

Configuration Manual

MSc Research Project
Data Analytics

Amol Upadhyay
Student ID: X20110812

School of Computing
National College of Ireland

Supervisor: Giovani Estrada

National College of Ireland
Project Submission Sheet
School of Computing



Student Name:	Amol Upadhyay
Student ID:	X20110812
Programme:	Data Analytics
Year:	2022
Module:	MSc Research Project
Supervisor:	Giovani Estrada
Submission Due Date:	15/12/2022
Project Title:	Configuration Manual
Word Count:	414
Page Count:	10

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

ALL internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

Signature:	
Date:	15th December 2022

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:

Attach a completed copy of this sheet to each project (including multiple copies).	<input type="checkbox"/>
Attach a Moodle submission receipt of the online project submission , to each project (including multiple copies).	<input type="checkbox"/>
You must ensure that you retain a HARD COPY of the project , both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer.	<input type="checkbox"/>

Assignments that are submitted to the Programme Coordinator office must be placed into the assignment box located outside the office.

Office Use Only	
Signature:	
Date:	
Penalty Applied (if applicable):	

Configuration Manual

Amol Upadhyay
X20110812

1 Introduction

A complete configuration manual contains critical information on the system, hardware, and software specs. It also displays the whole flow employed to execute the study on the "An investigation of e-commerce buying behaviour across the UK and Brazil". Section 2 of the handbook discusses system specifications, such as hardware and software setups. In Section 3, you'll learn how to set up the environment, import key libraries, and do pre-processing. The fourth segment discusses model development and assessment.

2 Configuration of System

System setup discussed the hardware and software requirements for the investigation.

2.1 Specification of Hardware

OS(operating System)	Windows 10
RAM	16GB
System Processor	i6intel
Speed in Hz	3.2GHz
Disk Memory	1GB approx.
GPU	NVIDIA

2.2 Specification of Software

Programming Language	Python 3.9 version
Other Softwares	Anaconda & jupyter
Web Browsers	Google Chrome

3 Setting Environment

3.1 Launching Jupyter on Anaconda

The very first stage is to activate the Anaconda app. Anaconda delivers a wealth of relevant software to meet your requirements. Jupyter Notebook used to implement the code above, and it comes with the most recent version of Python. There are many excellent Python libraries for analysis.

