaled = scaler.fit_transform(df_rfm_log_trans[['Recency', 'Frequency', 'Monetary']])
df_rfm_scaled

df_rfm_scaled = pd.DataFrame(df_rfm_scaled)
df_rfm_scaled.columns = ['Recency', 'Frequency', 'Monetary']
```

In [38]: df\_rfm\_scaled.head()

Out[38]:

	Recency	Frequency	Monetary
0	1.402988	-0.388507	-0.770922
1	-2.100874	0.967301	1.485132
2	0.392218	0.361655	0.364190
3	-0.552268	-1.138669	0.342970
4	1.368370	-1.138669	-0.527416

```
In [39]: # k-means with some arbitrary k
kmeans = KMeans(n_clusters=3, max_iter=50)
kmeans.fit(df_rfm_scaled)
```

Out[39]: KMeans(max\_iter=50, n\_clusters=3)

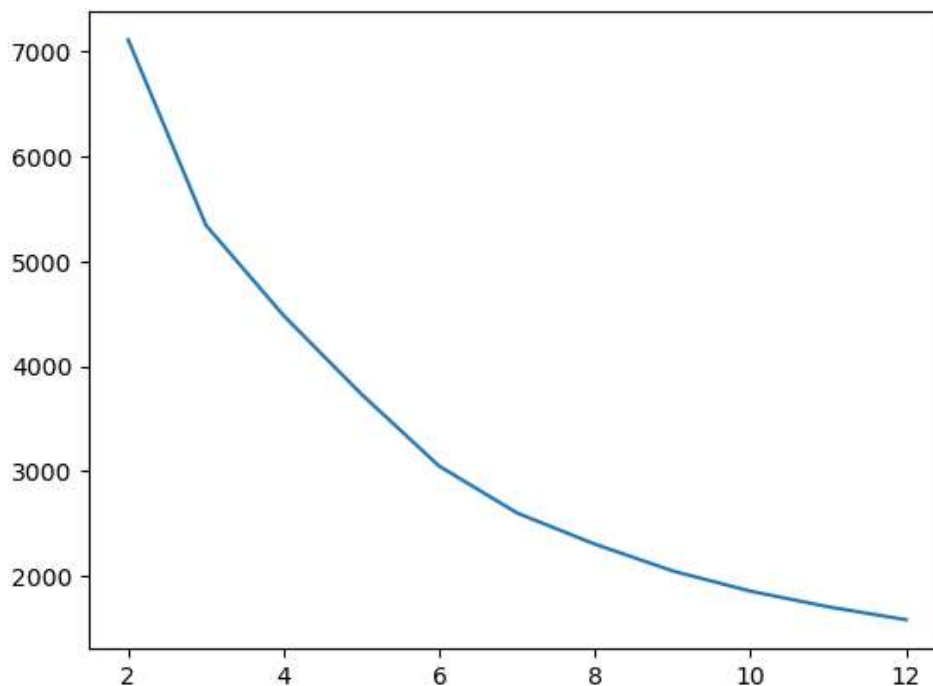
In [40]: kmeans.labels\_

Out[40]: array([1, 2, 0, ..., 0, 2, 0])

```
In [41]: ssd = []
range_n_clusters = [2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]
for num_clusters in range_n_clusters:
    kmeans = KMeans(n_clusters=num_clusters, max_iter=100)
    kmeans.fit(df_rfm_scaled)

    ssd.append(kmeans.inertia_)

# plot the SSDs for each n_clusters
plt.plot(range_n_clusters, ssd);
```



```
In [42]: # Creating dataframe for exporting to create visualization in tableau later
df_inertia = pd.DataFrame(list(zip(range_n_clusters, ssd)), columns=['clusters', 'intertia'])
df_inertia
```

Out[42]:

	clusters	intertia
0	2	7113.109513
1	3	5343.136928
2	4	4480.966244
3	5	3730.690486
4	6	3044.999774
5	7	2598.665907
6	8	2301.542620
7	9	2046.156974
8	10	1853.011464
9	11	1704.000984
10	12	1580.762814

