

## OUTPUT: Project 3: Retail

The JupyterLab interface displays the 'Project 3: Retail' notebook. The left sidebar shows a file explorer with a list of files: 'DS / Project 1\_Retail /', 'C:\Users\m...', 'Online Reta...', 'product\_de...', 'Project 1\_R...', and 'Project 3\_R...' (selected). The main area shows the notebook content for 'Project Task: Week 1'.

**Project Task: Week 1**

**Data Cleanings:**

1. Perform a preliminary data inspection and data cleaning.
- a. Check for missing data and formulate an apt strategy to treat them.
- b. Remove duplicate data records.
- c. Perform descriptive analytics on the given data.

**Data Transformation:**

2. Perform cohort analysis (a cohort is a group of subjects that share a defining characteristic). Observe how a cohort behaves across time and compare it to other cohorts.
- a. Create month cohorts and analyze active customers for each cohort.
- b. Analyze the retention rate of customers.

```
[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
```

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

```
[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	844068	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

The JupyterLab interface displays the 'Project 3: Retail' notebook. The left sidebar shows a file explorer with a list of files: 'DS / Project 1\_Retail /', 'C:\Users\m...', 'Online Reta...', 'product\_de...', 'Project 1\_R...', and 'Project 3\_R...' (selected). The main area shows the notebook content for 'Project Task: Week 1'.

**Project Task: Week 1**

**Data Cleanings:**

1. Perform a preliminary data inspection and data cleaning.
- a. Check for missing data and formulate an apt strategy to treat them.
- b. Remove duplicate data records.
- c. Perform descriptive analytics on the given data.

**Data Transformation:**

2. Perform cohort analysis (a cohort is a group of subjects that share a defining characteristic). Observe how a cohort behaves across time and compare it to other cohorts.
- a. Create month cohorts and analyze active customers for each cohort.
- b. Analyze the retention rate of customers.

```
[1]: # Missing values treatment:
```

```
[10]: # Check shape of data
df.shape
```

```
[10]: (486829, 7)
```

```
[4]: # Check feature details of data
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   486829 non-null  float64  
 7   Country      541909 non-null  object  
dtypes: datetime64[ns](1), float64(2), int64(1), object(4)
memory usage: 33.1+ MB
```

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

```
[5]: InvoiceNo      0
StockCode      0
Description    1454
Quantity       0
InvoiceDate    0
UnitPrice      0
CustomerID    135880
Country        0
dtype: int64
```

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

jupyter

Used 0 of 50 hours in Nov, 2022

Start Lab

End Lab

File

Edit

View

Run

Kernel

Tools

Settings

Help

Project 3\_Retail.ipynb

Country  
dtype: int64

0

Python 3

Filter files by name

/ DS / Project 1\_Retail /

Name	Last Modified
C:\Users\m...	29 minutes ago
Online Reta...	2 hours ago
product_de...	28 minutes ago
Project 1_R...	2 hours ago
Project 3_R...	29 minutes ago

[6]:

# Calculating the Missing Values % contribution in DF  
df\_null = round(df.isnull().sum()/len(df)\*100,2)  
df\_null

[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

As we can see two columns in data have missing values.  
  
Description - 0.27% (1454 nos.)  
  
CustomerID - 24.93% (135080)

CustomerID is important feature of our analysis since our analysis is centered around Customers only so we can not impute null values CustomerID with mean/ median/ mode in this case. We will check possibility to fill null values in CustomerID column by looking up for InvoiceNo of the row having null CustomerID in other rows where CustomerID is present. If there are still any null values in CustomerID after this process then we will drop complete row having missing CustomerID.

We can drop Description feature from our data since it is not going to contribute in our model.

[8]:

invoice\_null\_custid = set(df[df['CustomerID'].isnull()]['InvoiceNo'])  
df[df['InvoiceNo'].isin(invoice\_null\_custid) & (~df['CustomerID'].isnull())]

[8]:

InvoiceNo StockCode Description Quantity InvoiceDate UnitPrice CustomerID Country

We could not find any value to impute null values in CustomerID column since all entries for a particular InvoiceNo have missing CustomerID if that particular InvoiceNo has null CustomerID in even one entry. So we will drop all rows having null values in CustomerID.

[9]:

df = df.drop("Description", axis=1)  
df = df.dropna()  
df.shape

[9]:

(486829, 7)

Simple

0

3

No Kernel

idle

Mode: Command

Ln 1, Col 19

Project 3\_Retail.ipynb

Project 3\_Retail.ipynb

```

df.shape
[9]: (406829, 7)

[ ]: #Remove duplicate data records:

[11]: df = df.drop_duplicates()
df.shape
[11]: (401602, 7)

[ ]: #Perform descriptive analysis on the given data:

[12]: # CustomerID is 'float64', changing the datatype of CustomerID to string as Customer ID as numerical data does not make sense
df['CustomerID'] = df['CustomerID'].astype(str)

[13]: df.describe(datetime_is_numeric=True)

```

	Quantity	InvoiceDate	UnitPrice
count	401602.000000	401602	401602.000000
mean	12.182579	2011-07-10 12:08:08.129743104	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

Quantity: Average quantity of each product in transaction is 12.18. Also note that minimum value in Quantity column is negative. This implies that some customers had returned the product during our period of analysis. InvoiceDate: Our data has transaction between 01-12-2010 to 09-12-2011 UnitPrice: Average price of each product in transactions is 3.47

```

[14]: df.describe(include=['O'])

[14]:
InvoiceNo  StockCode  CustomerID  Country
count      401602      401602      401602      401602

```

Simple 0 3 No Kernel | Idle Mode: Command Ln 1, Col 19 Project 3\_Retail.ipynb

Project 3\_Retail.ipynb

```

[14]: df.describe(include=['O'])

[14]:
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

```

InvoiceNo: Total entries in preprocessed data are 4,01,602 but transactions are 22,190. Most number of entries (count of unique products) are in Invoice No. '576339' and is 542 nos. StockCode: There are total 3684 unique products in our data and product with stock code '85123A' appears most frequently (2065 times) in our data. CustomerID: There are 4372 unique customers in our final preprocessed data. Customer with ID '17841' appears most frequently in data (7812 times) Country: Company has customers across 37 countries. Most entries are from United Kingdom in our dataset (356726)

```

[ ]: #Perform Cohort Analysis
#Create month cohort of customers and analyze active customers in each cohort:

