

Customer Segmentation and Churn predictive



1. Customer retention (%)
2. RFM (Recency, Frequency and MonetaryValue) segmentation
3. EDA
4. Model selection and Training
5. K means & Random forest algorithm

Python ,pandas , numpy , matpoltlib , feature engineering , Machine learning

```
1 # Importing required libraries
2 import numpy as np
3 import pandas as pd
4 import scipy.stats as stats
5 from scipy.stats import norm
6 import matplotlib.pyplot as plt
7 import datetime as dt
8 import warnings
9 warnings.filterwarnings("ignore")
10 # Importing data
11 DATA_ = pd.read_csv('/content/dataset12M.csv')
12 DATA_.head()
13
```

🔍

Unnamed: 0	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID
0	416792	572558	POPPY'S PLAYHOUSE BEDROOM	6	2011-10-25	2.10	14286
1	482904	577485	VINTAGE LEAF MAGNETIC NOTEPAD	1	2011-11-20	1.45	16360

```
1 # Convert InvoiceDate to datetime
2 DATA_['InvoiceDate'] = pd.to_datetime(DATA_['InvoiceDate'])
3
4 # Create a function that truncates a date object to the first day of the month
5 def get_month(x): return dt.datetime(x.year, x.month, 1)
6
7 # Apply the function to the InvoiceDate column and create a new column called InvoiceMonth
8 DATA_['InvoiceMonth'] = DATA_['InvoiceDate'].apply(get_month)
9
10 # Group by CustomerID and select values of InvoiceMonth
11 grouping = DATA_.groupby('CustomerID')['InvoiceMonth']
12
13 # Use transform() along with min() to assign the earliest InvoiceMonth value to each customer
14 # CohortMonth is the month of the customer's first purchase
15 DATA_['CohortMonth'] = grouping.transform('min')
16
17 # Check the first 5 columns
18 DATA_.head()
19
```


Unnamed: 0	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country
0	416792	572558	POPPY'S PLAYHOUSE BEDROOM	6	2011-10-25	2.10	14286	United Kingdom
1	482904	577485	VINTAGE LEAF MAGNETIC NOTEPAD	1	2011-11-20	1.45	16360	United Kingdom

```
1 # Calculate time offset
2 # Create a function to extract integer values for years, months, and days
3 def get_date_int(df, column):
4     year = df[column].dt.year
5     month = df[column].dt.month
6     day = df[column].dt.day
7     return year, month, day
8
9 # Calculate the number of months between the first and last transaction for each customer
10 invoice_year, invoice_month, _ = get_date_int(DATA_, 'InvoiceMonth')
11 cohort_year, cohort_month, _ = get_date_int(DATA_, 'CohortMonth')
12
13 # Calculate the difference in years
```

```
14 years_diff = invoice_year - cohort_year
15
16 # Calculate the difference in months
17 months_diff = invoice_month - cohort_month
18
19 # Convert CohortMonth to 'date' format
20 DATA_['CohortMonth'] = pd.to_datetime(DATA_['CohortMonth']).dt.date
21
22 # Extract the difference in months from the first transaction/acquisition per customer
23 DATA_['CohortIndex'] = years_diff * 12 + months_diff
24 DATA_.head()
```

	Unnamed: 0	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country
0	416792	572558	22745	POPPY'S PLAYHOUSE BEDROOM	6	2011-10-25	2.10	14286	United Kingdom
1	482904	577485	23196	VINTAGE LEAF MAGNETIC NOTEPAD	1	2011-11-20	1.45	16360	United Kingdom