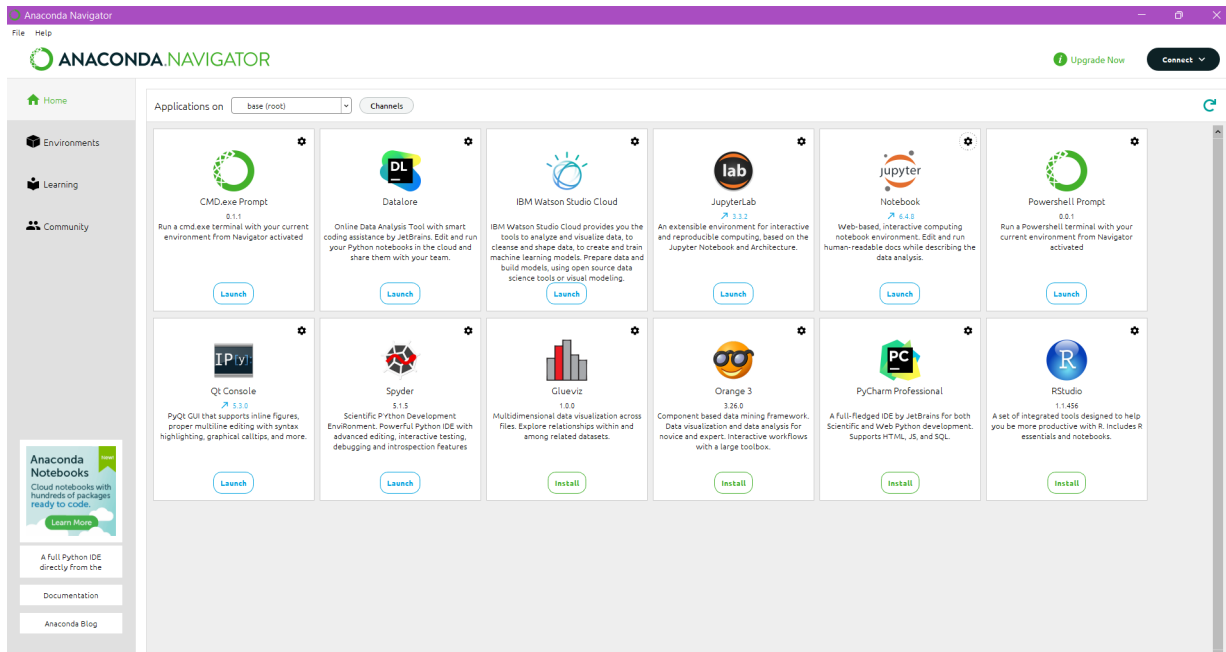


Figure 1: anaconda interface

This is the main page for Jupyter Notebook; to begin coding, create a new ipynb file.

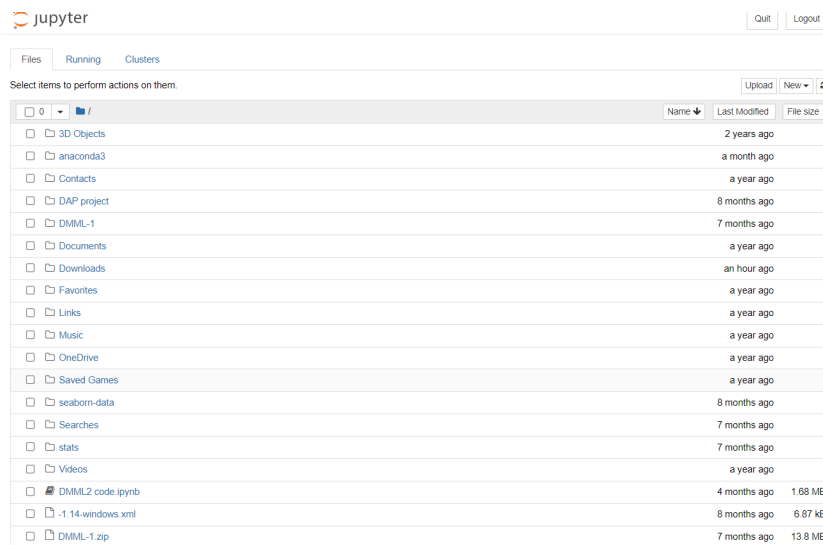


Figure 2: Jupyter

3.2 Data Preparation

Keep all the dataset files in one folder before diving into the code. All the required files are in one place in the folder MachineLearningCVE.

» This PC > Desktop > Research Project-Amol > Braziliandataset

Name	Status	Date modified	Type	Size
olist_customers_dataset	✔	04-11-2022 12:44	Microsoft Excel Com...	8,823 KB
olist_geolocation_dataset	✘	04-11-2022 12:44	Microsoft Excel Com...	59,838 KB
olist_order_items_dataset	✔	04-11-2022 12:44	Microsoft Excel Com...	15,077 KB
olist_order_payments_dataset	✔	04-11-2022 12:44	Microsoft Excel Com...	5,642 KB
olist_order_reviews_dataset	✔	04-11-2022 12:44	Microsoft Excel Com...	14,113 KB
olist_orders_dataset	✔	04-11-2022 12:44	Microsoft Excel Com...	17,242 KB
olist_products_dataset	✔	04-11-2022 12:44	Microsoft Excel Com...	2,324 KB
olist_sellers_dataset	✔	04-11-2022 12:44	Microsoft Excel Com...	171 KB
product_category_name_translation	✔	04-11-2022 12:44	Microsoft Excel Com...	3 KB