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

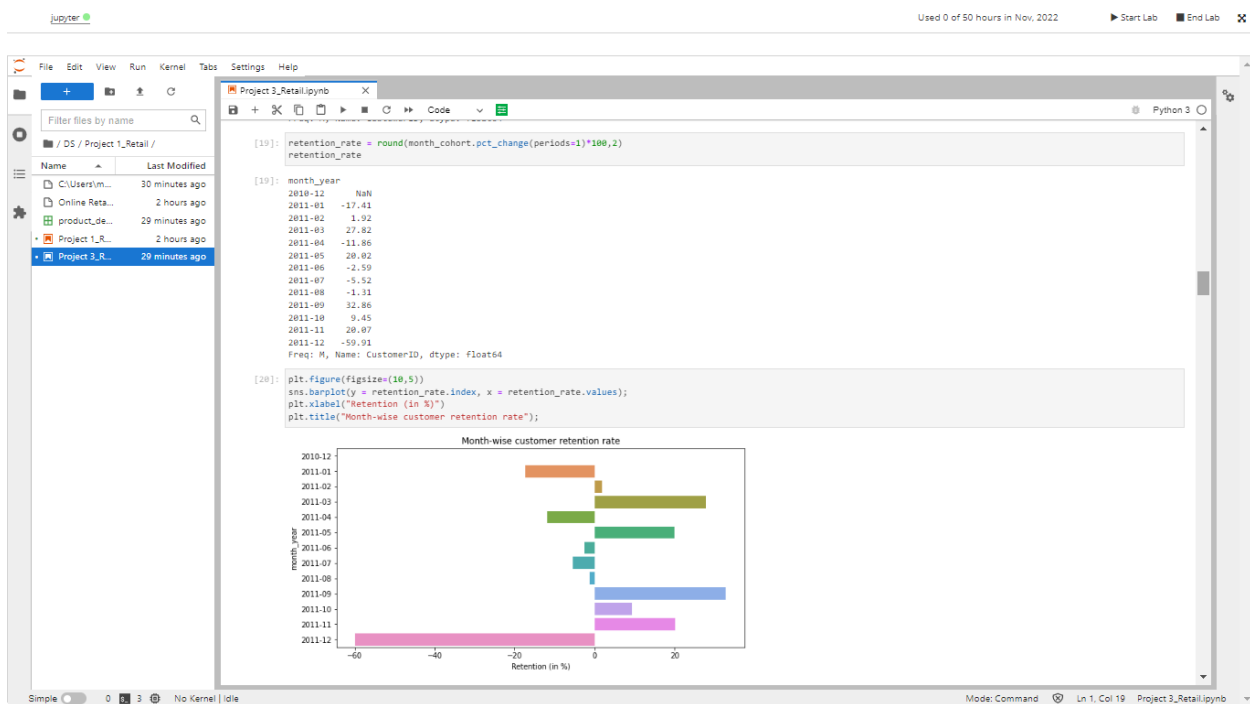
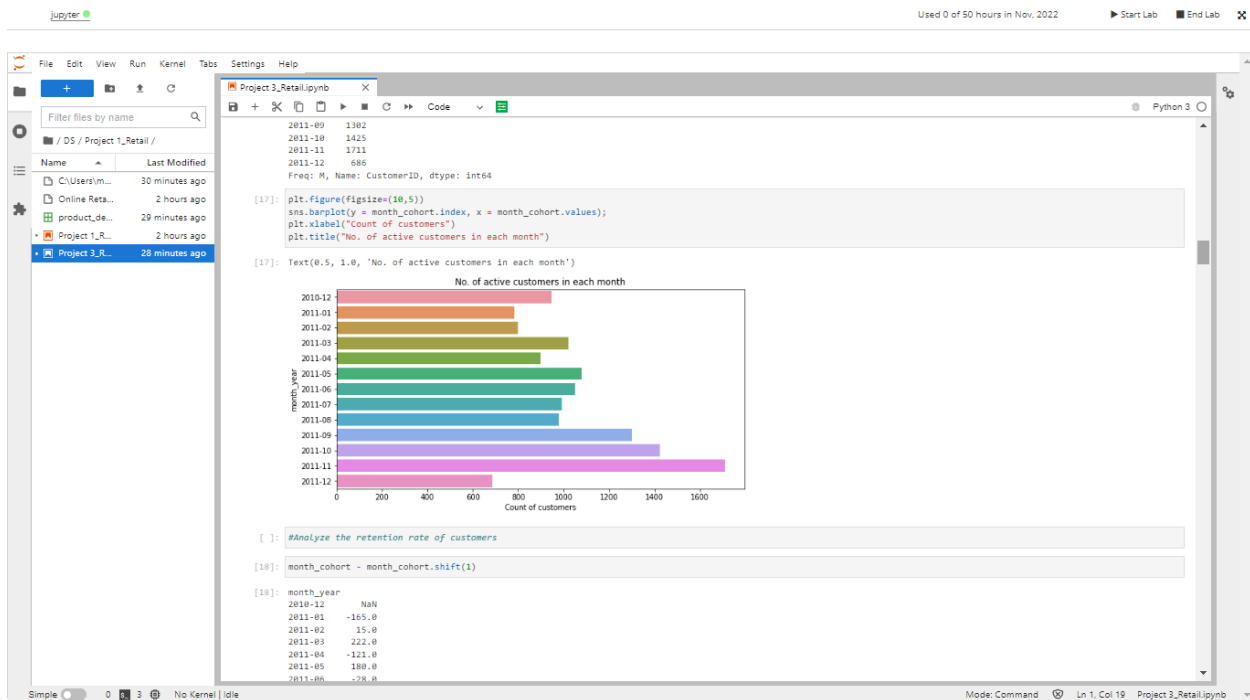
[15]: 13

[16]: month_cohort = df.groupby('month_year')['CustomerID'].nunique()
month_cohort

```

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	1505

Simple 0 3 No Kernel | Idle Mode: Command Ln 1, Col 19 Project 3\_Retail.ipynb



Project 3\_Retail.ipynb

Project Task: Week 2

Data Modelling :

1. Build a RFM (Recency Frequency Monetary) model; Recency means the number of days since a customer made the last purchase. Frequency is the number of purchase in a given period. It could be 3 months, 6 months or 1 year. Monetary is the total amount of money a customer spent in that given period. Therefore, big spenders will be differentiated among other customers such as MVP (Minimum Viable Product) or VIP.
2. Calculate RFM metrics.
3. Build RFM Segments. Give recency, frequency, and monetary scores individually by dividing them into quartiles.

b1. Combine three ratings to get a RFM segment (as strings).

b2. Get the RFM score by adding up the three ratings.

b3. Analyze the RFM segments by summarizing them and comment on the findings.

Note: Rate "recency" for customer who has been active more recently higher than the less recent customer, because each company wants its customers to be recent.

Note: Rate "frequency" and "monetary" higher, because the company wants the customer to visit more often and spend more money

```
[ ]: #Monetary analysis
[21]: df['amount'] = df['Quantity']*df['UnitPrice']
df.head()

[21]:
```

	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	844068	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

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

Project 3\_Retail.ipynb

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

[22]:
```

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

```
[ ]: #Frequency Analysis
[23]: df_frequency = df.groupby('CustomerID').nunique()[['InvoiceNo']].reset_index()
df_frequency = df.drop_duplicates('InvoiceNo').groupby('CustomerID').count()[['InvoiceNo']].reset_index()
df_frequency

[23]:
```

CustomerID	InvoiceNo
0	12346.0
1	12347.0
2	12348.0
3	12349.0

Project 3\_Retail.ipynb

Used 0 of 50 hours in Nov, 2022

Start Lab End Lab

File Edit View Run Kernel Tabs Settings Help

Filter files by name

/ DS / Project 1\_Retail /

Name	Last Modified
C:\Users\m...	30 minutes ago
Online Reta...	2 hours ago
product_de...	29 minutes ago
Project 1_R...	2 hours ago
Project 3_R...	29 minutes ago

[23]:

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

[25]: #Recency Analysis

[26]:

```
# We will fix reference date for calculating recency 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()
```

[26]:

InvoiceNo	StockCode	Quantity	InvoiceDate	UnitPrice	CustomerID	Country	month_year	amount	days_to_last_order
0	536965	85123A	6	2010-12-01 08:26:00	2.55	17850.0	United Kingdom	2010-12	15.30
1	536965	71053	6	2010-12-01 08:26:00	3.39	17850.0	United Kingdom	2010-12	20.34
2	536965	844068	8	2010-12-01 08:26:00	2.75	17850.0	United Kingdom	2010-12	22.00
3	536965	84029G	6	2010-12-01 08:26:00	3.39	17850.0	United Kingdom	2010-12	20.34
4	536965	84029E	6	2010-12-01 08:26:00	3.39	17850.0	United Kingdom	2010-12	20.34

[27]:

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

Simple 0 3 No Kernel | Idle

Mode: Command Ln 1, Col 19 Project 3\_Retail.ipynb

Project 3\_Retail.ipynb

Used 0 of 50 hours in Nov, 2022

Start Lab End Lab

File Edit View Run Kernel Tabs Settings Help

Filter files by name

/ DS / Project 1\_Retail /

Name	Last Modified
C:\Users\m...	30 minutes ago
Online Reta...	2 hours ago
product_de...	29 minutes ago
Project 1_R...	3 hours ago
Project 3_R...	29 minutes ago

[27]:

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

[ ]: #Calculate RFM metrics

[28]:

```
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()
```

[28]:

CustomerID	Recency	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

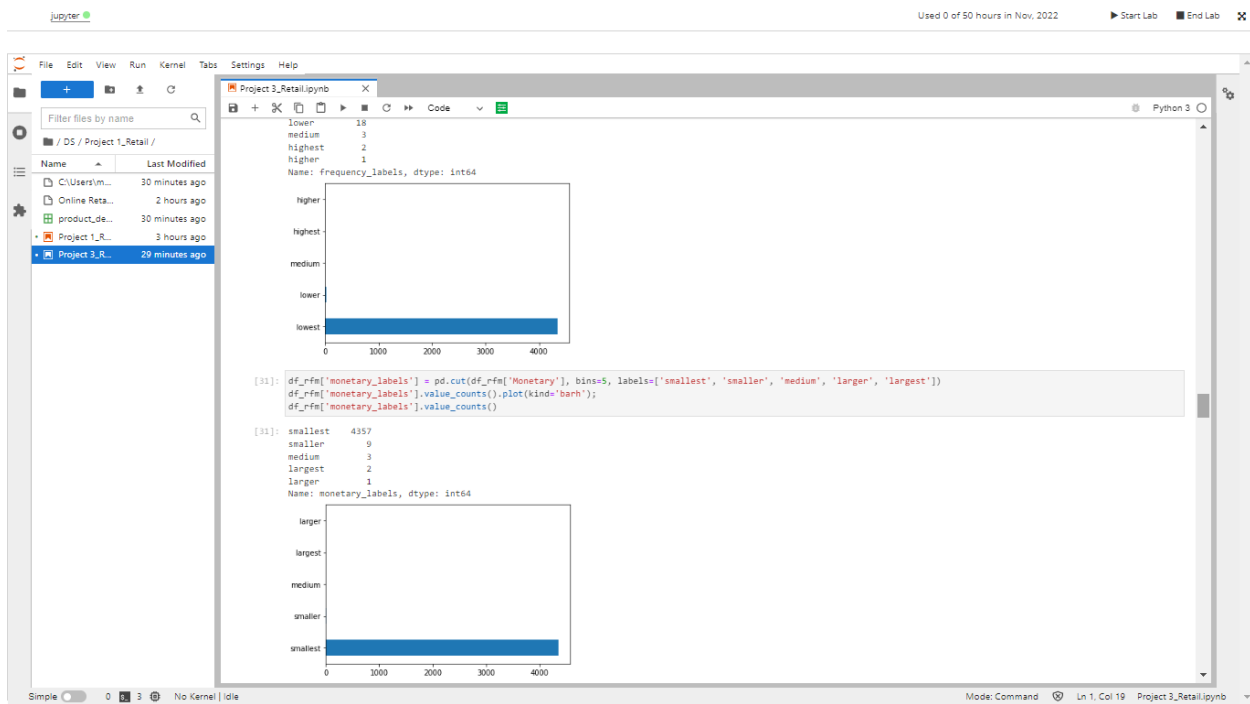
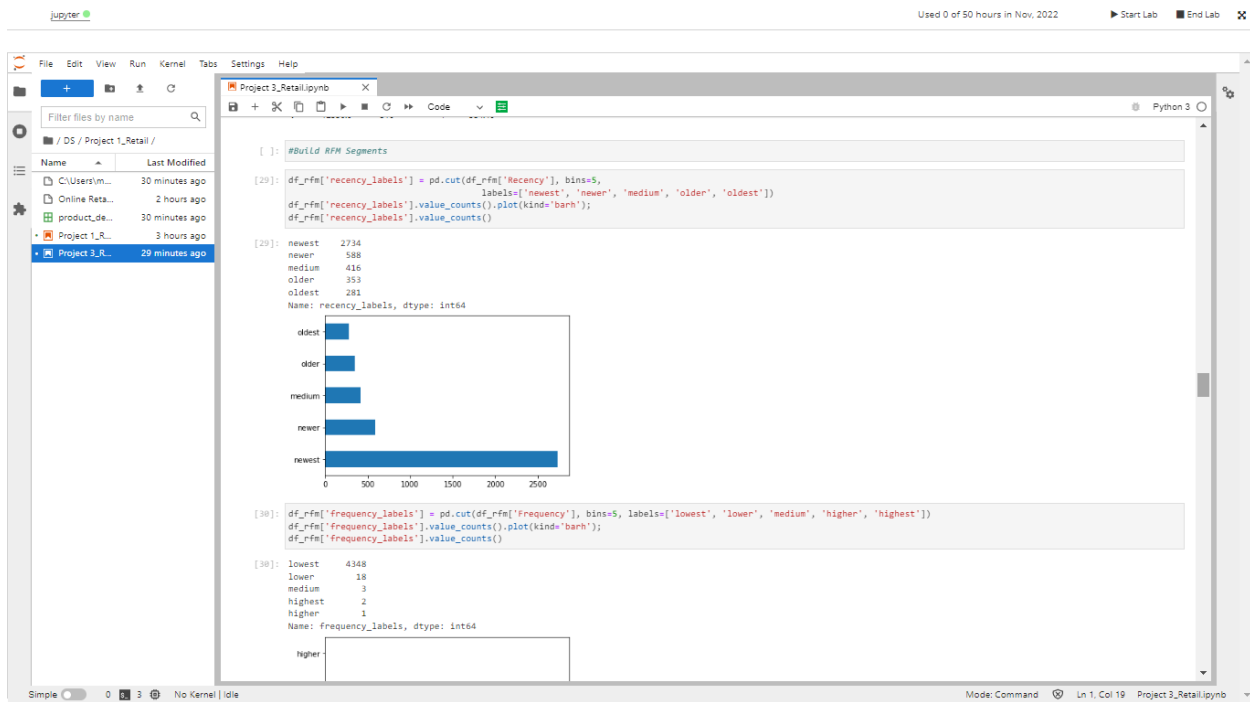
[ ]: #Build RFM Segments