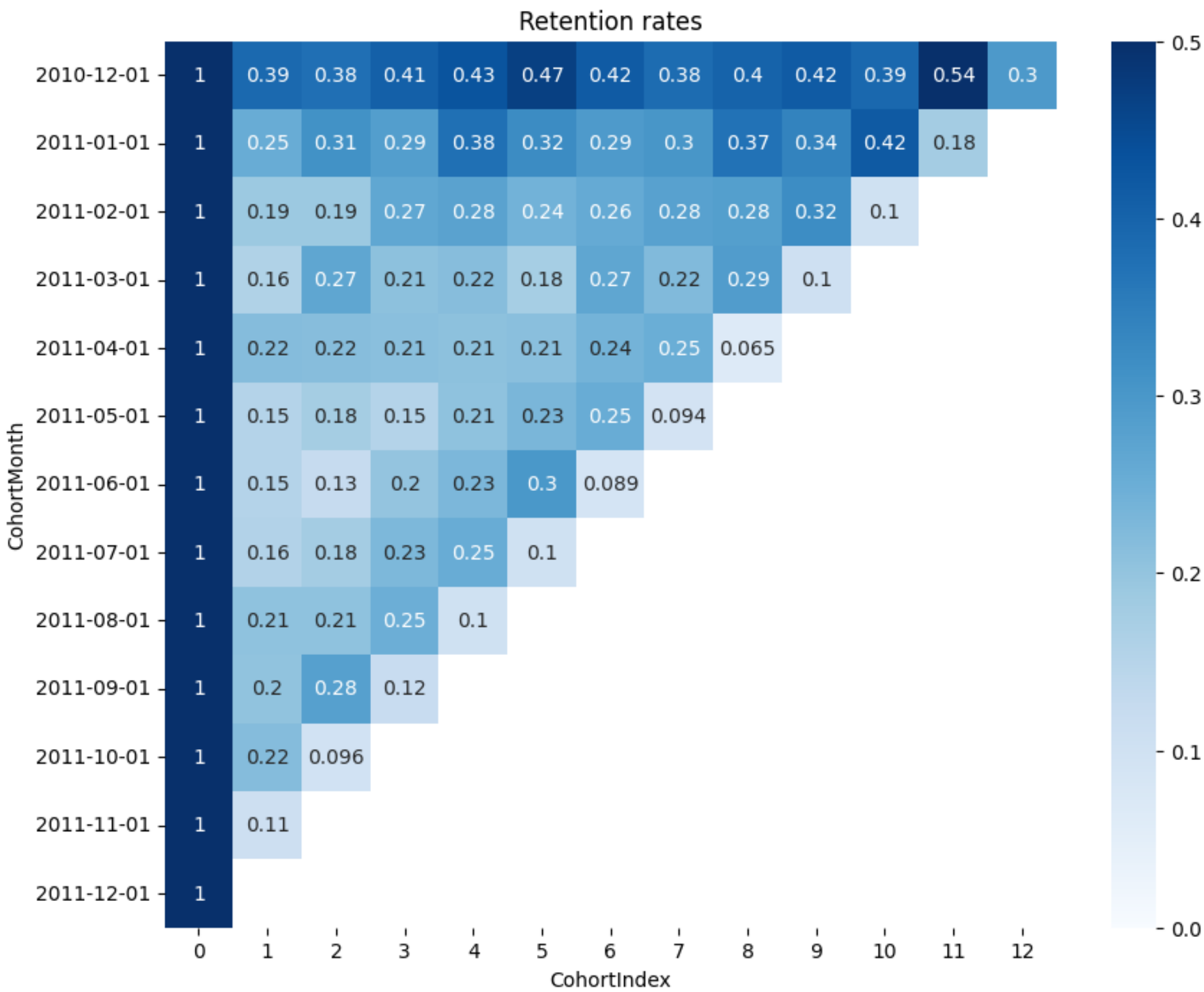
```
1 # Create DataFrame 'grouping' to group data based on CohortMonth and CohortIndex
2 grouping = DATA_.groupby(['CohortMonth', 'CohortIndex'])
3
4 # Count the number of unique values per CustomerID
5 cohort_data = grouping['CustomerID'].apply(pd.Series.nunique).reset_index()
6
7 # Create a pivot table
8 cohort_counts = cohort_data.pivot(index='CohortMonth', columns='CohortIndex', values='CustomerID')
9 cohort_counts
10
```