Figure 3: CSV Files-Brazil

» This PC > Desktop > Research Project-Amol > UKdataset

Name	Status	Date modified	Type	Size
data	✔	21-10-2022 00:32	Microsoft Excel Com...	44,513 KB

Figure 4: CSV File-UK

3.3 Importing Important Libraries

It displays all implementation-required libraries, and libraries that need installation may be installed with "pip."

```
import pandas as pd
import numpy as np
import urllib
import unicode
import matplotlib.pyplot as plt
import seaborn as sns
import datetime
from scipy import stats
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler
import warnings
warnings.filterwarnings("ignore")
```

Figure 5: Libraries for Brazil Market Analysis

```

import missingno as msno |

import gc
import datetime as dt

%matplotlib inline
color = sns.color_palette()

from matplotlib.ticker import PercentFormatter
from sklearn.preprocessing import StandardScaler
from sklearn.cluster import KMeans
from yellowbrick.cluster import KElbowVisualizer
from itertools import combinations
import statsmodels.api as sm
import warnings
from statsmodels.tsa.stattools import adfuller
from statsmodels.tsa.seasonal import seasonal_decompose
from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
from scipy import stats
from sklearn.metrics import r2_score, mean_squared_error

pd.options.mode.chained_assignment = None

plt.rcParams["axes.facecolor"] = "#A2A2A2"
plt.rcParams["axes.grid"] = 1

```

Figure 6: Libraries for UK Market Analysis

```

#Import all relevant libraries
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from mlxtend.frequent_patterns import apriori
from mlxtend.frequent_patterns import association_rules

```

Figure 7: Libraries for Market Basket Analysis for both Uk and Brazil

3.4 Importing Dataset

The pandas package is used to load each CSV file into a python data frame.

```

In [3]: df = pd.read_csv('C:/Users/Amo1/OneDrive/Desktop/Research Project-Amo1/UKdataset/data.csv', encoding='latin')
# df = pd.read_csv('data.csv', encoding='Latin')
df.head()

```

Figure 8: Importing UK Dataset

```

import missingno as msno |

import gc
import datetime as dt

%matplotlib inline
color = sns.color_palette()

from matplotlib.ticker import PercentFormatter
from sklearn.preprocessing import StandardScaler
from sklearn.cluster import KMeans
from yellowbrick.cluster import KElbowVisualizer
from itertools import combinations
import statsmodels.api as sm
import warnings
from statsmodels.tsa.stattools import adfuller
from statsmodels.tsa.seasonal import seasonal_decompose
from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
from scipy import stats
from sklearn.metrics import r2_score, mean_squared_error

pd.options.mode.chained_assignment = None

plt.rcParams["axes.facecolor"] = "#A2A2A2"
plt.rcParams["axes.grid"] = 1

```

Figure 9: Importing MBA UK Dataset

```

df = pd.read_csv('C:/Users/Amol/OneDrive/Desktop/Research Project-Amol/UKdataset/data.csv', encoding="ISO-8859-1",
                dtype={'CustomerID': str, 'InvoiceID': str}) # We converted these datasets into string data type.
print('DataFrame dimensions:', df.shape)
# Used to understand the dimension of the dataset.
df['InvoiceDate'] = pd.to_datetime(df['InvoiceDate'])

# gives some info on columns types and number of null values
tab_info=pd.DataFrame(df.dtypes).T.rename(index={0: 'column type'})
tab_info-tab_info.append(pd.DataFrame(df.isnull().sum()).T.rename(index={0: 'null values (nb)'}))
tab_info-tab_info.append(pd.DataFrame(df.isnull().sum()/df.shape[0]*100).T.rename(index={0: 'null values (%)'}))
display(tab_info)

# show first lines
display(df[:5])

```

Figure 10: Importing MBA Brazil Dataset

3.5 Data Pre-processing

In this part processes like eliminating missing numbers, outliers, etc are done. During this process, data is changed into a form that can be used in modeling.