```
In [43]: # 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 0.3777355515144021
For n_clusters=4, the silhouette score is 0.3624139517202161
For n_clusters=5, the silhouette score is 0.36419107385871785
For n_clusters=6, the silhouette score is 0.3447566670598209
For n_clusters=7, the silhouette score is 0.34295615020953824
For n_clusters=8, the silhouette score is 0.3395626912561983
For n_clusters=9, the silhouette score is 0.3463314430741466
For n_clusters=10, the silhouette score is 0.35612016121427464
```

```
In [44]: # Final model with k=3
kmeans = KMeans(n_clusters=3, max_iter=50)
kmeans.fit(df_rfm_scaled)
```

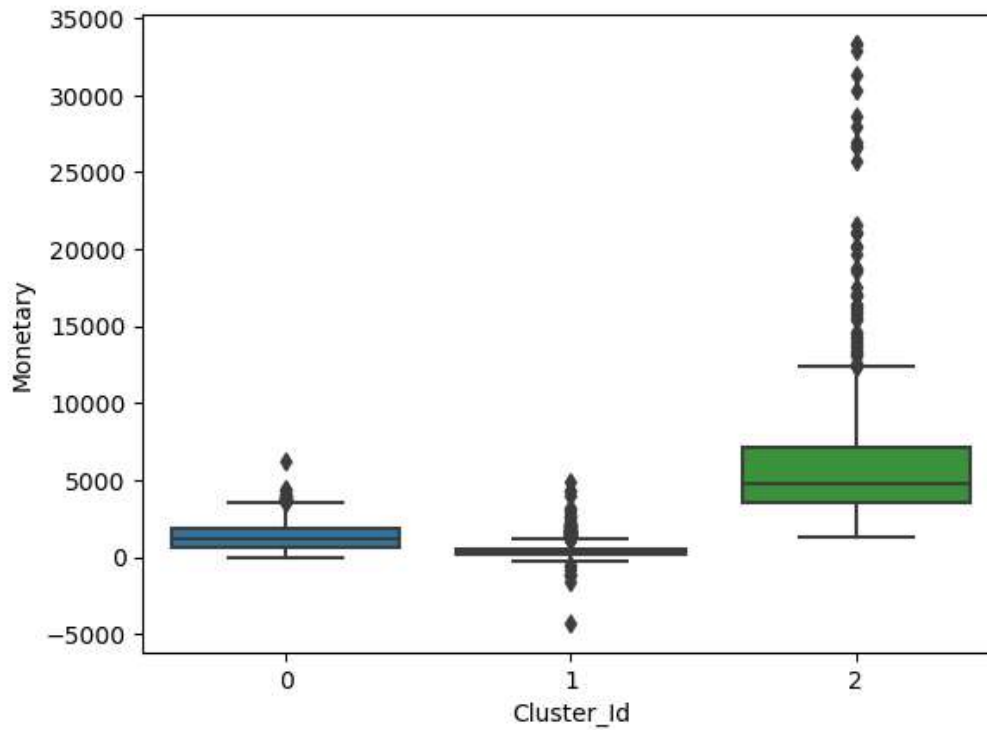
```
Out[44]: KMeans(max_iter=50, n_clusters=3)
```