CohortIndex	0	1	2	3	4	5	6	7	8	9	10	11	12	
CohortMonth														
2010-12-01	383.0	149.0	145.0	156.0	165.0	180.0	160.0	147.0	154.0	160.0	150.0	208.0	113.0	
2011-01-01	429.0	109.0	134.0	123.0	161.0	139.0	126.0	130.0	160.0	146.0	180.0	77.0	NaN	
2011-02-01	352.0	67.0	67.0	94.0	97.0	85.0	91.0	98.0	100.0	113.0	36.0	NaN	NaN	
2011-03-01	422.0	67.0	113.0	88.0	91.0	74.0	113.0	94.0	122.0	44.0	NaN	NaN	NaN	
2011-04-01	279.0	61.0	60.0	59.0	58.0	59.0	67.0	70.0	18.0	NaN	NaN	NaN	NaN	
2011-05-01	267.0	41.0	47.0	41.0	55.0	62.0	68.0	25.0	NaN	NaN	NaN	NaN	NaN	
2011-06-01	214.0	33.0	27.0	43.0	49.0	64.0	19.0	NaN	NaN	NaN	NaN	NaN	NaN	
2011-07-01	185.0	29.0	33.0	42.0	47.0	19.0	NaN	NaN	NaN	NaN	NaN	NaN	NaN	
2011-08-01	145.0	30.0	30.0	36.0	15.0	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	
2011-09-01	284.0	58.0	80.0	34.0	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	
2011-10-01	332.0	72.0	32.0	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	
2011-11-01	311.0	34.0	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	
2011-12-01	40.0	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	

```
1 # Select the first column and save it as cohort_sizes
2 cohort_sizes = cohort_counts.iloc[:, 0]
3
4 # Divide cohort_counts by cohort_sizes for all rows
5 retention = cohort_counts.divide(cohort_sizes, axis=0)
6 retention.round(3) * 100
7
```