3.5.1 Brazil Dataset Pre-processing

```

In [2]: dfs_paths = ["C:/Users/Amol/OneDrive/Desktop/Research Project-Amol/Brazilliandataset/olist_customers_dataset.csv", "C:/Users/Amol/OneDrive/Desktop/Research Project-Amol/Brazilliandataset/olist_order_items_dataset.csv", "C:/Users/Amol/OneDrive/Desktop/Research Project-Amol/Brazilliandataset/olist_order_reviews_dataset.csv", "C:/Users/Amol/OneDrive/Desktop/Research Project-Amol/Brazilliandataset/olist_products_dataset.csv", "C:/Users/Amol/OneDrive/Desktop/Research Project-Amol/Brazilliandataset/product_category_name_translation.csv"]

Data Pre-Processing(Cleaning and Transformation)

In [3]: customers_df = pd.read_csv("C:/Users/Amol/OneDrive/Desktop/Research Project-Amol/Brazilliandataset/olist_customers_dataset.csv")
geo_df = pd.read_csv("C:/Users/Amol/OneDrive/Desktop/Research Project-Amol/Brazilliandataset/olist_geolocation_dataset.csv")
orderitem_df = pd.read_csv("C:/Users/Amol/OneDrive/Desktop/Research Project-Amol/Brazilliandataset/olist_order_items_dataset.csv")
orderpay_df = pd.read_csv("C:/Users/Amol/OneDrive/Desktop/Research Project-Amol/Brazilliandataset/olist_order_payments_dataset.csv")
orderreviews_df = pd.read_csv("C:/Users/Amol/OneDrive/Desktop/Research Project-Amol/Brazilliandataset/olist_order_reviews_dataset.csv")
orders_df = pd.read_csv("C:/Users/Amol/OneDrive/Desktop/Research Project-Amol/Brazilliandataset/olist_orders_dataset.csv")
products_df = pd.read_csv("C:/Users/Amol/OneDrive/Desktop/Research Project-Amol/Brazilliandataset/olist_products_dataset.csv")
sellers_df = pd.read_csv("C:/Users/Amol/OneDrive/Desktop/Research Project-Amol/Brazilliandataset/olist_sellers_dataset.csv")
categoryname_df = pd.read_csv("C:/Users/Amol/OneDrive/Desktop/Research Project-Amol/Brazilliandataset/product_category_name_translation.csv")
pd.set_option('display.max_columns', 500)

In [4]: # Change cols names before merging
customers_df.rename(columns={"customer_zip_code_prefix": "zip_code"}, inplace=True)
geo_df.rename(columns={"geolocation_zip_code_prefix": "zip_code"}, inplace=True)

In [5]: # Join datasets
data = orders_df.merge(customers_df, on="customer_id").merge(orderitem_df, on="order_id").merge(products_df, on="product_id").mer

```

Figure 11: Data Pre-processing part 1

```

: # Null values
(data.isna().sum() / len(data)).sort_values(ascending=False)

```

Figure 12: Data Pre-processing part 2

```

le = LabelEncoder()
dframe['Label'] = le.fit_transform(dframe['Label'])

X = dframe.drop(['Label'], axis=1)
X[X < 0] = 0
print(X)
y = dframe['Label']
del dframe

X_new = SelectKBest(chi2, k=30).fit_transform(X, y)
X_new.shape

(2062474, 30)

X_train, X_test, y_train, y_test = train_test_split(X_new, y, test_size=0.25, random_state=1234)

St_scalar = StandardScaler()
X_train = St_scalar.fit_transform(X_train)
X_test = St_scalar.transform(X_test)

