

Data Science - Times Series

Description

Online shopping is, more than ever, part of people's daily lives. More than ever, using data to make decisions is the big bet of companies in the world.

Understanding the profile of the customer portfolio to generate personalized campaigns is imperative in a competitive world like today.

Furthermore, good sales planning is crucial for the supply sector chain to power the network and avoid problems such as sold out and overstocking.

Let's take on this challenge to work with a purchasing base via e-commerce that allows us to build views about customers and projections about future sales.

Goals

Customer segmentation

Construction of customer segmentation with a focus on value pyramid: Recency, Frequency and Value

Sales projection by segment

- Train a time series forecast model to forecast sales (Sales column) for the segment with the highest sales volume (across the entire history) - daily frequency
- Model validation: December 2014 (daily frequency)
- Use MAPE, SMAPE and RMSE metrics to calculate the result in Dec 2014.

Data Environment

- conda install pandas openpyxl matplotlib seaborn scipy scikit-learn xgboost

Import libraries

```
In [529]:
# Data Analysis Step
import datetime as dt
import matplotlib.pyplot as plt
import pandas as pd
import numpy as np
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')

# plot.parameters
%matplotlib inline
plt.rcParams['font.size'] = 14
plt.rcParams['figure.figsize'] = (9, 5)
plt.rcParams['figure.facecolor'] = '#00000000'

# outliers with mahalanobis
from scipy.spatial import distance

# Segmentation Step
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
from sklearn.cluster import KMeans

# Machine Learning
import xgboost as xgb
from sklearn.ensemble import RandomForestRegressor
from sklearn.metrics import mean_absolute_error, mean_squared_error, mean_squared_log_error
from statsmodels.tools.eval_measures import rmse

In [7]:
# dataset from the Excel file
dataset = pd.read_excel('.../data/dataset_desafio_datascience.xlsx')
df = dataset.copy()
```

Step 1 - Data Preparation and Investigation

Data Preparation

Notes:

- TARGET -> 'Sales', 'customer ID', 'segment', 'city', 'state', 'order date', 'ship date', 'profit'
- order and ship date are not in datetime format.

Understanding the data

```
In [10]:
df.summary(df):
'''
Display a summary of a DataFrame.

This function displays the first few rows of the DataFrame, basic information about the DataFrame's struc
and summary statistics including percentiles.

Parameters:
df (DataFrame): The pandas DataFrame to be summarized.

Returns:
None

display(f'Rows: {df.shape[0]}; Columns: {df.shape[1]}')
display(df.head())
print('-' * 100)
display(df.info())
print('-' * 100)
display(df.describe([0.01, 0.25, 0.50, 0.75, 0.99]))

In [11]:
summary(df)

'Rows: 51290; Columns: 23'
```

	Order ID	Order Date	Ship Date	Ship Mode	customer ID	Customer Name	Segment	City	State	Country	...	Product ID	Cate
0	ES-2011-5338028	28-12-2011	01-01-2012	Second Class	18115	Mick Hernandez	Home Office	Turin	Piedmont	Italy	...	FUR-BO-10003541	Furr
1	ID-2011-24160	29-12-2011	01-01-2012	Second Class	20170	Sarah Bern	Consumer	Gold Coast	Queensland	Australia	...	TEC-PH-10002601	Techn
2	CA-2011-141313	28-12-2011	01-01-2012	Standard Class	10780	Anthony Jacobs	Corporate	Beverly	Massachusetts	United States	...	OFF-AP-10002651	C
3	ID-2011-24160	29-12-2011	01-01-2012	Second Class	20170	Sarah Bern	Consumer	Gold Coast	Queensland	Australia	...	FUR-BO-10002308	Furr
4	CA-2011-104738	30-12-2011	01-01-2012	Second Class	20620	Stefania Perrino	Corporate	Laredo	Texas	United States	...	TEC-PH-10002468	Techn

5 rows * 23 columns

```
-----
<class pandas.core.frame.DataFrame>
RangeIndex: 51290 entries, 0 to 51289
Data columns (total 23 columns):
# Column Non-Null Count Dtype
---
0 Order ID 51290 non-null object
1 Order Date 51290 non-null object
2 Ship Date 51290 non-null object
3 Ship Mode 51290 non-null object
4 customer ID 51290 non-null int64
5 Customer Name 51290 non-null object
6 Segment 51290 non-null object
7 City 51290 non-null object
8 State 51290 non-null object
9 Country 51290 non-null object
10 Postal Code 9994 non-null float64
11 Market 51290 non-null object
12 Region 51290 non-null object
13 Product ID 51290 non-null object
14 Category 51290 non-null object
15 Sub-Category 51290 non-null object
16 Product Name 51290 non-null object
17 Sales 51290 non-null float64
18 Quantity 51290 non-null int64
19 Discount 51290 non-null float64
20 Profit 51290 non-null float64
21 Shipping Cost 51290 non-null float64
22 Order Priority 51290 non-null object
dtypes: float64(5), int64(12), object(16)
memory usage: 9.0+ MB
None
```

	customer ID	Postal Code	Sales	Quantity	Discount	Profit	Shipping Cost
count	51290.000000	9994.000000	51290.000000	51290.000000	51290.000000	51290.000000	51290.000000
mean	14037.790310	55190.379428	246.490581	3.476545	0.142908	28.610982	26.375915
std	5289.254612	32063.693350	487.565361	2.278756	0.212280	174.340972	57.296804

```
min 15.000000 1040.000000 0.444000 1.000000 0.000000 -6599.978000 0.000000
1 600.000000 2149.000000 3.690000 1.000000 0.000000 -351.505650 0.200000
25 11050.000000 23223.000000 30.758625 2.000000 0.000000 0.000000 2.610000
50% 14635.000000 56430.500000 85.053000 3.000000 0.000000 9.240000 7.790000
75% 18265.000000 90008.000000 251.053200 5.000000 0.200000 36.810000 24.450000
99% 21790.000000 98115.000000 2301.000000 11.000000 0.700000 587.359950 286.754300
max 21925.000000 99301.000000 22638.480000 14.000000 0.850000 8399.976000 933.570000
```

Selecting specific columns for analysis

```
In [12]:
df = df[['Order ID', 'Order Date', 'Ship Date', 'customer ID', 'Segment', 'City', 'State', 'Region', 'Country', 'Market', 'Product ID', 'Category']]

In [13]:
df.head()

Out[13]:
   Order ID  Order Date  Ship Date  customer ID  Segment  City  State  Region  Country  Market  Product ID  Category
0  ES-2011-5338028  28-12-2011  01-01-2012  18115  Home Office  Turin  Piedmont  South  Italy  EU  FUR-BO-10003541  Furniture
1  ID-2011-24160  29-12-2011  01-01-2012  20170  Consumer  Gold Coast  Queensland  Oceania  Australia  APAC  TEC-PH-10002601  Technology
2  CA-2011-141313  28-12-2011  01-01-2012  10780  Corporate  Beverly  Massachusetts  East  United States  US  OFF-AP-10002651  Office Supplies
3  ID-2011-24160  29-12-2011  01-01-2012  20170  Consumer  Gold Coast  Queensland  Oceania  Australia  APAC  FUR-BO-10002308  Furniture
4  CA-2011-104738  30-12-2011  01-01-2012  20620  Corporate  Laredo  Texas  Central  United States  US  TEC-PH-10002468  Technology
```

Normalize Date

```
In [14]:
# function, verify if is date format
def is_date_format(df, column):
    data_type = df[column].dtype

    if data_type == 'datetime64[ns]':
        print(f'The '{column}' column is in date format.")
    else:
        print(f'The '{column}' column is not in date format.")

In [15]:
is_date_format(df, 'Order Date')

The 'Order Date' column is not in date format.

In [16]:
is_date_format(df, 'Ship Date')

The 'Ship Date' column is not in date format.

In [17]:
# convert to date format
df['Order Date'] = pd.to_datetime(df['Order Date'], format='%d-%m-%Y', dayfirst=True)

is_date_format(df, 'Order Date')
df['Order Date']

The 'Order Date' column is in date format.

Out[17]:
0 2011-12-28
1 2011-12-29
2 2011-12-28
3 2011-12-29
4 2011-12-30
...
51285 2014-12-26
51286 2014-12-27
51287 2014-12-24
51288 2014-12-24
51289 2014-12-27
Name: Order Date, Length: 51290, dtype: datetime64[ns]

In [18]:
# convert to date format
df['Ship Date'] = pd.to_datetime(df['Ship Date'], format='%d-%m-%Y', dayfirst=True)

is_date_format(df, 'Ship Date')
df['Ship Date']

The 'Ship Date' column is in date format.

Out[18]:
0 2012-01-01
1 2012-01-01
2 2012-01-01
3 2012-01-01
4 2012-01-01
...
51285 2014-12-31
51286 2014-12-31
51287 2014-12-31
51288 2014-12-31
51289 2014-12-31
Name: Ship Date, Length: 51290, dtype: datetime64[ns]
```

Data Cleaning

Handling missing values

```
In [19]:
# Is there null values?
df.isna().sum()

Out[19]:
Order ID      0
Order Date    0
Ship Date     0
customer ID    0
Segment       0
City          0
State         0
Region        0
Country       0
Market        0
Product ID    0
Category      0
Product Name  0
Sales         0
Quantity      0
Profit        0
dtype: int64
```

Lucky! I didn't need to remove null values as they are only present in the 'Postal Code' column, which will be used in the analysis.