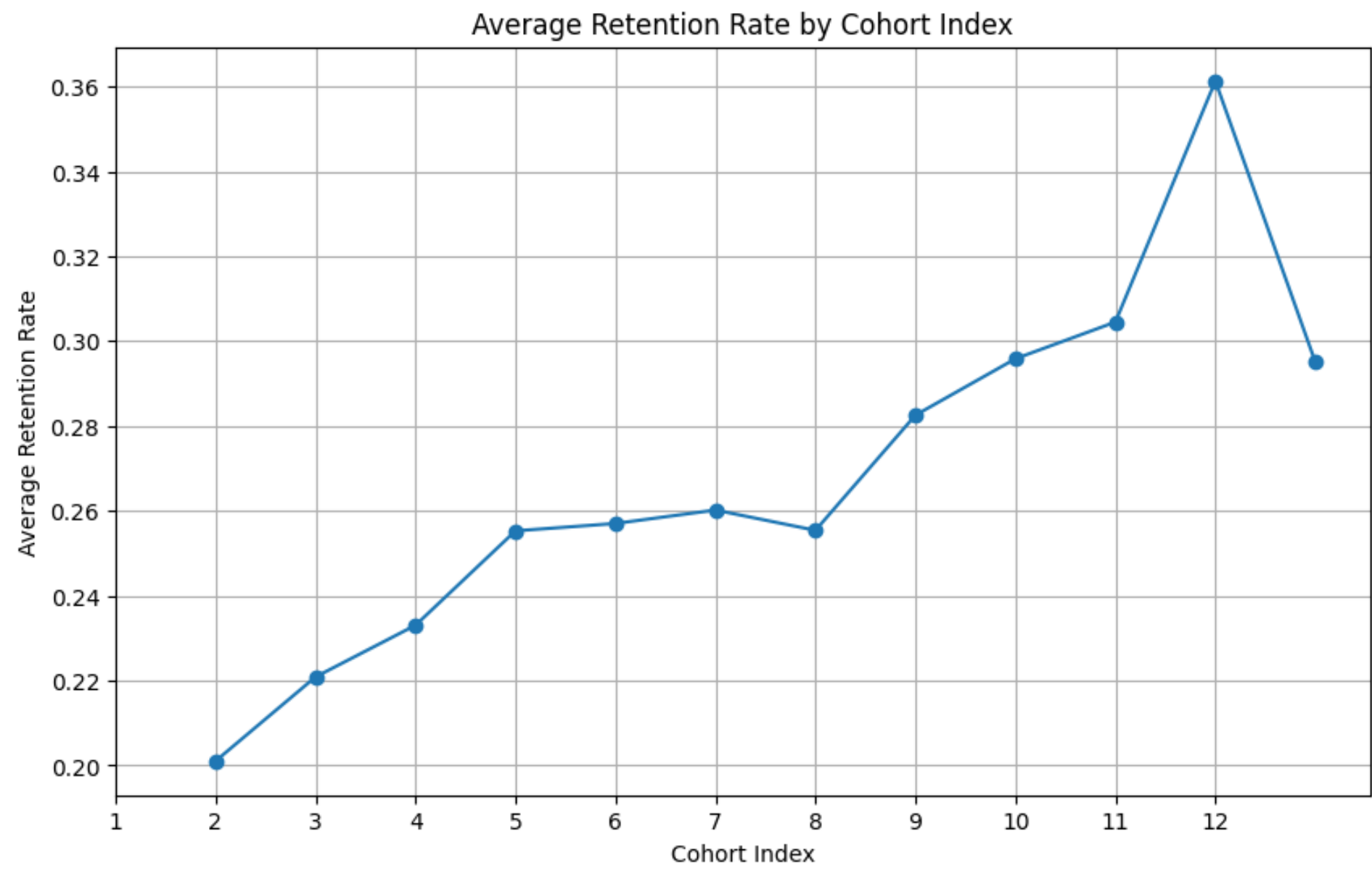
CohortIndex	0	1	2	3	4	5	6	7	8	9	10	11	12
CohortMonth													
2010-12-01	100.0	38.9	37.9	40.7	43.1	47.0	41.8	38.4	40.2	41.8	39.2	54.3	29.5
2011-01-01	100.0	25.4	31.2	28.7	37.5	32.4	29.4	30.3	37.3	34.0	42.0	17.9	NaN
2011-02-01	100.0	19.0	19.0	26.7	27.6	24.1	25.9	27.8	28.4	32.1	10.2	NaN	NaN
2011-03-01	100.0	15.9	26.8	20.9	21.6	17.5	26.8	22.3	28.9	10.4	NaN	NaN	NaN
2011-04-01	100.0	21.9	21.5	21.1	20.8	21.1	24.0	25.1	6.5	NaN	NaN	NaN	NaN
2011-05-01	100.0	15.4	17.6	15.4	20.6	23.2	25.5	9.4	NaN	NaN	NaN	NaN	NaN
2011-06-01	100.0	15.4	12.6	20.1	22.9	29.9	8.9	NaN	NaN	NaN	NaN	NaN	NaN
2011-07-01	100.0	15.7	17.8	22.7	25.4	10.3	NaN	NaN	NaN	NaN	NaN	NaN	NaN

```
1 # Visualize retention rates as a heatmap
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4
5 plt.figure(figsize=(10, 8))
6 plt.title('Retention rates')
7 sns.heatmap(data=retention,
8             annot=True,
9             vmin=0.0,
10            vmax=0.5,
11            cmap='Blues')
12 plt.show()
13
```



```
1 import numpy as np
2
3 # Calculate the mean retention rate for each CohortIndex
4 average_retention = retention.mean(axis=0)
5
6 # Exclude CohortIndex 0 because this is where we got first transaction from customer
7 average_retention = average_retention[1:]
8
9 # Create the line chart
```

```
10 plt.figure(figsize=(10, 6))
11 plt.plot(average_retention.index, average_retention.values, marker='o')
12 plt.xlabel('Cohort Index')
13 plt.ylabel('Average Retention Rate')
14 plt.title('Average Retention Rate by Cohort Index')
15 plt.xticks(np.arange(len(average_retention.index)), average_retention.index)
16 plt.grid(True)
17 plt.show()
18
```



FROM ABOVE PLOTS WE CAN FIND The retention rate appears to decrease with increasing cohort index(since 0 or initial index), indicating that customers tend to be less engaged and make fewer repeat purchases over time. Cohort 0 has the highest retention rate since it represents the initial cohort where all customers made their first transactions. After cohort 0, the retention rate drops initially but shows a relatively stable pattern afterward, suggesting that customers who remain engaged tend to maintain their level of activity. Monitoring retention rates by cohort index can help identify trends and assess the effectiveness of customer retention strategies over time. It is crucial to focus on strategies to retain customers beyond the initial cohort to ensure sustained business growth.