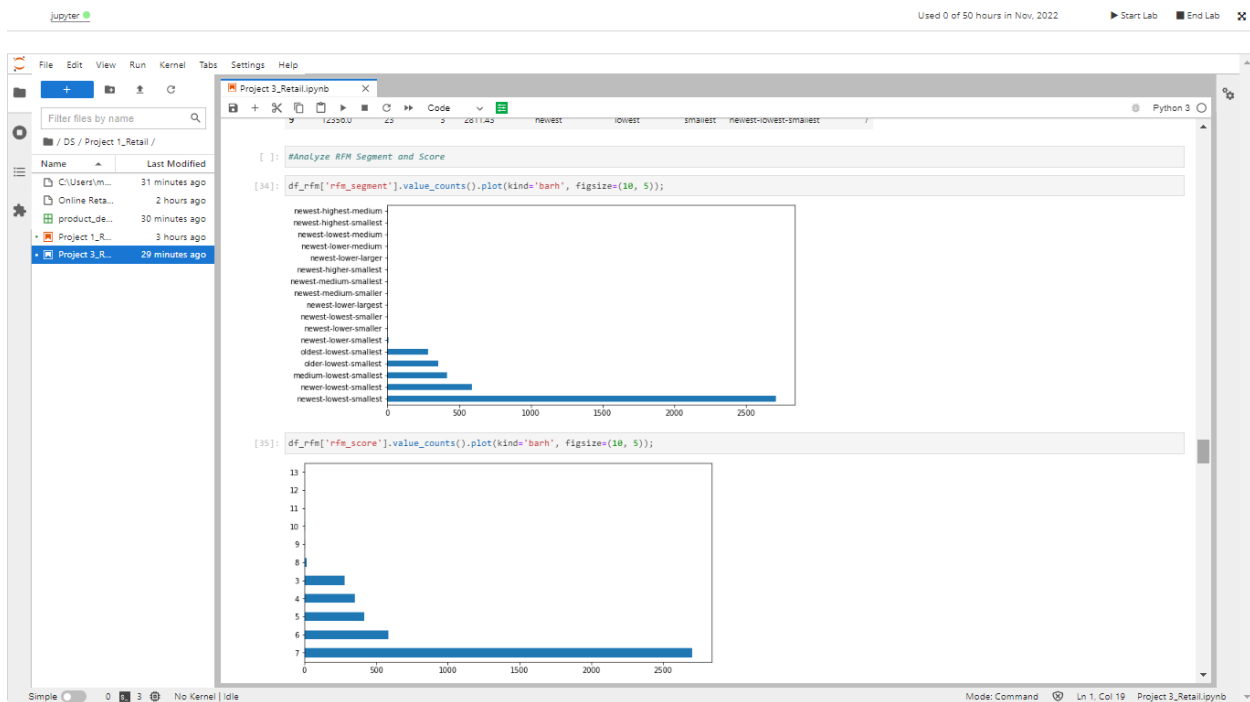
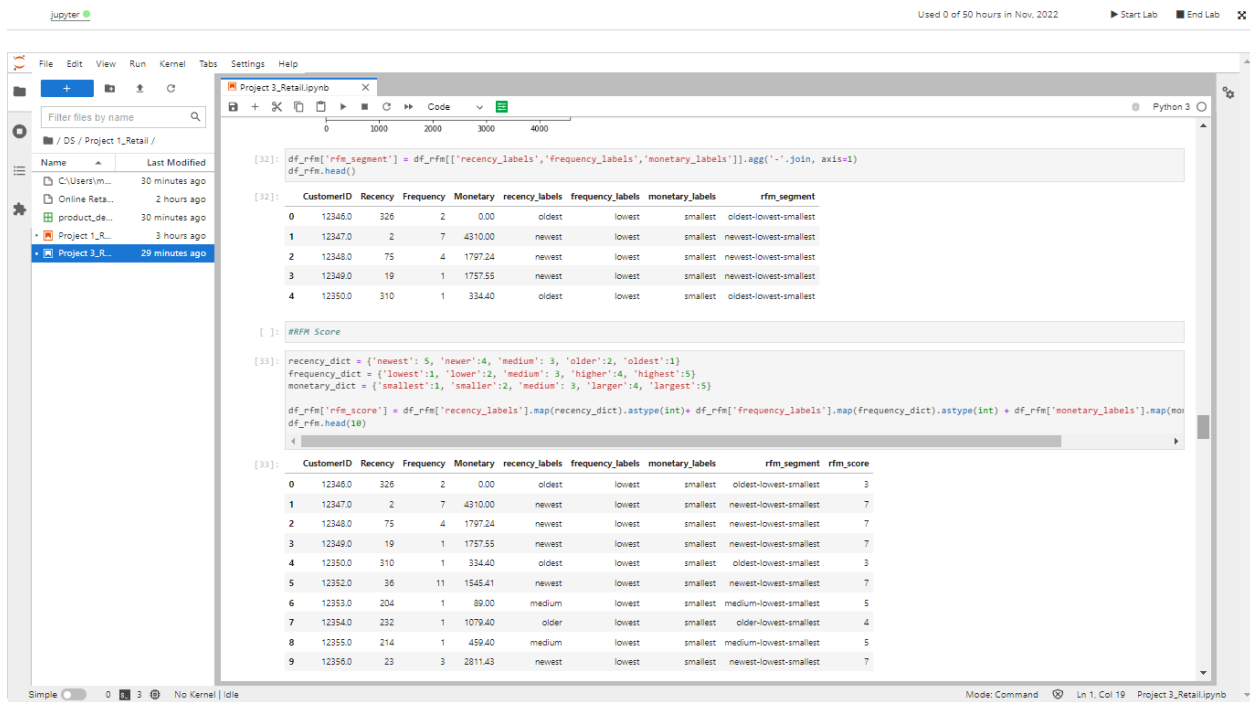
[29]:

```
df_rfm['recency_labels'] = pd.cut(df_rfm['Recency'], bins=5,
```