Detecting and Removing Outliers

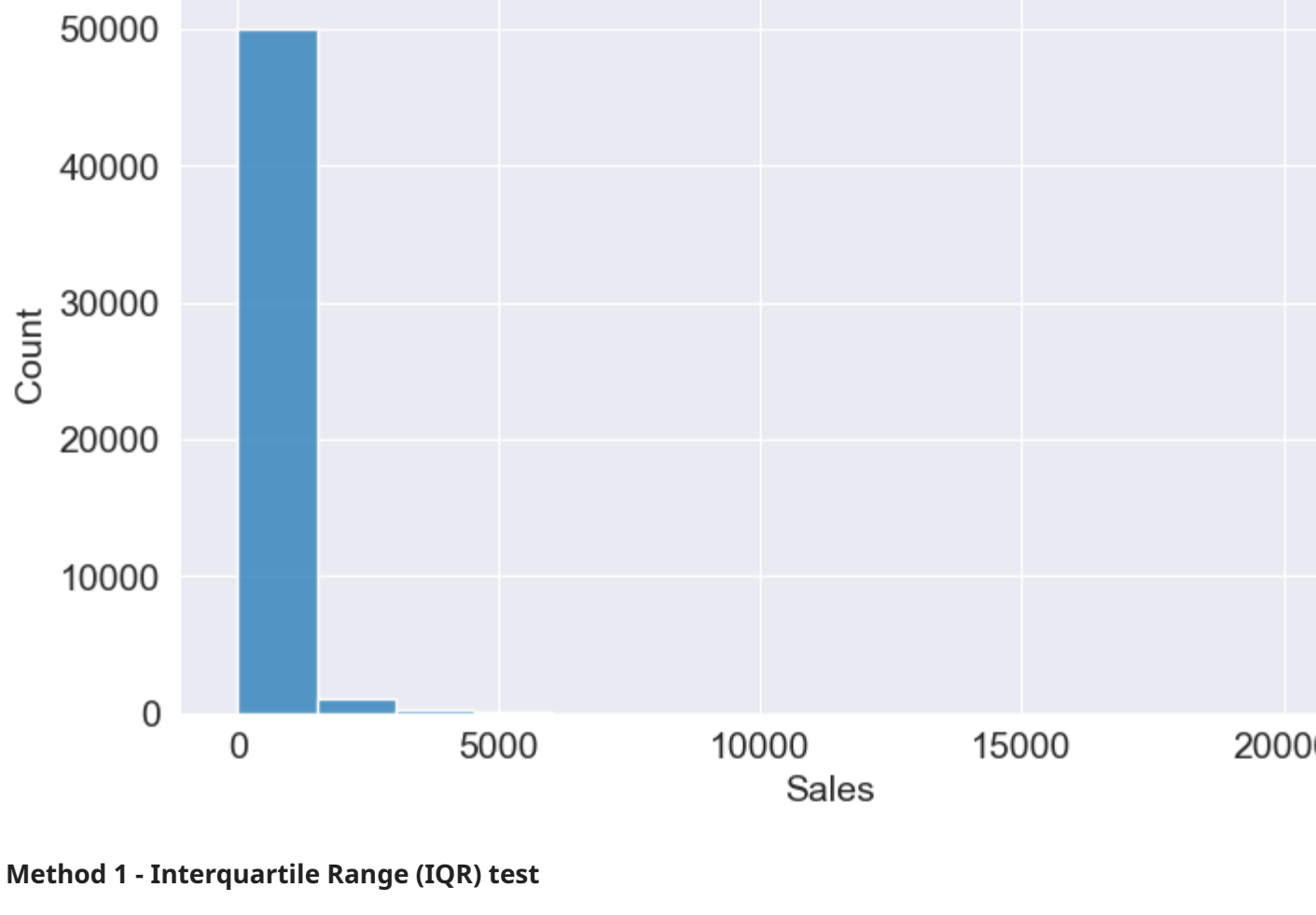
Box-plot

A box plot shows the distribution of the data and identify as any outliers as points that fall outside the "whiskers" of the plot.

```
In [25]:
# Box-plot function
def plot_boxplot(df, column):
    sns.boxplot(x = df[column])
    plt.show()

# Hist function
def plot_hist(df, column):
    sns.histplot(df[column], bins=15, kde=False)
    sns.set_style('darkgrid')
    plt.show()

In [21]:
plot_boxplot(df, 'Sales')
```



```
In [26]:
plot_hist(df, 'Sales')
```

Method 1 - Interquartile Range (IQR) test

The interquartile range test is a statistical method used to identify potential outliers in a dataset. It involves calculating the difference between the first quartile (Q1) and the third quartile (Q3) of the data distribution, which represents the middle 50% of the data.

The IQR is then multiplied by a constant factor, typically 1.5, to define the lower and upper limits of the expected range of data values. Data points that fall outside of this expected range are considered potential outliers.

To apply the interquartile test, we first sort the data in ascending order and then calculate the values of Q1 and Q3. These can be calculated as the median of the lower half of the data and the median of the upper half of the data, respectively. The IQR is then calculated as Q3 minus Q1.

The lower limit of the expected range is defined as Q1 minus 1.5 times the IQR, and the upper limit is defined as Q3 plus 1.5 times the IQR. Any data points that fall outside of these limits are considered potential outliers and should be further examined and verified.

One advantage of the interquartile test is that it is less sensitive to extreme values than the next test, the Z-score test, which makes it more suitable for datasets with skewed or non-normal distributions. However, it may not be as effective in detecting outliers in small sample sizes or datasets with multiple modes.

```
In [27]:
# Creating the series value with the target columns SALES
value = df['Sales']

# first quartile
q1 = value.quantile(.25)

# third quartile
q3 = value.quantile(.75)

# interquartile range
IIO = q3 - q1

# inferior limit
lower_limit = q1 - 1.5 * IIO

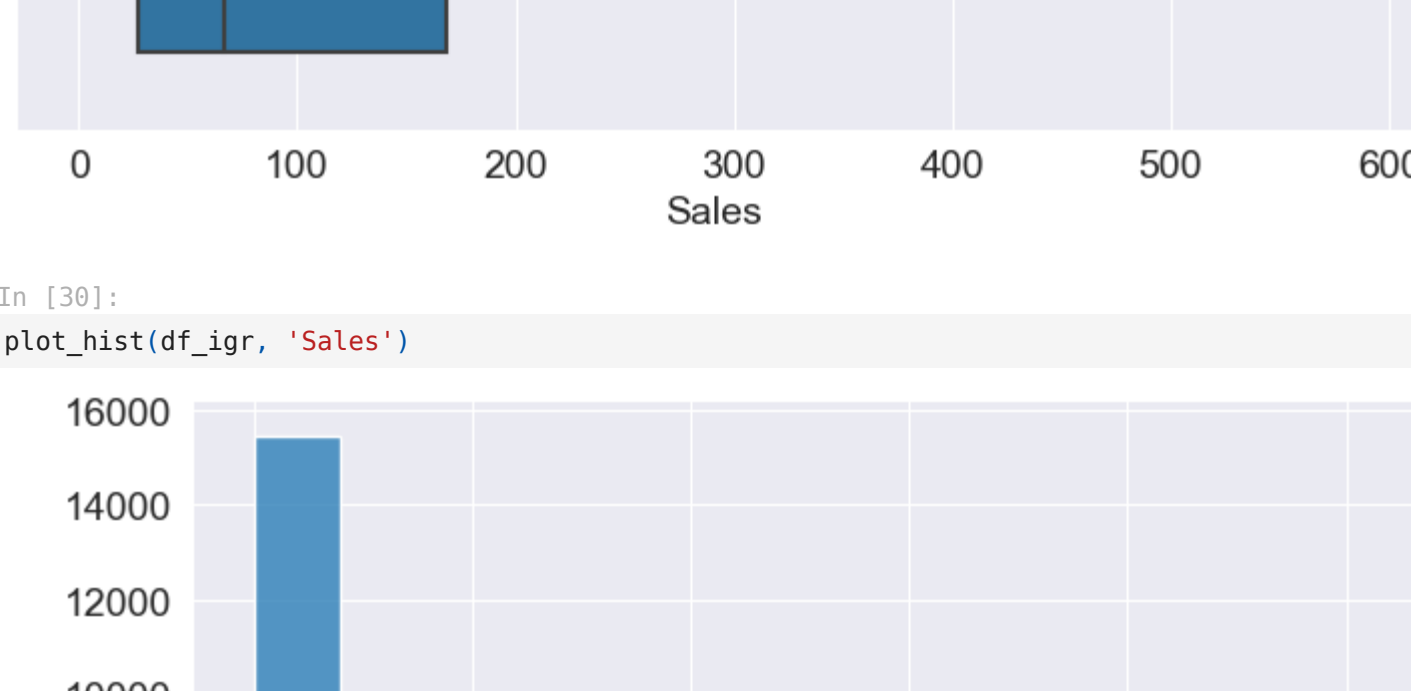
# upper limit
upper_limit = q3 + 1.5 * IIO

print(f'''
[01] -> {q1}
[03] -> {q3}
[IIO] -> {IIO:.2f}
[01 - 1.5 * IIO] -> {lower_limit:.2f}
[03 + 1.5 * IIO] -> {upper_limit:.2f}
''')

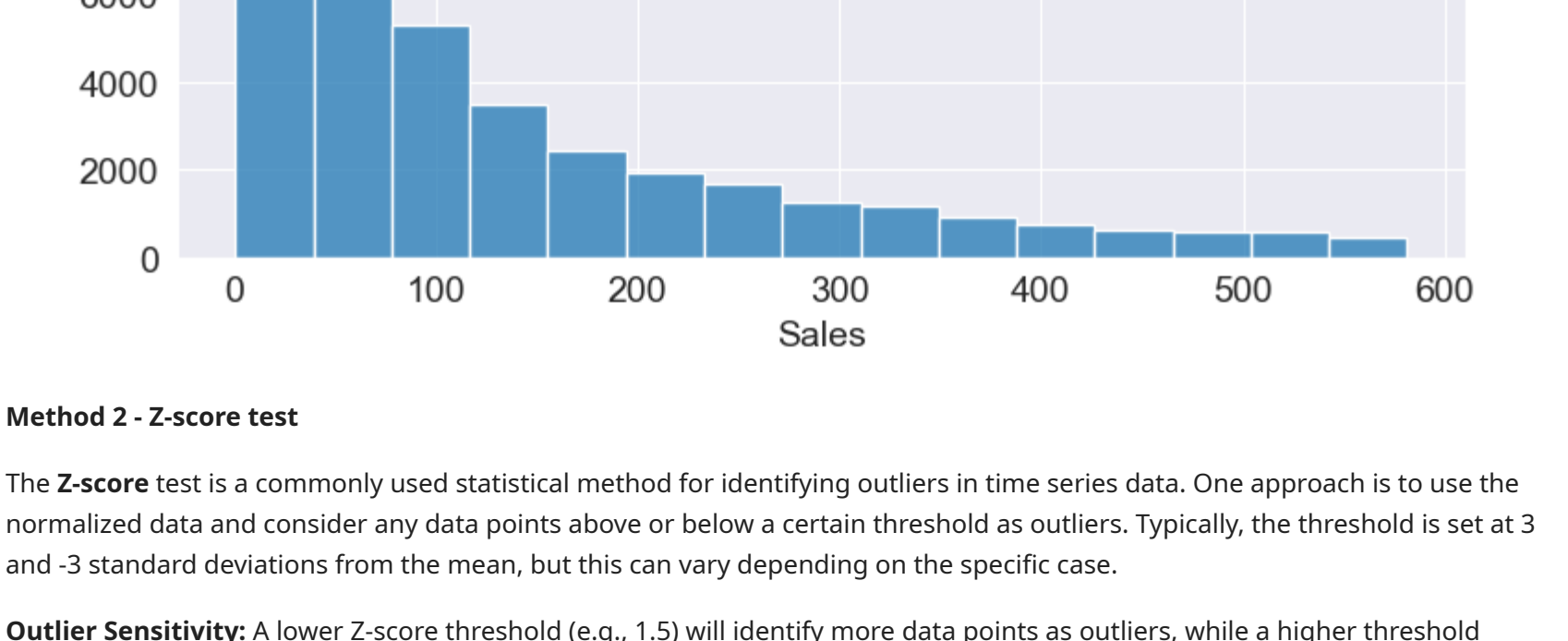
[01] -> 30.758625000000002
[03] -> 251.0532
[01 - 1.5 * IIO] -> -228.29
[03 + 1.5 * IIO] -> 299.68
[03 + 1.5 * IIO] -> 581.58

In [28]:
# applying selection = (value <= lower_limit) & (value <= upper_limit)
df_igr = df[selection]

In [29]:
plot_boxplot(df_igr, 'Sales')
```



```
In [30]:
plot_hist(df_igr, 'Sales')
```