2. Preprocess Data

```
1 # Print the minimum and maximum dates in the 'InvoiceDate' column
2 print('Min: {}; Max: {}'.format(min(DATA_.InvoiceDate), max(DATA_.InvoiceDate)))
3 DATA_.head()
4
```

Min: 2010-12-10 00:00:00; Max: 2011-12-09 00:00:00

Unnamed: 0									
	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country	
0	416792	572558	22745 POPPY'S PLAYHOUSE BEDROOM	6	2011-10-25	2.10	14286	United Kingdom	
1	482904	577485	23196 VINTAGE LEAF MAGNETIC NOTEPAD	1	2011-11-20	1.45	16360	United Kingdom	

RFM

```
1 # Create the 'TotalSum' column
2 DATA_['TotalSum'] = DATA_['Quantity'] * DATA_['UnitPrice']
3
4 # Convert 'InvoiceDate' to the 'datetime' format
5 DATA_['InvoiceDate'] = pd.to_datetime(DATA_['InvoiceDate'])
6
```



```
7 # Create 'snapshot_date' to select the most recent date in the entire dataset and add 1 to simulate future date
8 snapshot_date = max(DATA_['InvoiceDate']) + dt.timedelta(days=1)
9
10 # Group the data by 'CustomerID'
11 datamart = DATA_.groupby(['CustomerID']).agg({
12     'InvoiceDate': lambda x: (snapshot_date - x.max()).days, # Days between the analysis date and the last invoice date
13     'InvoiceNo': 'count', # Number of invoices (transactions) per customer
14     'TotalSum': 'sum' # Total monetary value (spending) per customer
15 })
16
17 # Rename columns for easy interpretation
18 datamart.rename(columns={'InvoiceDate': 'Recency',
19                          'InvoiceNo': 'Frequency',
20                          'TotalSum': 'MonetaryValue'}, inplace=True)
21
22 # Datamart is a table where each row represents a customer with their recency, frequency, and monetary value
23 datamart.head()
24
```

	Recency	Frequency	MonetaryValue
CustomerID			
12747	3	25	948.70
12748	1	888	7046.16
12749	4	37	813.45
12820	4	17	268.02
12822	71	9	146.15

```
1 r_labels = range(4, 0, -1) # Customers who have been more recent will be better than less recent
2 f_labels = range(1, 5)
3 m_labels = range(1, 5)
4
```

```
1 r_groups = pd.qcut(datamart['Recency'], q=4, labels=r_labels)
2 f_groups = pd.qcut(datamart['Frequency'], q=4, labels=f_labels)
3 m_groups = pd.qcut(datamart['MonetaryValue'], q=4, labels=m_labels)
4
```

```
1 datamart = datamart.assign(R=r_groups.values, F=f_groups.values, M=m_groups.values)
```

```
1 # Convert 'R', 'F', and 'M' columns to numeric data types (integer)
2 datamart['R'] = datamart['R'].astype(int)
3 datamart['F'] = datamart['F'].astype(int)
4 datamart['M'] = datamart['M'].astype(int)
5
6 # Calculate RFM scores using the modified formula
7 datamart['RFM_Score'] = 0.5 * datamart['R'] + 0.3 * datamart['F'] + 0.2 * datamart['M']
8
9 # Create the 'RFM_Decile' column with unique bin edges using 'duplicates' parameter
10 datamart['RFM_Decile'] = pd.qcut(datamart['RFM_Score'], q=10, labels=False, duplicates='drop')
11 datamart
```

Recency Frequency MonetaryValue R F M RFM_Score RFM_Decile



CustomerID

12747	3	25	948.70	4	4	4	4.0	8
12748	1	888	7046.16	4	4	4	4.0	8
12749	4	37	813.45	4	4	4	4.0	8

```
1 # Define the threshold for churn based on 'RFM_Decile'
2 churn_threshold = 4
3
4 # Create the 'Churn' column based on the 'RFM_Decile' threshold
5 datamart['Churn'] = (datamart['RFM_Decile'] <= churn_threshold).astype(int)
6
7 # Print the DataFrame to check the updated columns
8 datamart
```