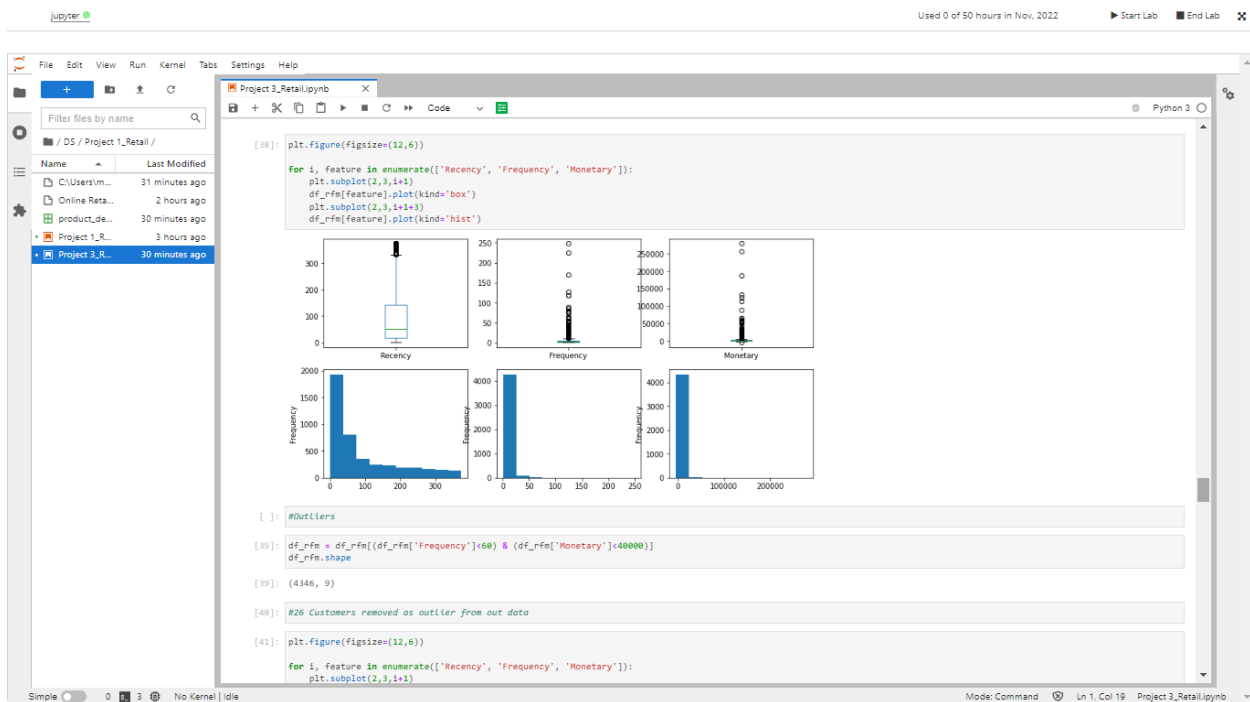
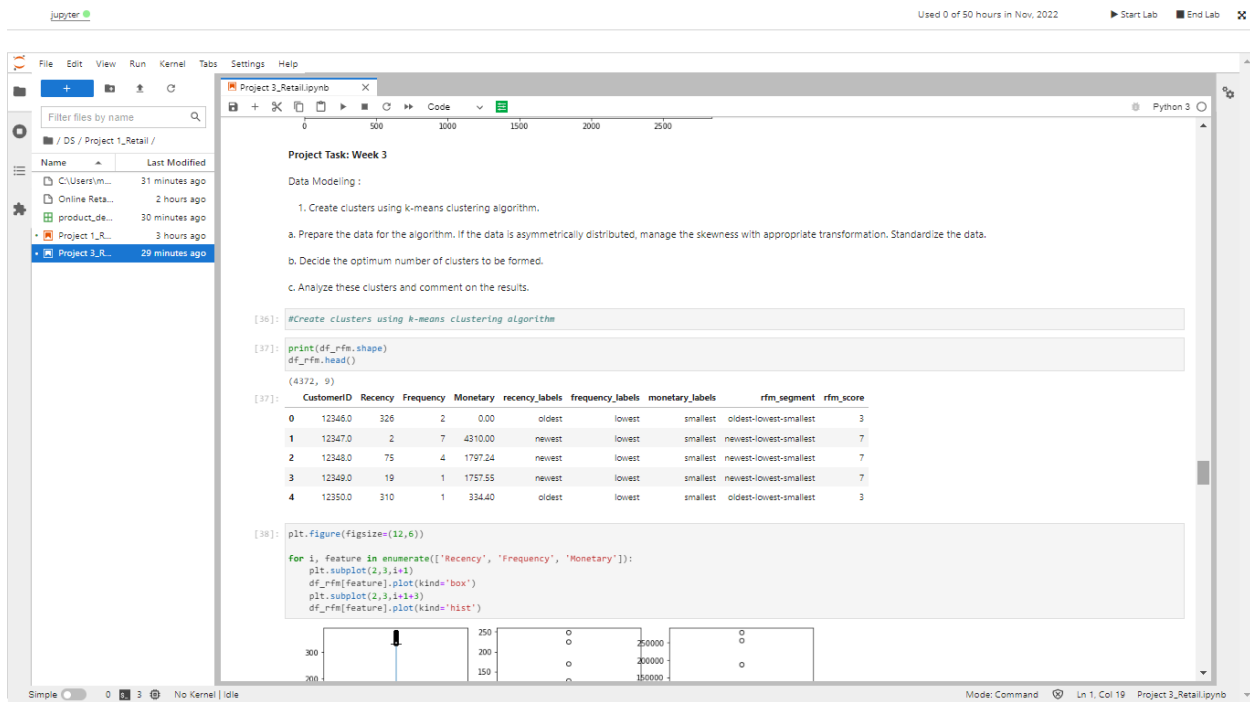
Simple 0 3 No Kernel | Idle