```

Figure 13: More Data Pre-processing Steps

3.5.2 UK Dataset Pre-processing

```

In [6]: df.rename(columns={"InvoiceNo": "invoice_num",
                        "StockCode": "stock_code",
                        "Description": "description",
                        "Quantity": "quantity",
                        "InvoiceDate": "invoice_date",
                        "UnitPrice": "unit_price",
                        "CustomerID": "customer_id",
                        "Country": "country"}, inplace=True)

In [7]: df['invoice_date'] = pd.to_datetime(df.invoice_date, format='%m/%d/%Y %H:%M')

In [8]: df['description'] = df.description.str.lower()

In [9]: df.isnull().sum()

```

Figure 14: Data Pre-processing part 1

```

In [20]: df_miss["day"] = df_miss['invoice_date'].map(lambda x: x.day)
df_miss["month"] = df_miss['invoice_date'].map(lambda x: x.month)
df_miss["year"] = df_miss['invoice_date'].map(lambda x: x.year)

In [21]: df_miss['daymonth'] = df_miss['day'].astype(str)+'/'+df_miss['month'].astype(str)
df_miss['daymonthyear'] = df_miss['daymonth'].astype(str)+'/'+df_miss['year'].astype(str)
df_miss['monthyear'] = df_miss['month'].astype(str)+'/'+df_miss['year'].astype(str)

```

Figure 15: Data Pre-processing part 2

```

In [33]: df_new.duplicated().sum()
Out[33]: 5225

In [34]: df_new.drop_duplicates(inplace=True)

In [35]: df_new.duplicated().sum()
Out[35]: 0

In [36]: df_new['customer_id'] = df_new.customer_id.astype('int64')

```

Figure 16: Data Pre-processing Part 3


```
In [49]: import re
spec_list=[]
for code in df_new.stock_code:
    x=re.findall(r"^\w{1}$|\D[A-Z]+\D[A-Z]\d", code)
    if x not in spec_list:
        if len(x) >0 :
            spec_list.append(x)
spec_list

Out[49]: [['POST'], ['PADS'], ['M'], ['DOT'], ['C2'], ['BANK ', 'CHARGES']]

In [50]: spec_list[5] = ['BANK CHARGES']

In [51]: spec_list2=[item for sublist in spec_list for item in sublist]
spec_list2

Out[51]: ['POST', 'PADS', 'M', 'DOT', 'C2', 'BANK CHARGES']

In [52]: df_new[df_new['stock_code'].apply(lambda x: x in spec_list2)]
```

Figure 17: Data Pre-processing Part 4

3.5.3 MBA Data Pre-processing for Both Brazil & UK

```
#If you want to cross check the data you can replace 10 with 80.
item_freq = data['product_id'].value_counts()
data = data[data.isin(item_freq.index[item_freq >= 10]).values]
data['product_id'].value_counts()

422079e10f46682990e24d770e7f83d    137
aca2eb7d09e0a1a708ebd4e68314663af    130
99a4788cb2485695c36a24e339b06058    125
308c6c730842d78016a0823897a3720b    113
380u1104dcf3045d311335e4909c6b     104
...
65206b2da20004d0e0c5c2d30b7859e     10
4c1a109ecdf50432e3e5217cef404c      10
657247f6f60543b93e3cc788a8e6329c     10
98d61056e8560a048e5d78038790e77      10
c7fd135e515bffdab85508012842e0b       10
Name: product_id, Length: 329, dtype: int64

#Average products purchased per transaction
data['order_id'].value_counts().mean()

1.154075540719682

#Create a basket of all products, orders and quantity
basket_brazil = (data.groupby(['order_id', 'product_id'])['quantity']).sum().unstack().fillna(0).set_index('order_id')
basket_brazil.head()
```