Method 2 - Z-score test

The Z-score test is a commonly used statistical method for identifying outliers in time series data. One approach is to use the normalized data and consider any data points above or below a certain threshold as outliers. Typically, the threshold is set at 3 and -3 standard deviations from the mean, but this can vary depending on the specific case.

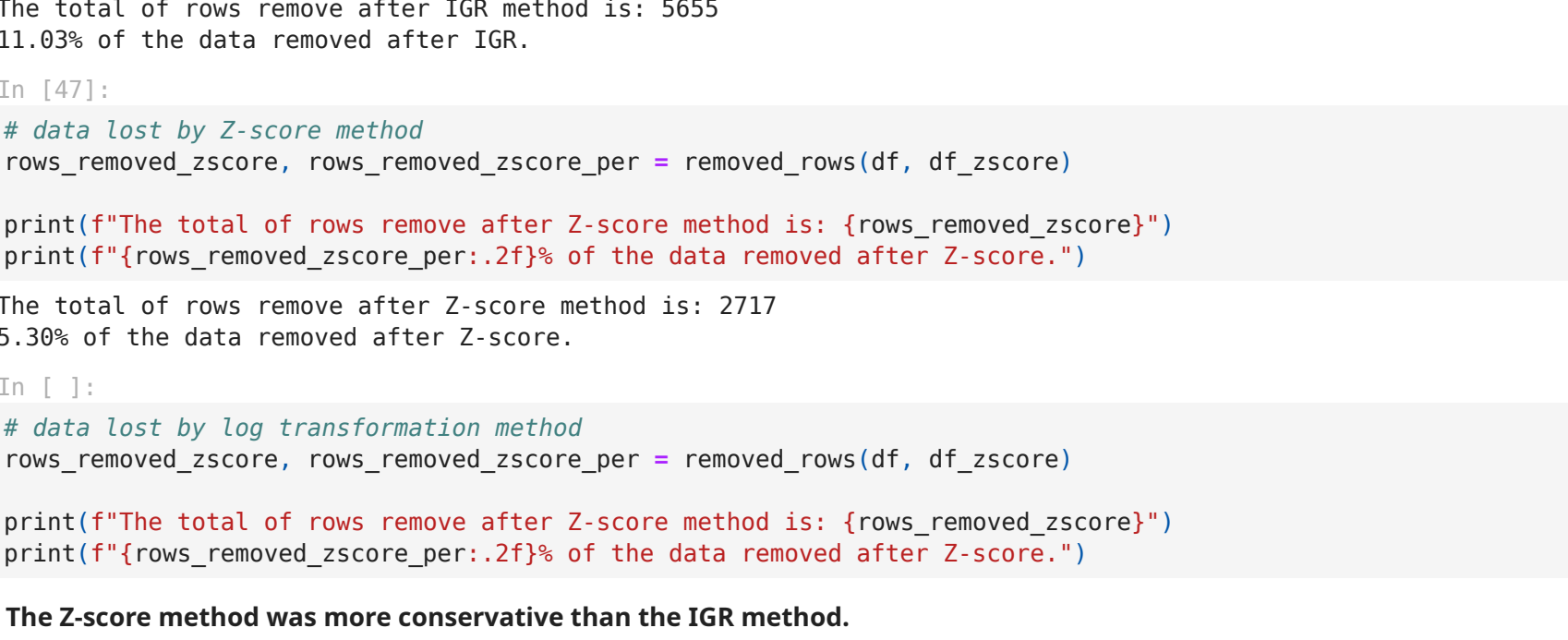
Outlier Sensitivity: A lower Z-score threshold (e.g., 1.5) will identify more data points as outliers, while a higher threshold (e.g., 3) will be more conservative in identifying outliers. The threshold choice should align with your sensitivity to outliers.

```
In [45]:
# Z-scores for the 'Sales' column
z_scores = np.abs(stats.zscore(df['Sales']))

# define a z-score threshold to identify outliers
z_score_threshold = 1.5

# new df with outliers removed using the z-score
df_zscore = df[z_scores < z_score_threshold]

In [46]:
plot_hist(df_zscore, 'Sales')
```



How much of the data we lost after remove outliers?

```
In [34]:
lost_by_igr = df_igr.shape[0]
rows_removed_igr = int(df.shape[0]) - lost_by_igr

print(f'The total of rows remove after IGR method is: {rows_removed_igr}')
print(f'The total of rows removed after IGR method is: {rows_removed_igr} / {df.shape[0]} = 11.03% of the data removed after IGR.')

In [47]:
# data lost by Z-score method
rows_removed_zscore, rows_removed_zscore_per = removed_rows(df, df_zscore)

print(f'The total of rows remove after Z-score method is: {rows_removed_zscore}')
print(f'The total of rows removed after Z-score method is: {rows_removed_zscore} / {df.shape[0]} = 5.30% of the data removed after Z-score.')

The total of rows removed after IGR method is: 5655
11.03% of the data removed after IGR.

The total of rows removed after Z-score method is: 2717
5.30% of the data removed after Z-score.
```

The Z-score method was more conservative than the IGR method.

Data Investigation - EDA

Collect Large Numbers Describing the Customer Base and E-commerce Segments

```
In [51]:
df_igr.head()

Out[51]:
   Order ID  Order Date  Ship Date  customer ID  Segment  City  State  Region  Country  Market  Product ID  Category
1  2011-12-29  01-01-2012  20170  Consumer  Gold Coast  Queensland  Oceania  Australia  APAC  TEC-PH-10002601  Technology  Sma
3  2011-12-29  01-01-2012  20170  Consumer  Gold Coast  Queensland  Oceania  Australia  APAC  FUR-BO-10002308  Furniture  Floa
4  2011-12-30  01-01-2012  20620  Corporate  Laredo  Texas  Central  United States  US  TEC-PH-10002468  Technology  PI
104738  2011-12-30  01-01-2012  20620  Corporate  Laredo  Texas  Central  United States  US  TEC-PH-10000576  Technology  A
```

Specific Requests

1-Average Delivery Time in All Countries

```
In [52]:
# Calculate 'Delivery Time' as the difference in days using .loc
df_igr.loc[:, 'Delivery Time'] = (df_igr['Ship Date'] - df_igr['Order Date']).dt.days

# Calculate the average delivery time
average_delivery_time_by_country = df_igr.groupby('Country')['Delivery Time'].mean()
print(f'Average Delivery Time in All Countries: {round(average_delivery_time)} days')

Average Delivery Time in All Countries: 4 days

In [53]:
# delivery time per country
average_delivery_time_by_country = df_igr.groupby('Country')['Delivery Time'].mean().reset_index().round()
average_delivery_time_by_country['Delivery Time'].unique()

Out[53]:
array([4., 5., 3., 2.])
```

2 - Most Profitable Customer Segment Each Year

```
In [54]:
# Extract year from Order Date
df_igr.loc[:, 'Year'] = df_igr['Order Date'].dt.year

most_profitable_segment = df_igr.groupby(['Year', 'Segment'])['Profit'].max().reset_index()

most_profitable_segment

Out[54]:
   Year  Segment  Profit
0  2011  Consumer  263.040
1  2011  Corporate  243.300
2  2011  Home Office  274.995
3  2012  Consumer  279.360
4  2012  Corporate  284.220
5  2012  Home Office  274.995
6  2013  Consumer  284.220
7  2013  Corporate  248.040
8  2013  Home Office  270.480
9  2014  Consumer  277.560
10 2014  Corporate  284.220
11 2014  Home Office  270.720
```

3-Segment with the Highest Sales Volume

```
In [55]:
# barplot function for visualization
def bar_plot(x, y, data, title, xlabel, ylabel):
    plt.figure(figsize=(10, 6))
    sns.barplot(x=x, y=y, data=data, palette='viridis')
    plt.title(title, fontsize=16)
    plt.xlabel(xlabel, fontsize=12)
    plt.ylabel(ylabel, fontsize=12)
    plt.xticks(rotation=45) # Rotate x-axis labels for better visibility
    plt.show()

sales_count_by_segment = df_igr.groupby('Segment')['Sales'].count().reset_index()
bar_plot('Segment', 'Sales', sales_count_by_segment, 'Sales Volume by Segment', 'Segment', 'Number of Sales')
```