```
In [45]: # assign the Label
df_rfm['Cluster_Id'] = kmeans.labels_
df_rfm.head()
```

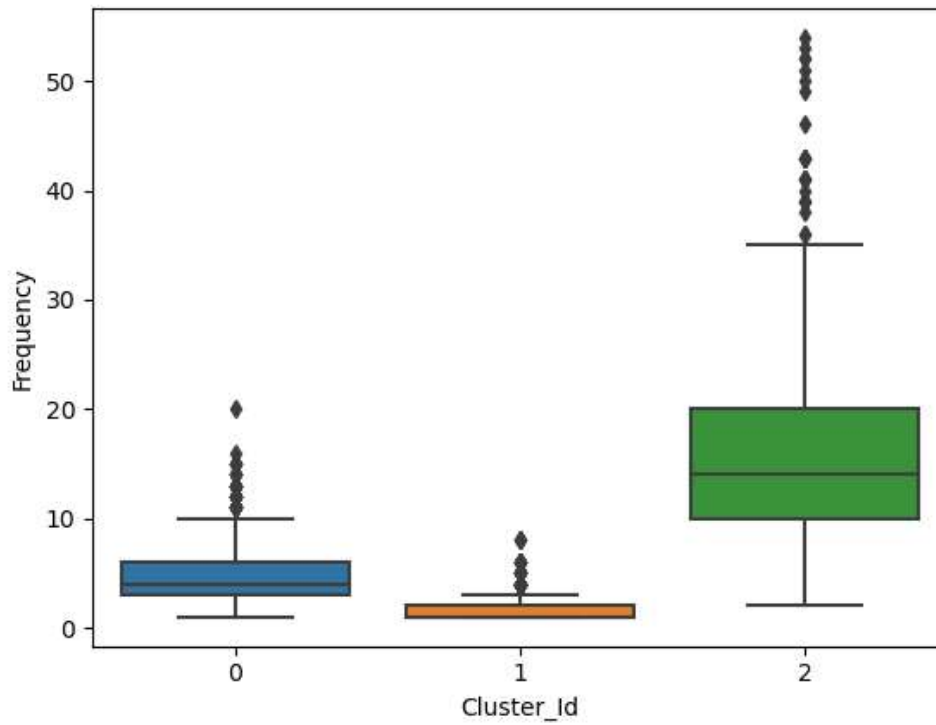
```
Out[45]:
```

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

```
In [46]: # Box plot to visualize Cluster Id vs Monetary
sns.boxplot(x='Cluster_Id', y='Monetary', data=df_rfm);
```

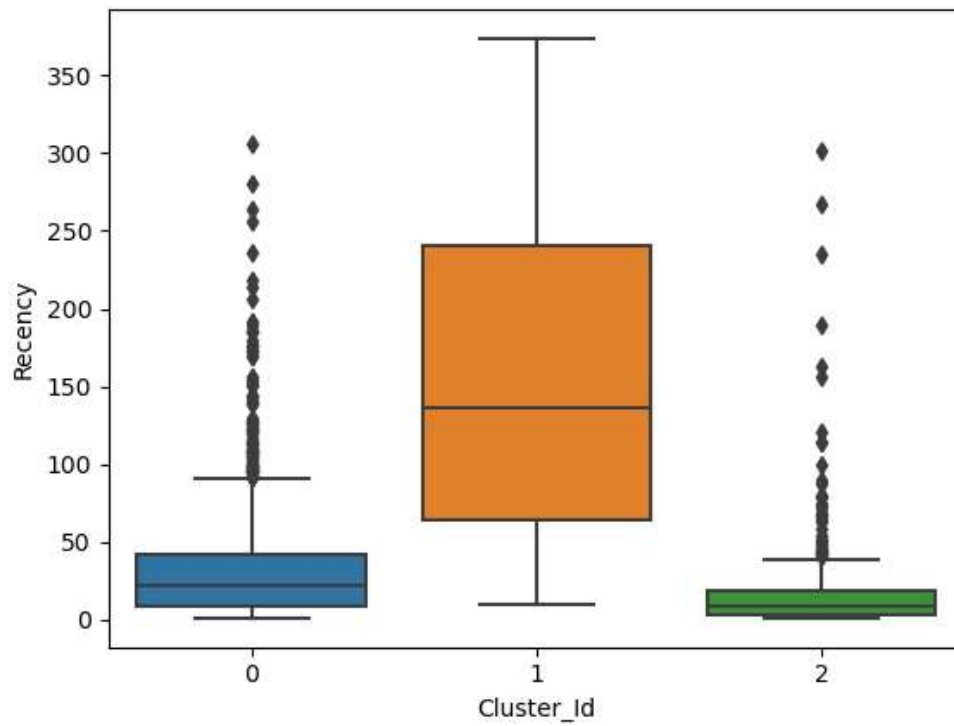


```
In [47]: # Box plot to visualize Cluster Id vs Frequency
sns.boxplot(x='Cluster_Id', y='Frequency', data=df_rfm);
```





```
In [48]: # Box plot to visualize Cluster Id vs Recency  
sns.boxplot(x='Cluster_Id', y='Recency', data=df_rfm);
```



```
In [ ]:
```