Figure 18: Data Pre-processing for MBA Brazil

```
df.dropna(axis = 0, subset = ['CustomerID'], inplace = True)
print('Dataframe dimensions:', df.shape)
#
# gives some info on columns types and number of null values
tab_info=pd.DataFrame(df.dtypes).T.rename(index={0:'column type'})
tab_info.tab_info.append(pd.DataFrame(df.isnull().sum()).T.rename(index={0:'null values (nb)'}))
tab_info=tab_info.append(pd.DataFrame(df.isnull().sum()/df.shape[0]*100).T.rename(index={0:'null values (%)'}))
display(tab_info)
#We deleted ~25% of null present in the CustomerID variable, as all these products are not associated with any Customers.
```

Figure 19: Data Pre-processing MBA UK part 1

```
#After deleting the Null values previously, we have now dropped all the duplicate values present in the dataset.
print('Entrées dupliquées: {}'.format(df.duplicated().sum()))
df.drop_duplicates(inplace = True)

Entrées dupliquées: 5225

#Remove rows with invoices that contain a "C"
df = df[~df['InvoiceNo'].str.contains("C")]

len(df)
```

Figure 20: Data Pre-processing MBA UK Part 2

4 Data Visualization

All of the conclusions obtained from the data are presented here.

4.1 Exploratory Data Analysis

4.1.1 EDA and RFM analysis of UK Data

```
customer_counts = df_new.customer_id.value_counts().sort_values(ascending=False).head(10)

fig, ax = plt.subplots(figsize = (10, 8))

sns.barplot(y = customer_counts.index, x = customer_counts.values, orient = "h",
            ax = ax, order = customer_counts.index, palette = "Reds_r")

plt.title("Customers that have most transactions:Top 10")
plt.ylabel("Customers")
plt.xlabel("Transaction Count")

plt.show()
```

Figure 21: EDA part 1

```
ax = df_new.groupby('invoice_num')['day'].unique().value_counts().sort_index().plot(kind='bar',color=color[0],figsize=(10,5))
ax.set_xlabel('Day',fontsize=15)
ax.set_ylabel('Number of Orders',fontsize=15)
ax.set_title('Number of orders for different Days',fontsize=15)
ax.set_xticklabels(['Mon', 'Tue', 'Wed', 'Thur', 'Fri', 'Sun'], rotation='horizontal', fontsize=15)
plt.show()
```

Figure 22: EDA Part 2

```
group_country_orders = df_new.groupby('country')['invoice_num'].count().sort_values()
# del group_country_orders['United Kingdom']

# plot number of unique customers in each country (with UK)
plt.subplots(figsize=(15,8))
group_country_orders.plot(kind='barh', fontsize=12, color=color[0])
plt.xlabel('Number of Orders', fontsize=12)
plt.ylabel('Country', fontsize=12)
plt.title('Number of Orders for different Countries', fontsize=12)
plt.show()
```

Figure 23: EDA Part 3

4.1.2 EDA and RFM analysis of Brazil Data

```
: other_state_geolocation = geo_df.groupby(['zip_code'])['geolocation_state'].nunique().reset_index(name='count')
other_state_geolocation[other_state_geolocation['count']>= 2].shape
max_state = geo_df.groupby(['zip_code', 'geolocation_state']).size().reset_index(name='count').drop_duplicates(subset = 'zip_code')

: geo_updated = geo_df.groupby(['zip_code', 'geolocation_city', 'geolocation_state'])[['geolocation_lat', 'geolocation_lng']].median()
geo_updated = geo_updated.merge(max_state, on=['zip_code', 'geolocation_state'], how='inner')

: geo_updated = customers_df.merge(geo_updated, left_on='zip_code', right_on='zip_code', how='inner')

: def plot_brazil_map(data):
brazil = plt.imread(urllib.request.urlopen("https://i.pinimg.com/originals/3a/0c/e1/3a0ce10b3c842748c255bc0aa445ad41.jpg").".")
ax = data.plot(kind="scatter", x="geolocation_lng", y="geolocation_lat", figsize=(10,10), alpha=0.3,s=0.3,c='blue')
plt.axis('off')
plt.imshow(brazil, extent=[-73.98283055, -33.8, -33.75116944, 5.4])
plt.show()
```