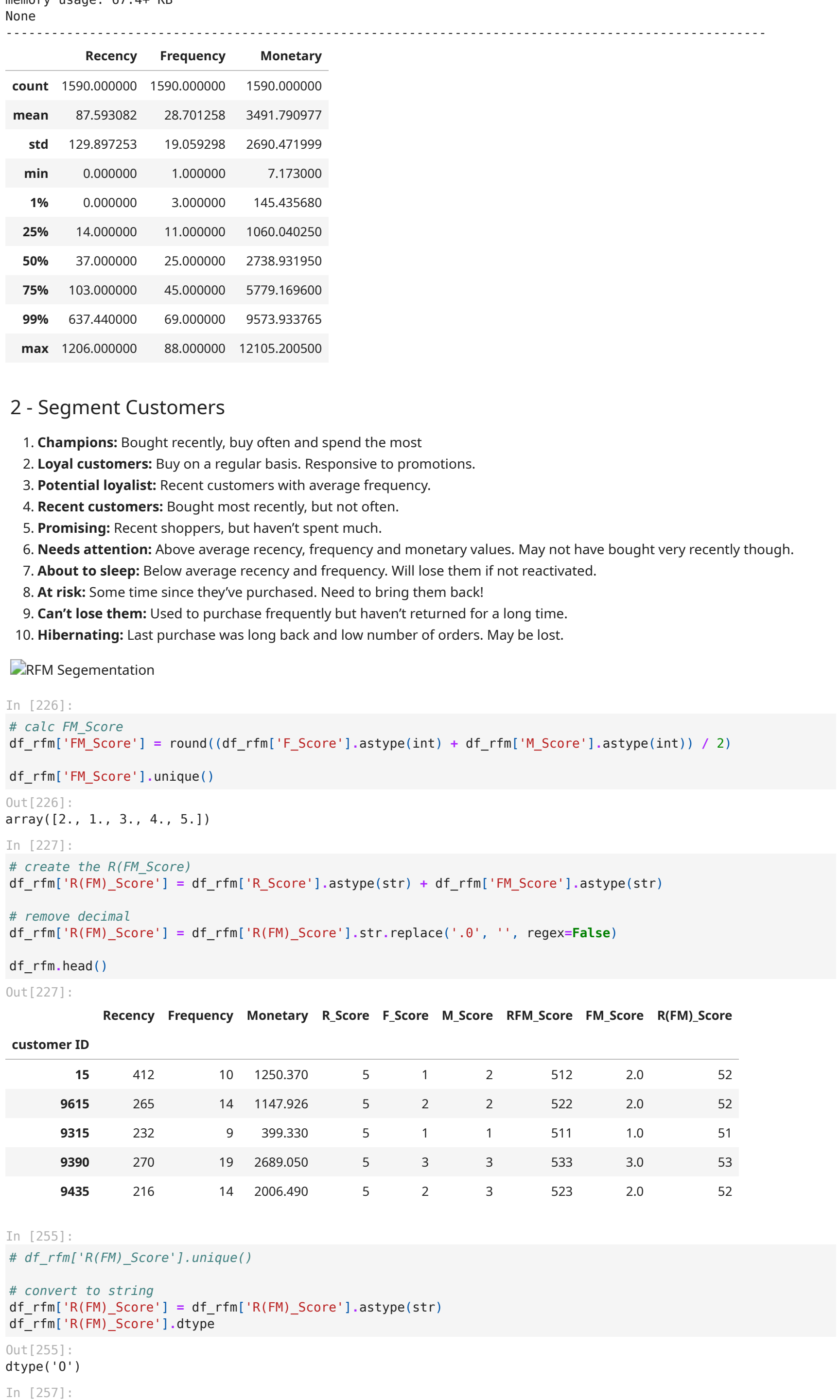
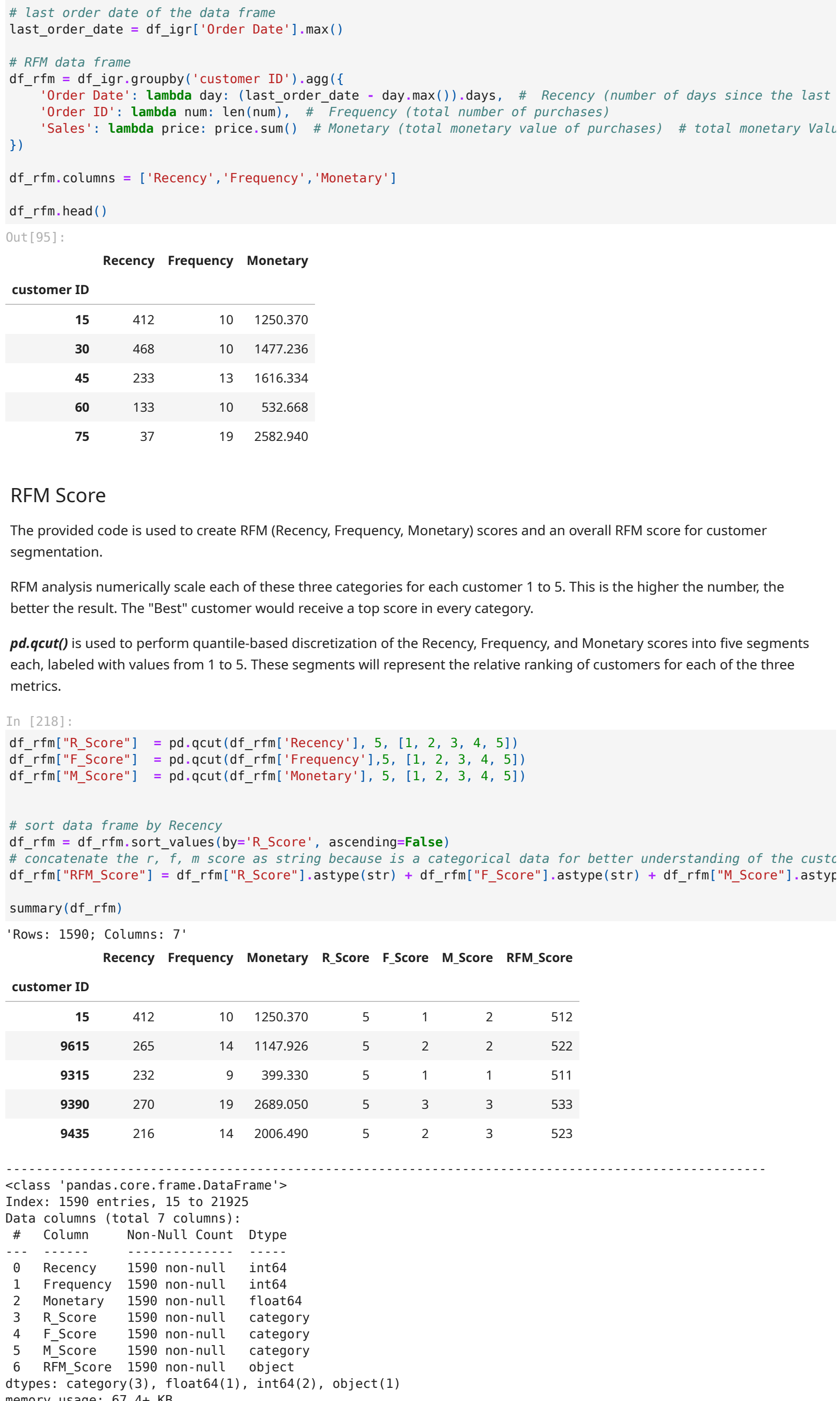
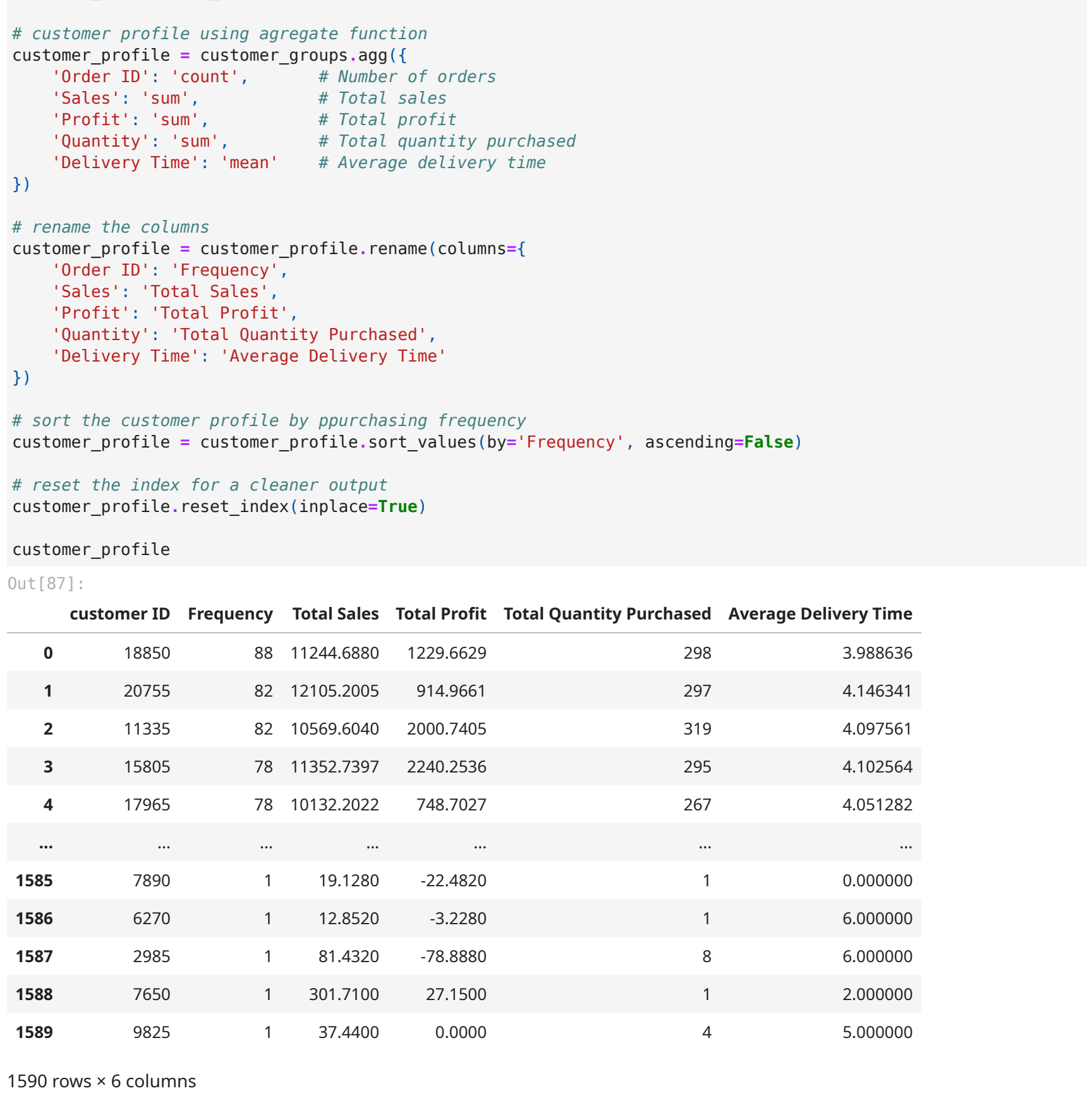
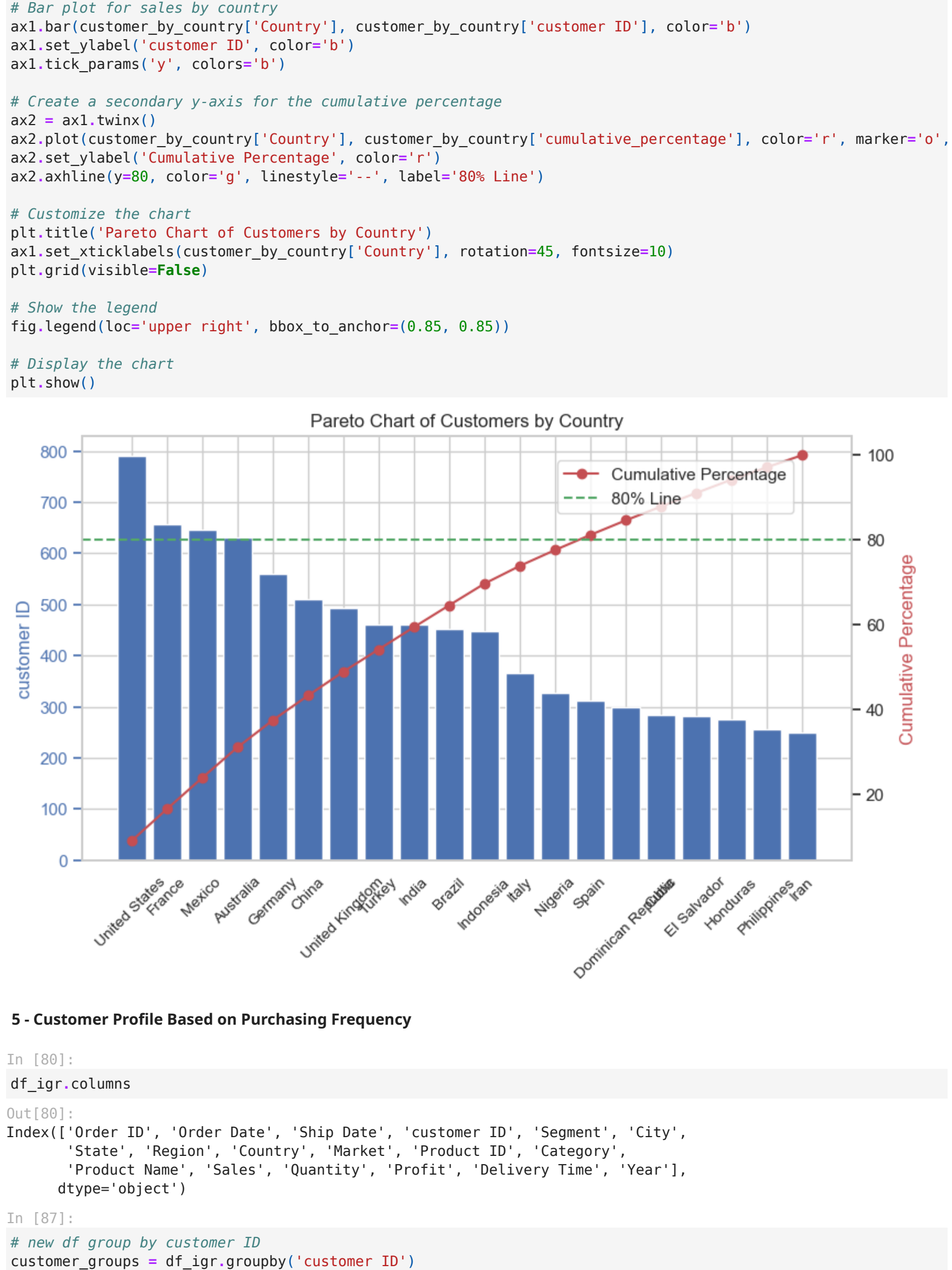
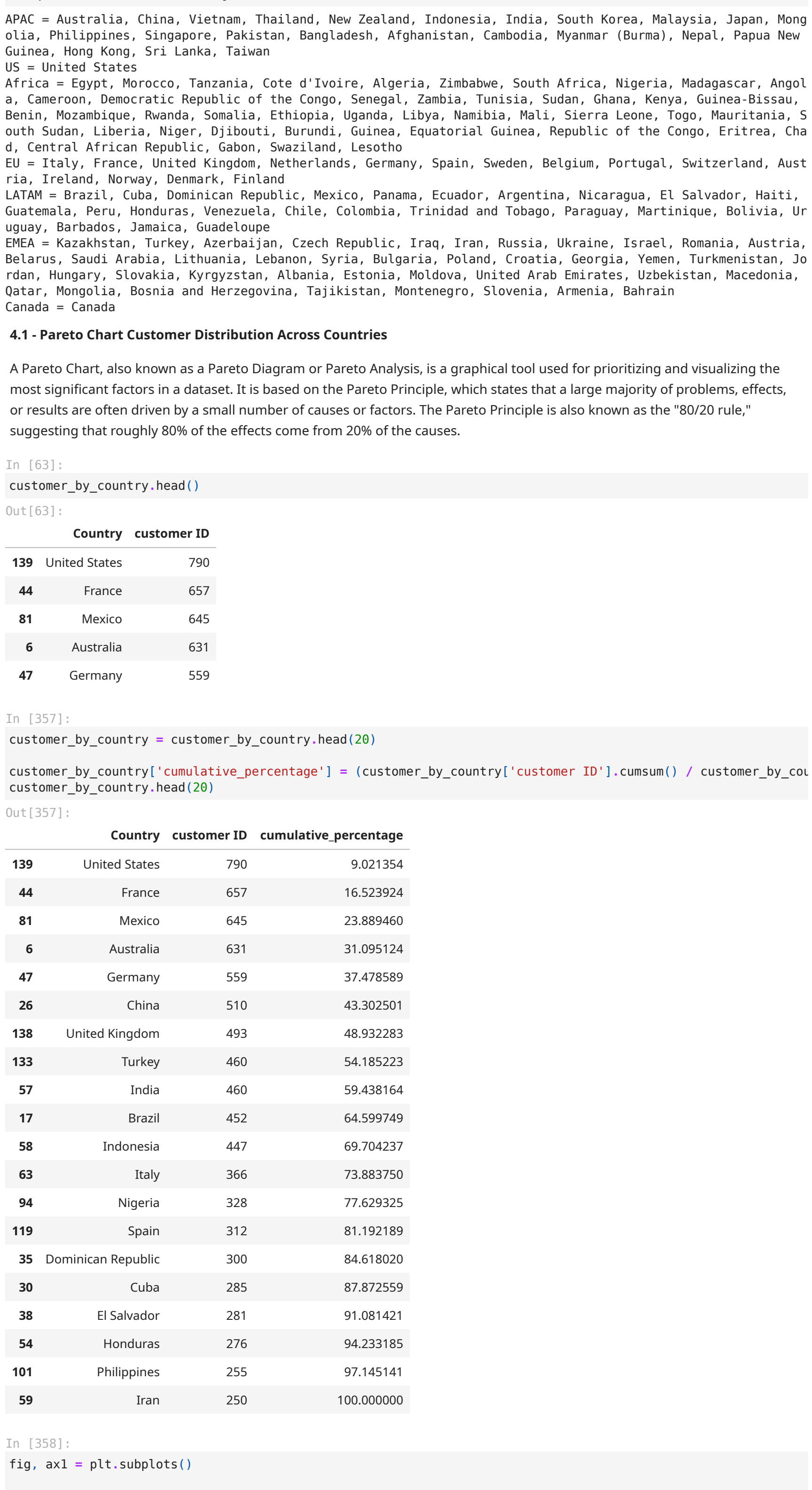
4-Customer Distribution Across Countries

```
In [353]:
customer_by_country = df_igr.groupby('Country')['customer ID'].nunique().reset_index()
customer_by_country = customer_by_country.sort_values(by='customer ID', ascending=False)

In [349]:
customer_by_country.head(20)

Out[349]:
   Country  customer ID
139  United States      790
44  France            657
8  Mexico            645
61  Australia        631
47  Germany          559
26  People's China    510
138  United Kingdom  493
57  Turkey           460
13  India            460
17  Brazil           452
58  Indonesia        447
63  Italy            366
94  Nigeria          328
119 Spain            312
35  Dominican Republic 300
30  Cuba             285
38  El Salvador       281
54  Honduras          276
101 Philippines       255
109 Iran              250
```

```
In [350]:
customer_by_country.tail()
```