Mode: Command Ln 1, Col 19 Project 3\_Retail.ipynb

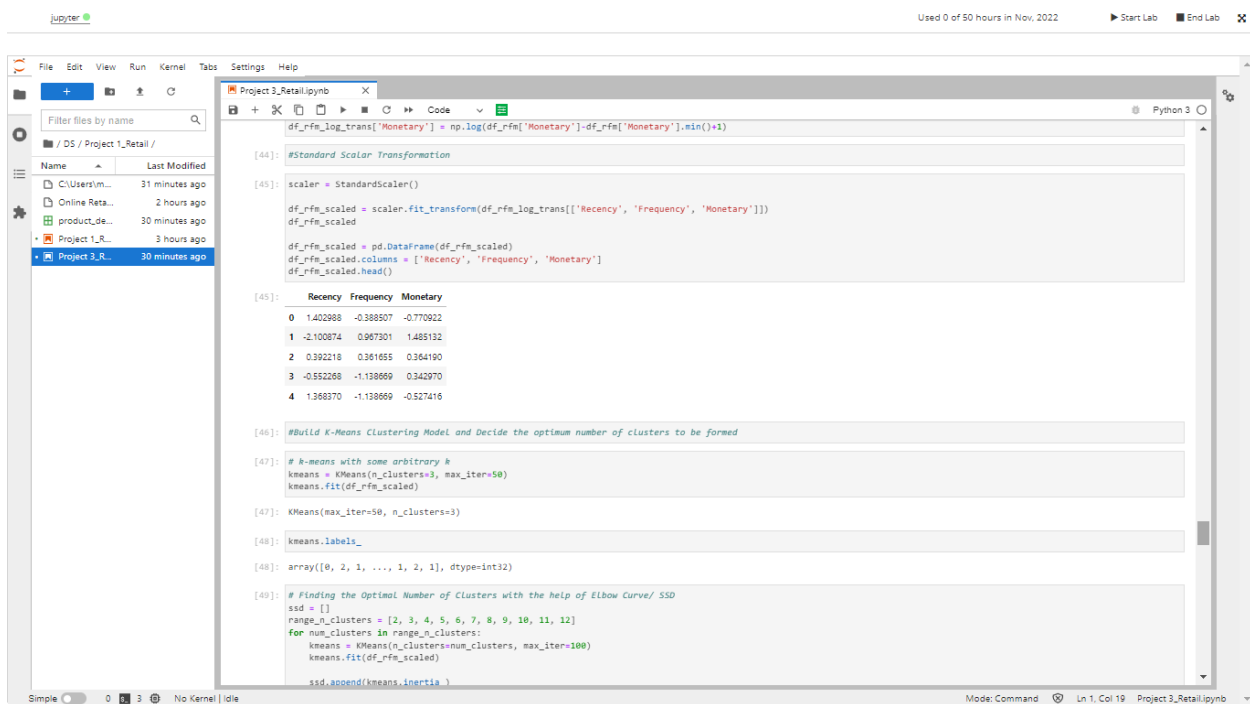
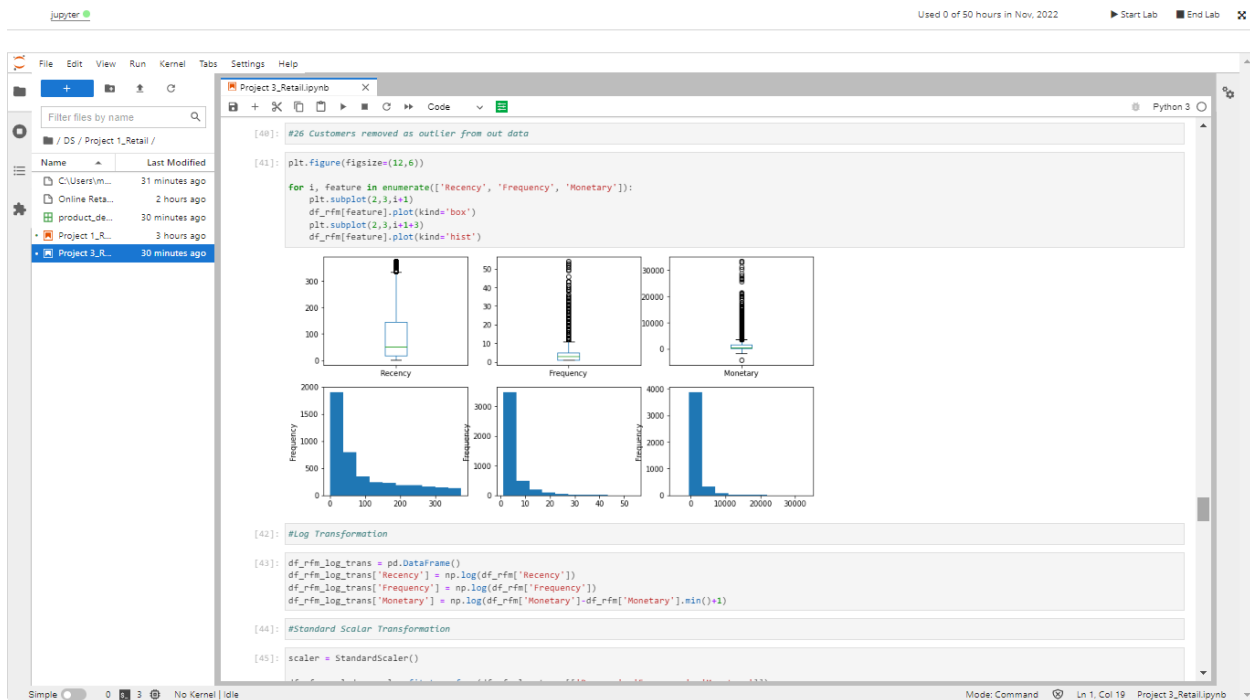


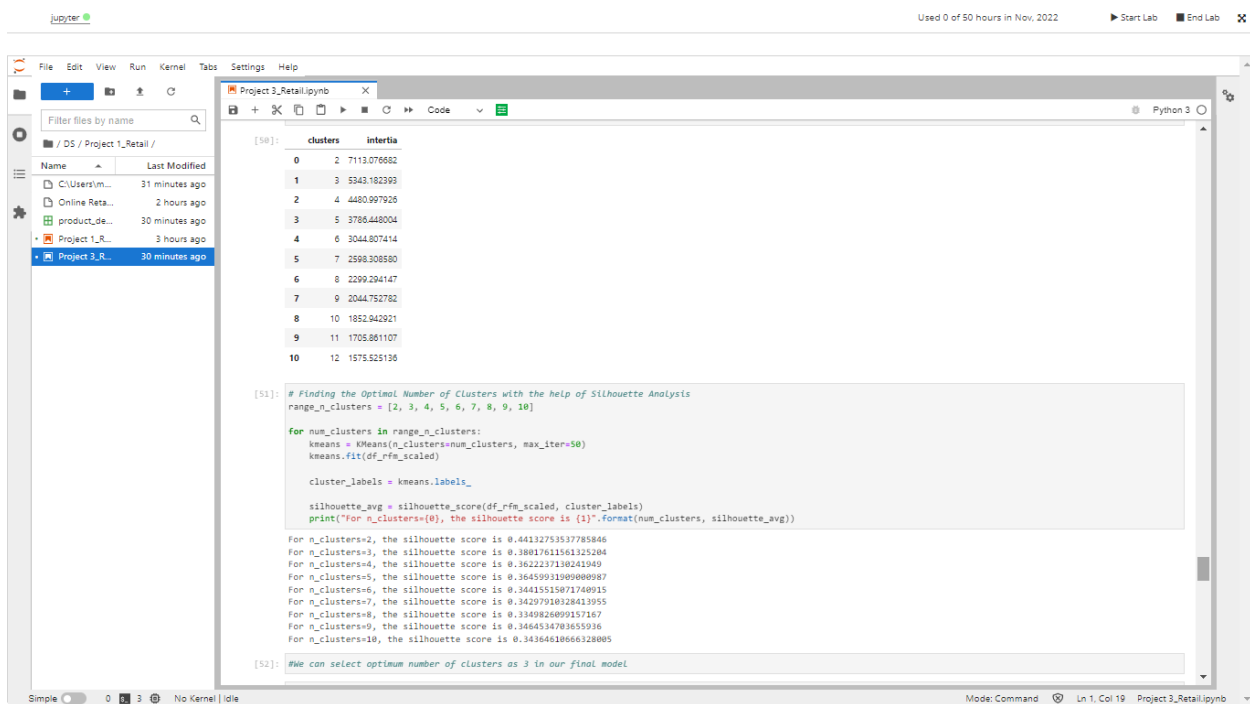
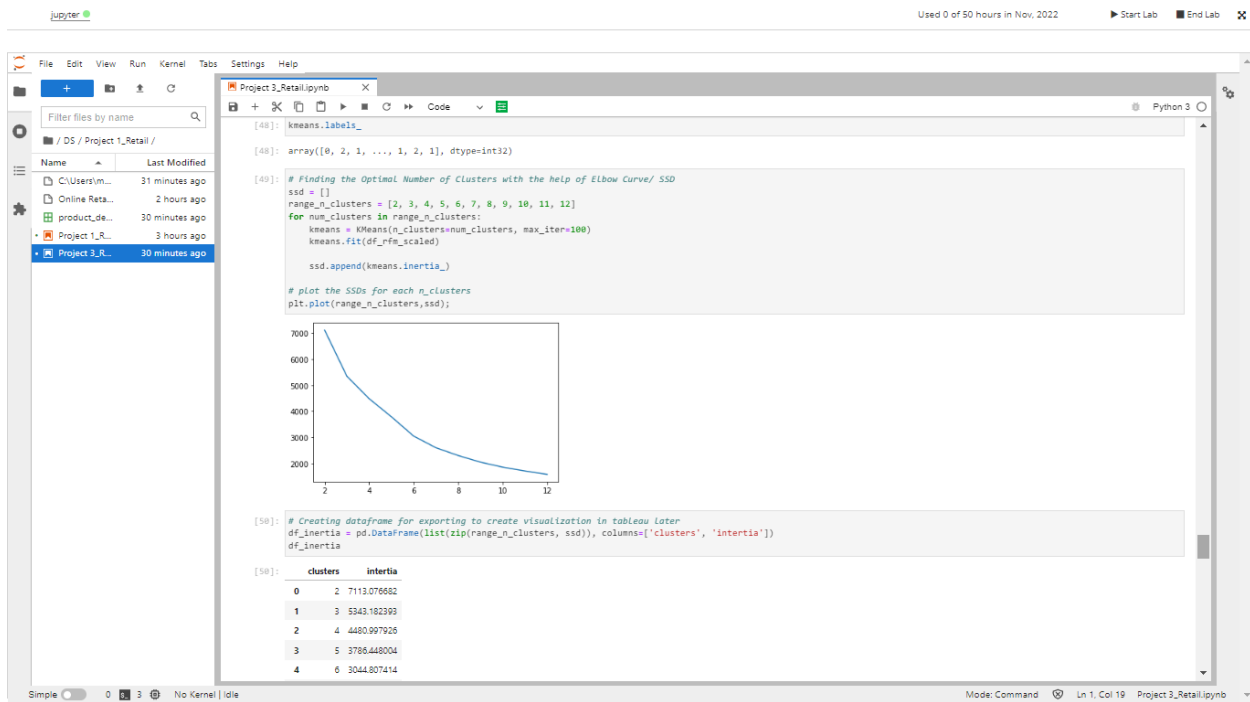






```
[ ]: #Outliers
[39]: df_rfm = df_rfm[(df_rfm['Frequency'] < 60) & (df_rfm['Monetary'] < 40000)]
df_rfm.shape
[39]: (4346, 9)
[40]: #26 Customers removed as outlier from out data
[41]: plt.figure(figsize=(12,6))
for i, feature in enumerate(['Recency', 'Frequency', 'Monetary']):
    plt.subplot(2,3,i+1)
```