Figure 24: EDA part 1

```
: sns.countplot(data["review_score"])
data["review_score"].value_counts() / data["review_score"].count() * 100
```

Figure 25: EDA Part 2

```

top_prod_catg = data.groupby('product_category_name_english')['order_id'].nunique().sort_values(ascending=False).head(10)
top_prod_catg

plt.figure(figsize=(15,7))

sns.barplot(y=top_prod_catg.index, x=top_prod_catg, palette="Set2")
plt.xlabel('Number of orders')
plt.title('Top 10 Product categories in terms of number of order ')

```

Figure 26: EDA Part 3

5 Segmentation Technique and MBA techniques on UK and Brazil Data

5.1 Segmentation Technique on Brazil and UK Data

```

# Train the model on 4 clusters
kmean_model = KMeans(n_clusters=4, random_state=5)
kmean_y = kmean_model.fit_predict(RFM_Table_scaled)
# Add labels to df
rfm_df['Cluster'] = kmean_model.labels_

# Function to visualize clusters
def rfm_values(df):
    df_new = df.groupby(['Cluster']).agg({
        'Recency': 'mean',
        'Frequency': 'mean',
        'Monetary': ['mean', 'count']
    }).round(0)
    return df_new

rfm_values(rfm_df)

```

Figure 27: RFM Technique on Brazil Data

```

In [354]: # Train the model on 6 clusters
kmean_model = KMeans(n_clusters=6, random_state=5)
kmean_y = kmean_model.fit_predict(RFM_Table_scaled)
# Add Labels to df
rfm_df['Cluster'] = kmean_model.labels_

In [355]: # Function to visualize clusters
def rfm_values(df):
    rfm_df = df.groupby(['Cluster']).agg({
        'Recency': 'mean',
        'Frequency': 'mean',
        'Monetary': ['mean', 'count']
    }).round(0)
    return rfm_df

In [356]: rfm_values(rfm_df)

```

Figure 28: RFM Technique on UK Data

5.2 Market Basket Analysis on UK and Brazil data

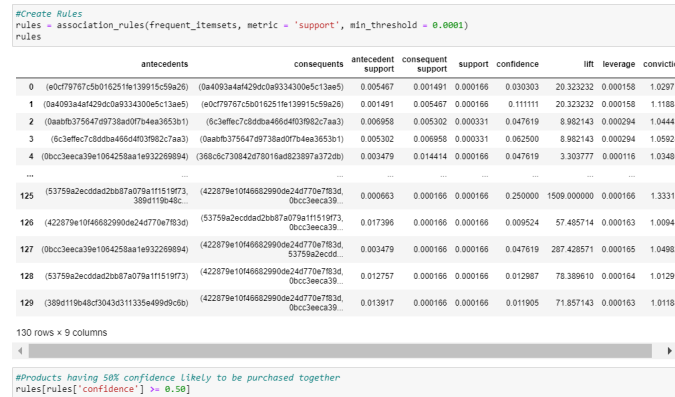


Figure 29: MBA for Brazil

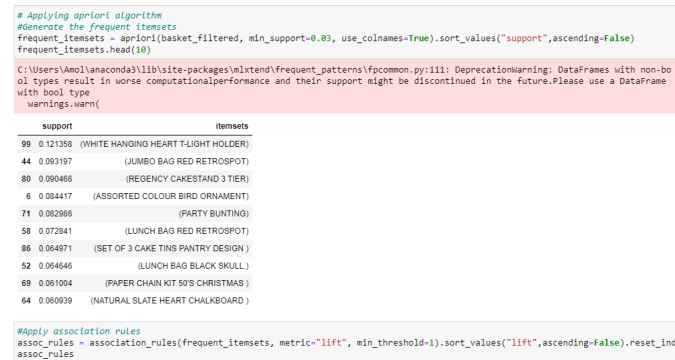


Figure 30: MBA for UK