Step 2 - Customer Segmentation

1 - Calculate Recency, Frequency, and Value:

RFM (Recency, Frequency, Monetary Value) analysis is a customer segmentation technique that uses past purchase behavior to divide customers into groups. RFM helps divide customers into various categories or clusters to identify customers who are more likely to respond to promotions and also for future personalization services.

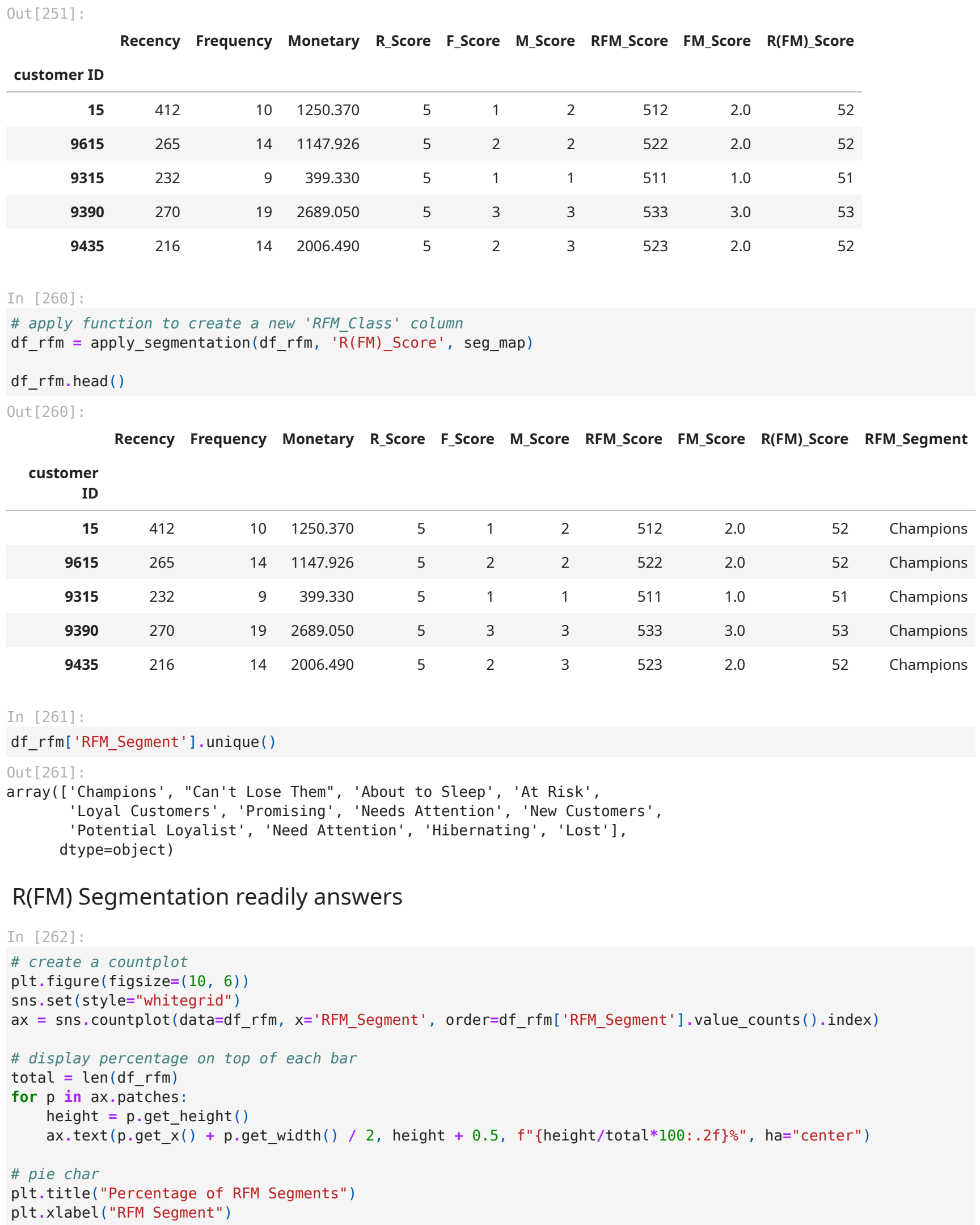
RECENTY (R): Days since last purchase

FREQUENCY (F): Total number of purchases

MONETARY VALUE (M): Total money this customer spent.