Used 0 of 50 hours in Nov, 2022 Start Lab End Lab

Project 3\_Retail.ipynb

```
[50]: clusters    inertia
      0      2  7113.076682
      1      3   5343.182393
      2      4  4480.997926
      3      5  3786.448004
      4      6  3044.807414
      5      7  2598.308580
      6      8  2299.294147
      7      9  2044.752782
      8     10  1852.942921
      9     11  1705.861107
     10     12  1575.525136

[51]: # 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 0.38817611561325204
For n_clusters=4, the silhouette score is 0.3622371382433499
For n_clusters=5, the silhouette score is 0.36459931090900987
For n_clusters=6, the silhouette score is 0.34415515871748915
For n_clusters=7, the silhouette score is 0.342979180328413955
For n_clusters=8, the silhouette score is 0.3349824899515167
For n_clusters=9, the silhouette score is 0.3464534703655936
For n_clusters=10, the silhouette score is 0.34364618666328005

[52]: #We can select optimum number of clusters as 3 in our final model
```

Simple 0 3 No Kernel | Idle Mode: Command Ln 1, Col 19 Project 3\_Retail.ipynb

Used 0 of 50 hours in Nov, 2022 Start Lab End Lab

Project 3\_Retail.ipynb

```
[52]: #We can select optimum number of clusters as 3 in our final model

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

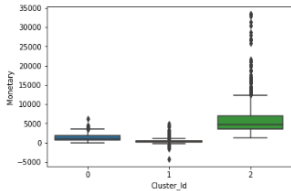
[54]: KMeans(max_iter=50, n_clusters=3)

[54]: #Analyze these clusters and comment on the results

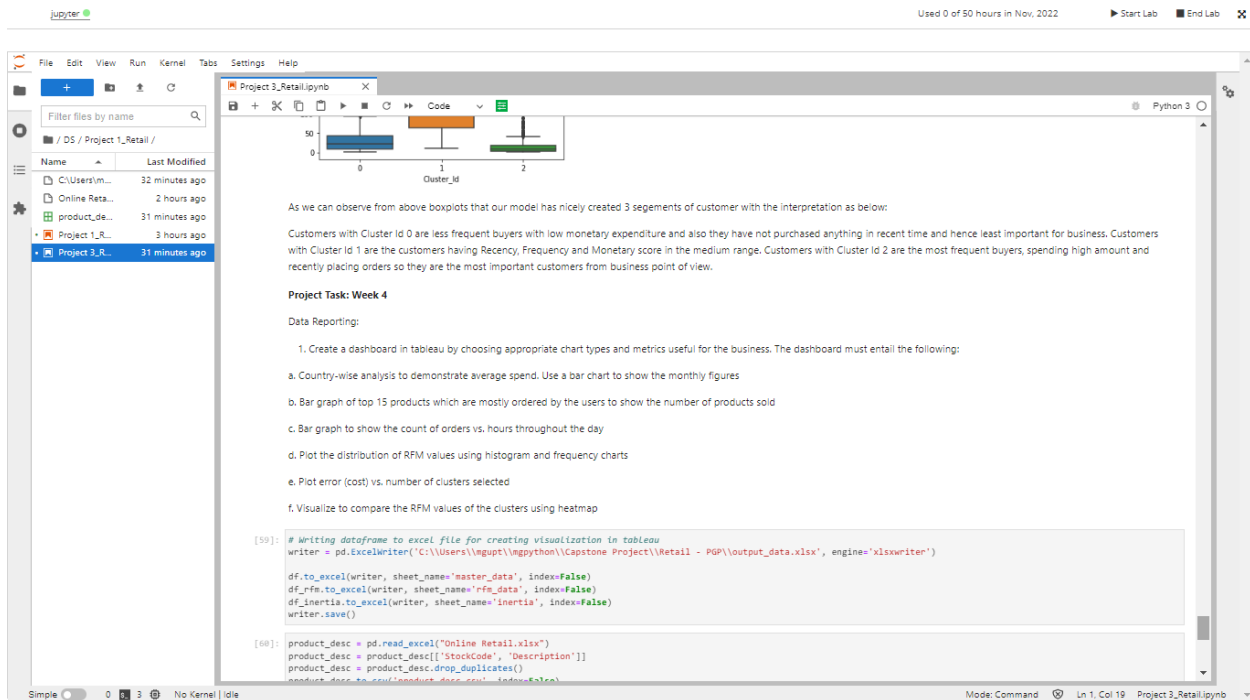
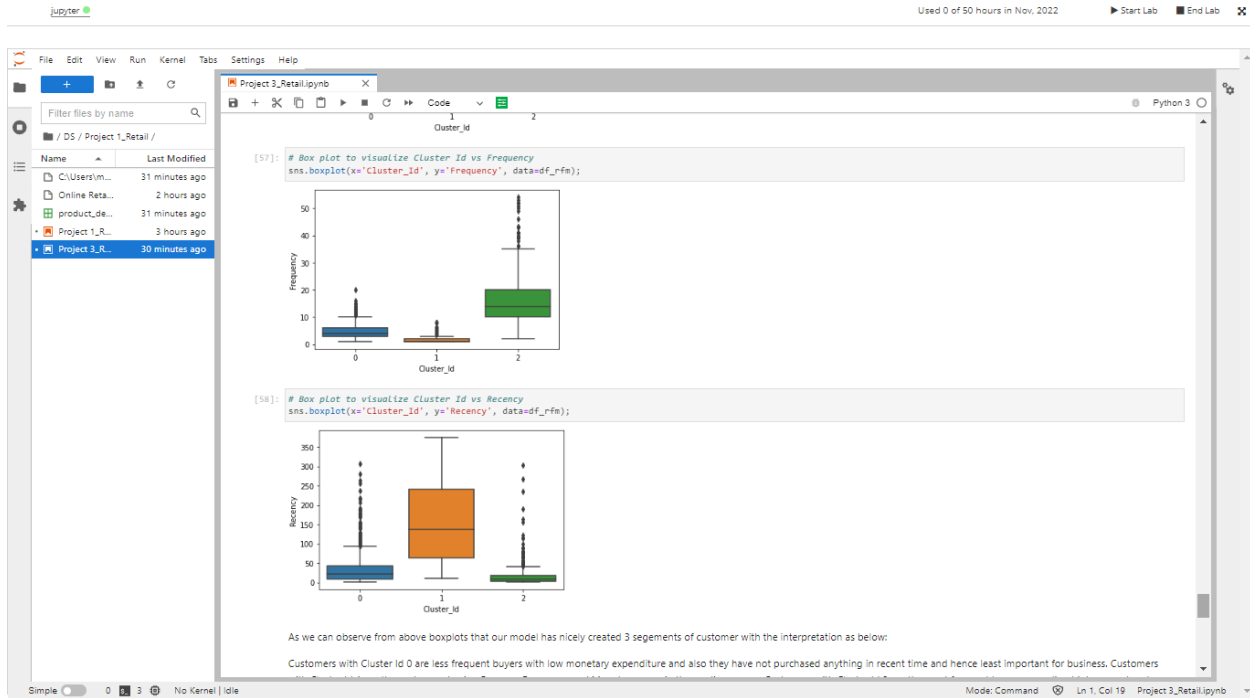
[55]: # assign the Label
df_rfm['Cluster_Id'] = kmeans.labels_
df_rfm.head()

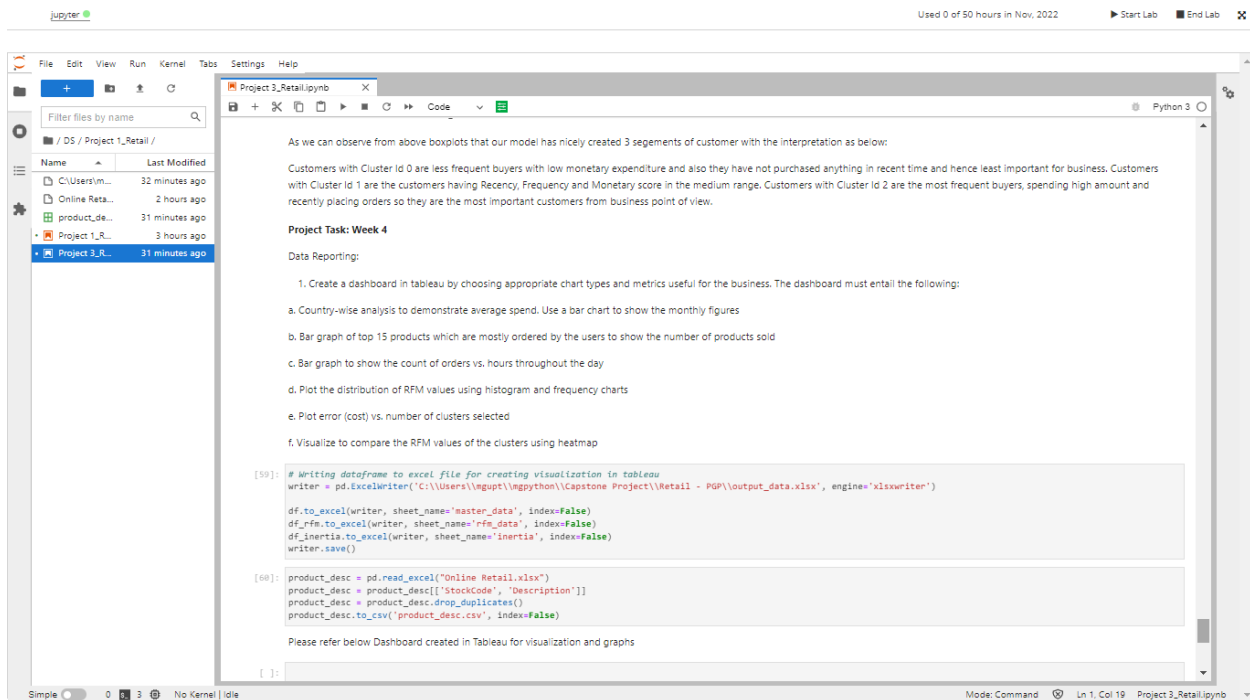
[55]: CustomerID  Recency  Frequency  Monetary  recency_labels  frequency_labels  monetary_labels  rfm_segment  rfm_score  Cluster_Id
0      12345.0    326         2         0.00      oldest        lowest        smallest      oldest-lowest-smallest    3         1
1      12347.0         2         7    4310.00      newest        lowest        smallest      newest-lowest-smallest    7         2
2      12348.0         75         4    1797.24      newest        lowest        smallest      newest-lowest-smallest    7         0
3      12349.0        19         1    1757.55      newest        lowest        smallest      newest-lowest-smallest    7         1
4      12350.0        310         1     334.40      oldest        lowest        smallest      oldest-lowest-smallest    3         1

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



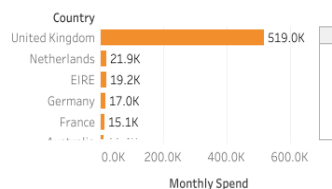
Simple 0 3 No Kernel | Idle Mode: Command Ln 1, Col 19 Project 3\_Retail.ipynb



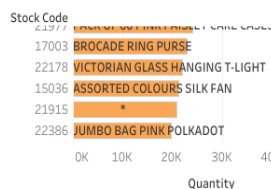


## DASHBOARD FOR CUSTOMER ANALYSIS

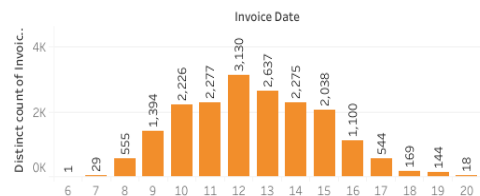
COUNTRYWISE MONTHLY SPEND



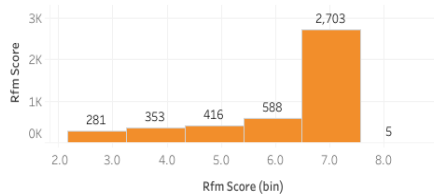
TOP 15 PRODUCTS



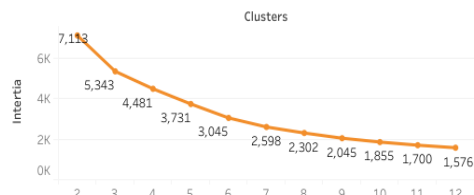
NUMBER OF ORDERS PLACED AT DIFFERENT HOURS



DISTRIBUTION OF RFM SCORES



INERTIA VS. NUMBER OF CLUSTERS



RFM SCORE HEATMAP