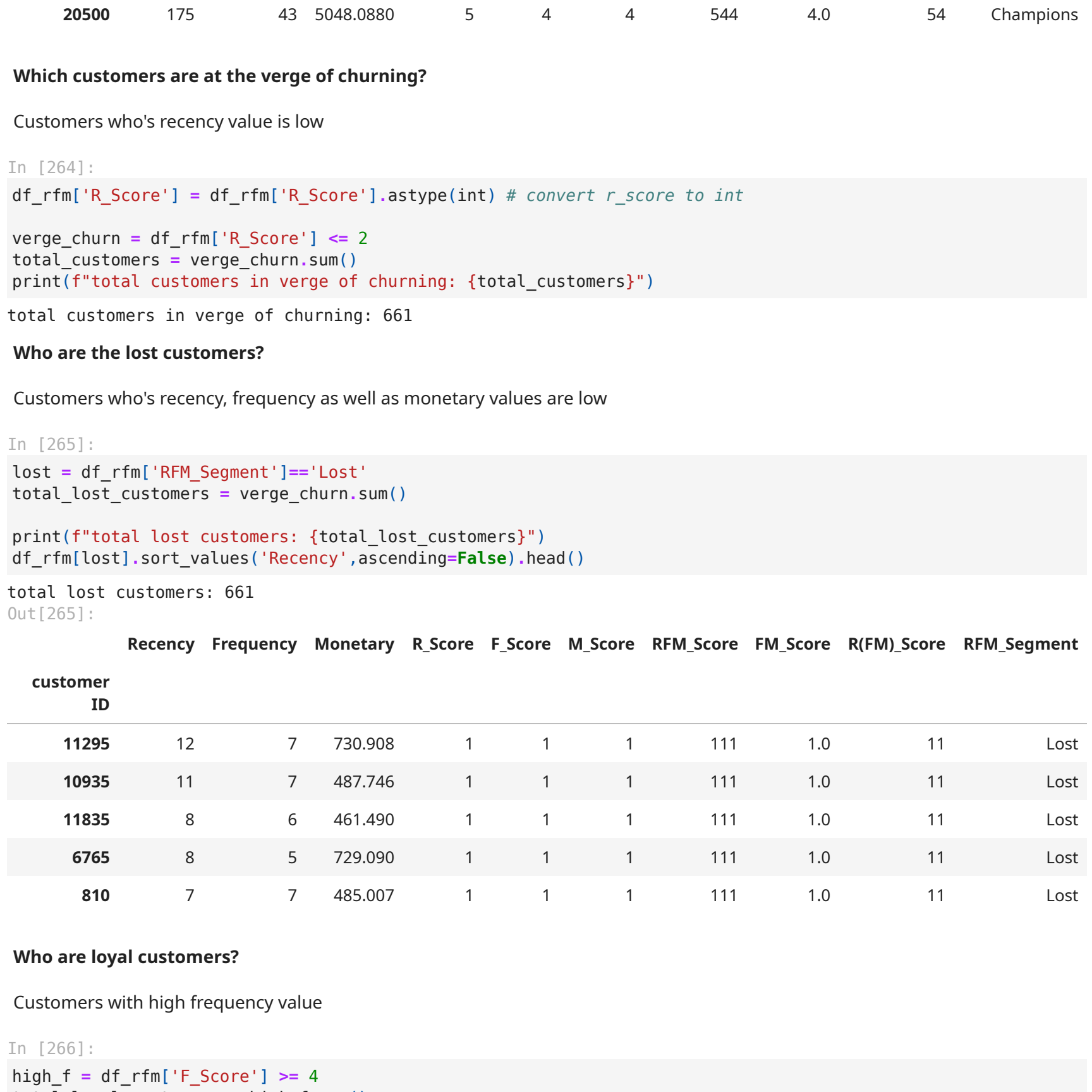
RFM Analysis

RFM Table

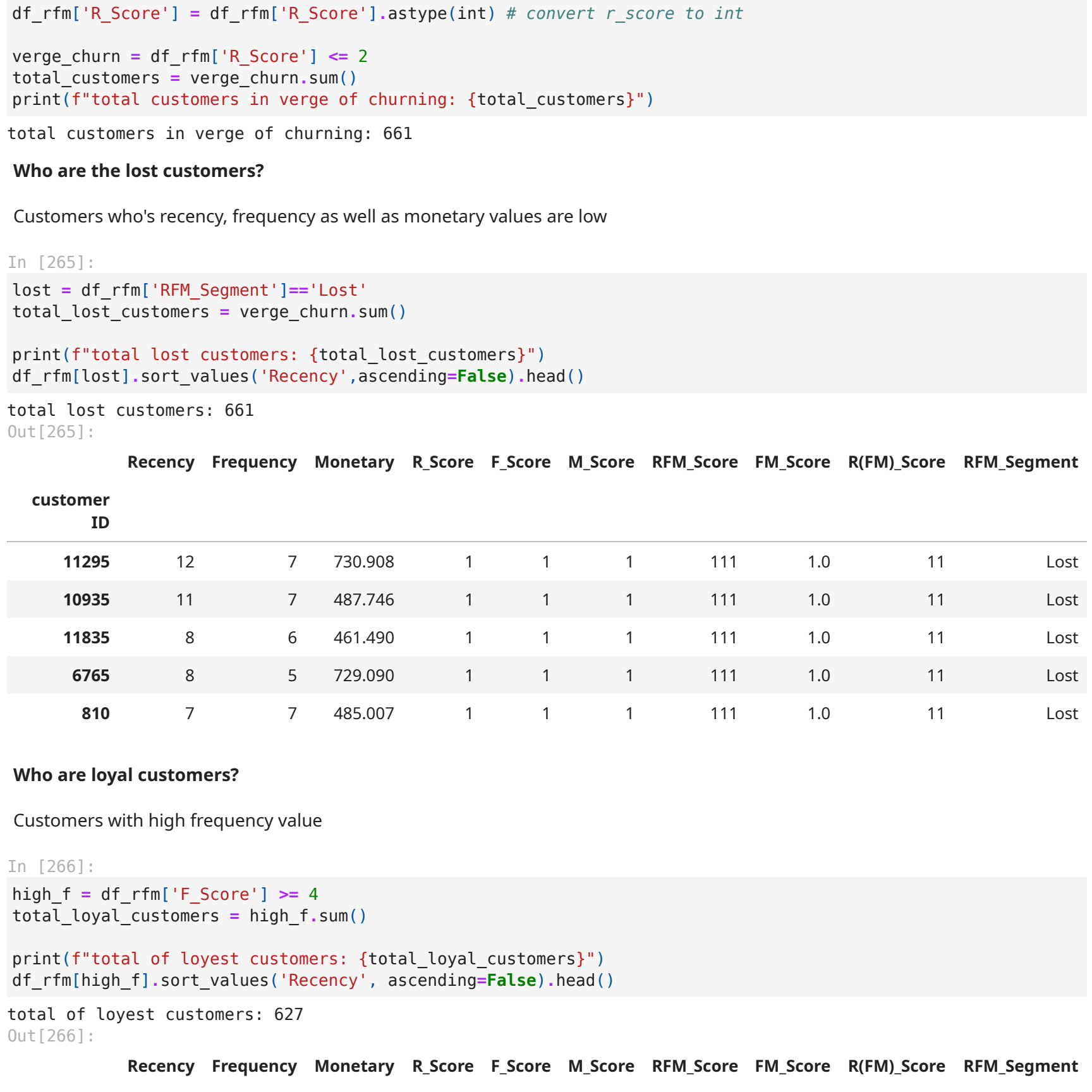


2. Segments Customers

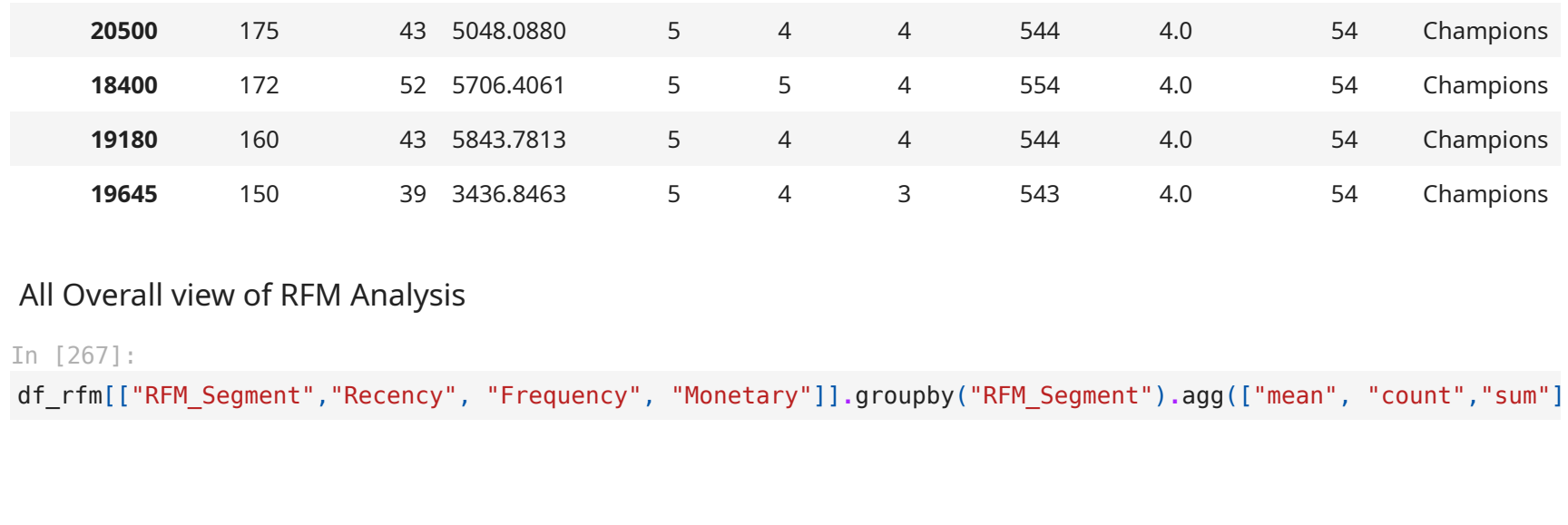
- Champions:** Bought recently, buy often and spend the most
- Loyal customers:** Buy on a regular basis. Responsive to promotions.
- Potential loyalist:** Recent customers with average frequency.
- Recent customers:** Bought most recently, but not often.
- Promising:** Recent shoppers, but haven't spent much.
- Needs attention:** Above average recency, frequency and monetary values. May not have bought very recently though.
- About to sleep:** Below average recency and frequency. Will lose them if not reactivated.
- At risk:** Some time since they've purchased. Need to bring them back for a long time.
- Can't lose them:** Used to purchase frequently but haven't returned for a while.
- Hibernating:** Last purchase was long back and low number of orders. May be lost.



RFM Segmentation ready answers



Who are my best customers?(Champions)



Which customers are at the verge of churning?

Customers who's recency value is low

Who are the lost customers?

Customers who's recency, frequency as well as monetary values are low

Who are loyal customers?

Customers with high frequency value

All Overall view of RFM Analysis

Out[267]:

RFM Segment	Recency			Frequency			Monetary		
	mean	count	sum	mean	count	sum	mean	count	sum
About to Sleep	55.947137	227	12700	16.110132	227	3657	1676.953625	227	3.806685e+05
At Risk	62.786325	117	7346	13.897436	117	1626	1346.521212	117	1.575430e+05
Can't Lose Them	41.927778	180	7547	12.096111	180	7733	5405.154811	180	9.729779e+05
Champions	293.741325	317	93116	11.000000	317	3487	1112.850698	317	3.527773e+05
Hibernating	5.508475	59	325	13.186441	59	778	1311.485203	59	7.737786e+04
Lost	5.600000	15	84	6.666667	15	100	532.353600	15	7.985304e+03
Loyal Customers	35.813725	102	3653	57.519608	102	5867	7698.135336	102	7.852098e+04
Need Attention	19.241935	62	1193	13.048387	62	809	1228.157419	62	7.614576e+04
Needs Attention	40.466667	30	1214	6.966667	30	209	452.687300	30	1.358062e+04
New Customers	17.073529	136	2322	58.000000	136	7888	7820.916243	136	1.063645e+06
Potential Loyalist	20.140000	50	1007	27.760000	50	1388	3122.507861	50	1.561254e+05
Promising	29.715254	295	8766	40.993220	295	12093	5111.746810	295	1.507956e+06

Step 3 - Sales projection by segment

still working

Train a time series forecast model to forecast sales (Sales column) for the segment with the highest sales volume (across the entire

- history (daily frequency)
- Model validation: December 2014 (daily frequency)
- Use MAPE, SMAPE and RMSE metrics to calculate the result in Dec 2014.

Duvida, Segment aqui é referente a coluna do dado original ou Segmentos relacionado a analise RFM?

Mantive o forecasting na coluna Segment do dado original.

Machine Learning Model

Traditional models like ARIMA are more time spending and require more knowledge about the data.

Machine learning models are easier and faster to apply.

Algorithm:

RandomForestRegressor for time series regression tasks such as sales forecasting to predict numeric values.

XGBoost

Data Frame Pre-processing

```
In [269]:
# create a df
df_ml = df_sgr.copy()
df_ml.head()
```

Out[269]:

	Order ID	Order Date	Ship Date	customer ID	Segment	City	State	Region	Country	Market	Product ID	Category
1	ID-2011-24160	2011-12-29	2012-01-01	20170	Consumer	Gold Coast	Queensland	Oceania	Australia	APAC	TEC-PH-10002601	Technology
3	ID-2011-24160	2011-12-29	2012-01-01	20170	Consumer	Gold Coast	Queensland	Oceania	Australia	APAC	FUR-BO-10002308	Furniture
4	CA-2011-104738	2011-12-30	2012-01-01	20620	Corporate	Laredo	Texas	Central	United States	US	TEC-PH-10002468	Technology
5	EG-2011-590	2011-12-28	2012-01-01	5370	Consumer	Sohag	Suhaj	Africa	Egypt	Africa	OFF-SME-10003752	Office Supplies
6	CA-2011-104738	2011-12-30	2012-01-01	20620	Corporate	Laredo	Texas	Central	United States	US	TEC-PH-10000576	Technology

```
In [299]:
# segment with the highest sales
highest_sales_segment = df_ml.groupby('Segment')['Sales'].sum().idxmax()
df_ml_segment = df_ml[df_ml['Segment'] == highest_sales_segment]
df_ml_segment = df_ml_segment.sort_values(by='Segment')
df_ml_segment.head()
```

Out[299]:

	Order ID	Order Date	Ship Date	customer ID	Segment	City	State	Region	Country	Market	Product ID	Cate
1	ID-2011-24160	2011-12-29	2012-01-01	20170	Consumer	Gold Coast	Queensland	Oceania	Australia	APAC	TEC-PH-10002601	Techn
34096	IN-2014-82652	2014-07-16	2014-07-21	11140	Consumer	Wollongong	New South Wales	Oceania	Australia	APAC	OFF-SU-10001488	C
34095	IN-2014-82652	2014-07-16	2014-07-21	11140	Consumer	Wollongong	New South Wales	Oceania	Australia	APAC	OFF-BI-10004562	C
34093	IN-2014-82652	2014-07-16	2014-07-21	11140	Consumer	Wollongong	New South Wales	Oceania	Australia	APAC	OFF-LA-10002086	C
34087	CG-2014-870	2014-07-19	2014-07-21	705	Consumer	Kinshasa	Kinshasa	Africa	Democratic Republic of the Congo	Africa	OFF-BIC-10001211	C

```
In [301]:
# reset index for resample('D')
df_ml_segment.set_index('Order Date', inplace=True)
df_ml_segment.head()
```

Out[301]:

	Order ID	Ship Date	customer ID	Segment	City	State	Region	Country	Market	Product ID	Category
2011-12-29	ID-2011-24160	2012-01-01	20170	Consumer	Gold Coast	Queensland	Oceania	Australia	APAC	TEC-PH-10002601	Technology
2014-07-16	IN-2014-82652	2014-07-21	11140	Consumer	Wollongong	New South Wales	Oceania	Australia	APAC	OFF-SU-10001488	Office Supplies
2014-07-16	IN-2014-82652	2014-07-21	11140	Consumer	Wollongong	New South Wales	Oceania	Australia	APAC	OFF-BI-10004562	Office Supplies
2014-07-16	IN-2014-82652	2014-07-21	11140	Consumer	Wollongong	New South Wales	Oceania	Australia	APAC	OFF-LA-10002086	Office Supplies
2014-07-19	CG-2014-870	2014-07-21	705	Consumer	Kinshasa	Kinshasa	Africa	Democratic Republic of the Congo	Africa	OFF-BIC-10001211	Office Supplies

```
In [301]:
# resample on daily frequency and calculate daily sales
daily_sales_df = df_ml_segment['Sales'].resample('D').sum()
daily_sales_df = daily_sales_df[daily_sales_df.index < '2014-12-01']
daily_sales_df.head()
```

Out[301]:

	Date	Sales
0	2011-01-01	763.69800
1	2011-01-02	314.22000
2	2011-01-03	1465.6262
3	2011-01-04	353.92800
4	2011-01-05	48.78000

```
In [308]:
daily_sales_df.set_index('Date', inplace=True)
daily_sales_df.head()
```

Out[308]:

	Date	Sales
2011-01-01	2011-01-01	763.69800
2011-01-02	2011-01-02	314.22000
2011-01-03	2011-01-03	1465.6262
2011-01-04	2011-01-04	353.92800
2011-01-05	2011-01-05	48.78000

```
In [310]:
# training and testing data sets for 2014
# as required the testing data for December of 2014, train for the rest of 2014
train_data = daily_sales_df[daily_sales_df.index < '2014-12-01']
test_data = daily_sales_df[daily_sales_df.index >= '2014-12-01']

train_data.shape, test_data.shape
```

Out[310]:

((1430, 1), (30, 1))

```
In [312]:
test_data # only for december
```

Out[312]:

	Date	Sales
2014-12-01	2014-12-01	4536.16100
2014-12-02	2014-12-02	8693.51350
2014-12-03	2014-12-03	4970.56520
2014-12-04	2014-12-04	4941.87410
2014-12-05	2014-12-05	3661.41000
2014-12-06	2014-12-06	1401.27200
2014-12-07	2014-12-07	828.89000
2014-12-08	2014-12-08	4660.00010
2014-12-09	2014-12-09	5801.84480
2014-12-10	2014-12-10	5412.58660
2014-12-11	2014-12-11	5189.48248
2014-12-12	2014-12-12	4830.93040
2014-12-13	2014-12-13	3405.40050
2014-12-14	2014-12-14	991.62400
2014-12-15	2014-12-15	5864.21100
2014-12-16	2014-12-16	4550.02600
2014-12-17	2014-12-17	2196.65710
2014-12-18	2014-12-18	5179.34300
2014-12-19	2014-12-19	4143.44720
2014-12-20	2014-12-20	1104.91520
2014-12-21	2014-12-21	17.76000
2014-12-22	2014-12-22	3538.54120
2014-12-23	2014-12-23	5342.80070
2014-12-24	2014-12-24	4026.46700
2014-12-25	2014-12-25	4306.27640
2014-12-26	2014-12-26	7106.10940
2014-12-27	2014-12-27	4138.38976
2014-12-28	2014-12-28	545.20800
2014-12-29	2014-12-29	8123.24728
2014-12-30	2014-12-30	7813.22930
2014-12-31	2014-12-31	2625.78100

Creating the target variable and the lagged variable (Sales)

Daily prediction: complex. The idea is to make the sale value for the previous day and predict it for the next day?

```
In [314]:
train_data['Sales']
```

Out[314]:

Date	Sales
2011-01-01	763.69800
2011-01-02	314.22000
2011-01-03	1465.62620
2011-01-04	353.92800
2011-01-05	48.78000
...	...
2014-11-26	6991.84824
2014-11-27	5452.68672
2014-11-28	6714.80742
2014-11-29	1500.29100
2014-11-30	1055.88800

```
Name: Sales, Length: 1430, dtype: float64
```

```
In [315]:
train_data['Sales'].shift(-1)
train_data['Sales'] # lagged?
```

Out[315]:

Date	Sales
2011-01-01	763.69800
2011-01-02	314.22000
2011-01-03	1465.62620
2011-01-04	353.92800
2011-01-05	48.78000
...	...
2014-11-26	6991.84824
2014-11-27	5452.68672
2014-11-28	6714.80742
2014-11-29	1500.29100
2014-11-30	1055.88800

```
Name: Sales, Length: 1430, dtype: float64
```

'Target' column represents the sales data for the next time period, which you want to predict.

shift(-1) is a Pandas method that shifts the values of a Series (or column) by a specified number of periods. In this case, shift(-1) shifts the 'Sales' column one period into the future.

```
In [317]:
train_data['Target'] = train_data['Sales'].shift(-1)
train_data.head()
```

Out[317]:

	Date	Sales	Target
2011-01-01	2011-01-01	763.69800	314.22000
2011-01-02	2011-01-02	314.22000	1465.62620
2011-01-03	2011-01-03	1465.62620	353.92800
2011-01-04	2011-01-04	353.92800	48.78000
2011-01-05	2011-01-05	48.78000	465.46020

```
In [318]:
# one row will have nan value
train_data.tail()
```

Out[318]:

	Date	Sales	Target
2014-11-26	2014-11-26	6991.84824	5452.68672
2014-11-27	2014-11-27	5452.68672	6714.00742
2014-11-28	2014-11-28	6714.00742	1500.29100
2014-11-29	2014-11-29	1500.29100	1055.88800
2014-11-30	2014-11-30	1055.88800	NaN

```
In [319]:
# remove this row
train_data = train_data.dropna()
train_data.tail()
```

Out[319]:

	Date	Sales	Target
2014-11-25	2014-11-25	8687.30604	6991.84824
2014-11-26	2014-11-26	6991.84824	5452.68672
2014-11-27	2014-11-27	5452.68672	6714.00742
2014-11-28	2014-11-28	6714.00742	1500.29100
2014-11-29	2014-11-29	1500.29100	1055.88800

```
In [320]:
test_data['Target'] = test_data['Sales'].shift(-1)
test_data = test_data.dropna()
test_data.tail()
```

Out[320]:

	Date	Sales	Target
2014-12-26	2014-12-26	7106.10940	4138.38976
2014-12-27	2014-12-27	4138.38976	545.20800
2014-12-28	2014-12-28	545.20800	8123.24728
2014-12-29	2014-12-29	8123.24728	7813.22930
2014-12-30	2014-12-30	7813.22930	2625.78100

Prepare for ML training and validation

```
In [321]:
X_train_data = train_data.loc[:, ['Sales']].values
y_train_data = train_data.loc[:, ['Target']].values
X_test_data = test_data.loc[:, ['Sales']].values
y_test_data = test_data.loc[:, ['Target']].values

X_train_data.shape, y_train_data.shape, X_test_data.shape, y_test_data.shape
```

Out[321]:

((1429, 1), (1429, 1), (30, 1), (30, 1))

Random Forest Regressor

```
In [323]:
rf_model = RandomForestRegressor(n_estimators=100, random_state=0)
rf_model.fit(X_train_data, reshape(-1, 1), y_train_data)
```

```
# Forecast for December 2014
rf_forecast = rf_model.predict(X_test_data, reshape(-1, 1))
```

```
In [344]:
plt.figure(figsize=(12, 6))
plt.plot(test_data.index, test_data, label='Test Data', marker='o')
plt.plot(test_data.index, rf_forecast, label='RFR Forecast', linestyle='--')
plt.xlabel('Date')
plt.ylabel('Sales')
plt.title('Random Forest Regressor Forecast vs. Test Data for December 2014')
plt.legend()
plt.grid(True)
plt.show()
```

Metrics for Forest Regressor

```
In [ ]:
# Calculate the metrics (MAE, MSE, RMSE)
mse = mean_squared_error(test_data[1:], rf_forecast)
rmse = np.sqrt(mean_squared_log_error(test_data[1:], rf_forecast))
```

```
print(f'MAE: {mse:.2f}')
print(f'MSE: {mse:.2f}')
print(f'RMSE: {rmse:.4f}')
```

XGBoost Model

```
X_train = train_data.values[:,1:].reshape(-1, 1) # Use all data except the last data point for training
y_train = train_data.values[:,1] # Predict the next day's sales
X_train_data.shape, y_train_data.shape
```

```
Out[336]:
((1429, 1), (1429, 1))
```

```
In [337]:
xgb_model = xgb.XGBRegressor(objective='reg:squarederror', n_estimators=1000)
xgb_model.fit(X_train, y_train)
```

```
Out[337]:
XGBRegressor(base_score=None, booster=None, callbacks=None,
              colsample_bylevel=None, colsample_bynode=None,
              colsample_bytree=None, device=None, early_stopping_rounds=None,
              enable_categorical=False, eval_metric=None, feature_types=None,
              gamma=None, grow_policy=None, importance_type=None,
              interaction_constraints=None, learning_rate=None, max_bin=None,
              max_cat_threshold=None, max_cat_to_onehot=None,
              max_delta_step=None, max_depth=None, max_leaves=None,
              min_child_weight=None, missing=nan, monotone_constraints=None,
```

```
In [343]:
plt.figure(figsize=(12, 6))
plt.plot(test_data.index, test_data, label='Test Data', marker='o')
plt.plot(test_data.index, xgb_forecast, label='XGBoost Forecast', linestyle='--')
plt.xlabel('Date')
plt.ylabel('Sales')
plt.title('XGBoost Forecast vs. Test Data for December 2014')
plt.legend()
plt.grid(True)
plt.show()
```

```
In [339]:
# Forecast for December 2014
X_valid = test_data.values[:,1:].reshape(-1, 1) # Use all data except the last data point for validation
xgb_forecast = xgb_model.predict(X_valid)
```

```
In [341]:
# Calculate the metrics (MAE, MSE, RMSE)
xgb_forecast = xgb_forecast[1:len(test_data)]
mae = mean_absolute_error(test_data[1:], xgb_forecast)
mse = mean_squared_error(test_data[1:], xgb_forecast)
rmse = np.sqrt(mean_squared_log_error(test_data[1:], xgb_forecast))
```

```
print(f'MAE: {mae:.2f}')
print(f'MSE: {mse:.2f}')
print(f'RMSE: {rmse:.4f}')
```

```
ValueError                                Traceback (most recent call last)
~/mnt/Arquivos/Documentos/PROJETOS/MentoriaBoiTata/notebooks/desafio_1.eroi_merlin.ipynb Cell 134 line 4
--> 4 <a href="vscode-notebook-cell:/mnt/Arquivos/Documentos/PROJETOS/MentoriaBoiTata/notebooks/desafio_1.eroi_merlin.ipynb#Y40IsZmZlZ2Q3ODI3line0">1</a> # Calculate the metrics (MAE, MSE, RMSE)
      <a href="vscode-notebook-cell:/mnt/Arquivos/Documentos/PROJETOS/MentoriaBoiTata/notebooks/desafio_1.eroi_merlin.ipynb#Y40IsZmZlZ2Q3ODI3line1">2</a> xgb_forecast = xgb_forecast[1:len(test_data)]
--> 2 <a href="vscode-notebook-cell:/mnt/Arquivos/Documentos/PROJETOS/MentoriaBoiTata/notebooks/desafio_1.eroi_merlin.ipynb#Y40IsZmZlZ2Q3ODI3line3">4</a> mae = mean_absolute_error(test_data[1:], xgb_forecast)
      <a href="vscode-notebook-cell:/mnt/Arquivos/Documentos/PROJETOS/MentoriaBoiTata/notebooks/desafio_1.eroi_merlin.ipynb#Y40IsZmZlZ2Q3ODI3line5">6</a> rmse = np.sqrt(mean_squared_log_error(test_data[1:], xgb_forecast))
```

```
File ~/local/lib/python3.11/site-packages/sklearn/metrics/_regression.py:196, in mean_absolute_error(y_true, y_pred, sample_weight, multioutput)
    142     y_true, y_pred, *, sample_weight=None, multioutput='uniform_average'
    143 ):
    144     """Mean absolute error regression loss.
    145     Read more in the :ref:`User Guide <mean_absolute_error>`.
    146     (...)
    194     0.85 ...
    195     """
--> 196     y_type, y_true, y_pred, multioutput = _check_reg_targets(
    197         y_true, y_pred, multioutput
    198     )
    199     check_consistent_length(y_true, y_pred, sample_weight)
    200     output_errors = np.average(np.abs(y_pred - y_true), weights=sample_weight, axis=0)
```

```
File ~/local/lib/python3.11/site-packages/sklearn/utils/validation.py:397, in check_consistent_length(*arrays)
    395     uniques = np.unique(lengths)
    396     if len(uniques) > 1:
--> 397         raise ValueError(
    398             "Found input
```