Recency Frequency MonetaryValue R F M RFM_Score RFM_Decile Churn

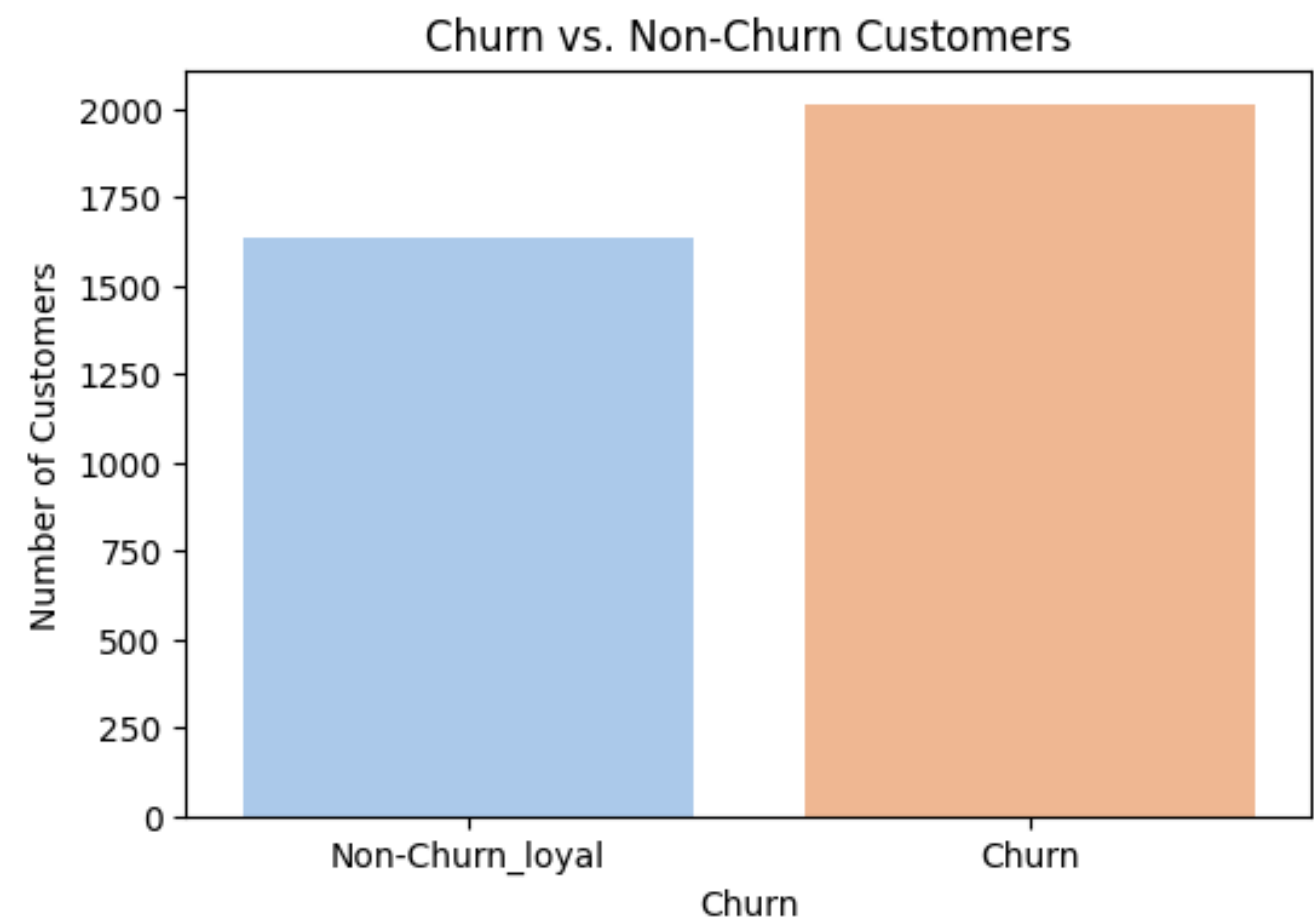


CustomerID

12747	3	25	948.70	4	4	4	4.0	8	0
12748	1	888	7046.16	4	4	4	4.0	8	0
12749	4	37	813.45	4	4	4	4.0	8	0
12820	4	17	268.02	4	3	3	3.5	7	0
12822	71	9	146.15	2	2	3	2.2	3	1
...
18280	278	2	38.70	1	1	1	1.0	0	1
18281	181	2	31.80	1	1	1	1.0	0	1
18282	8	2	30.70	4	1	1	2.5	4	1
18283	4	152	432.93	4	4	4	4.0	8	0
18287	43	15	395.76	3	3	4	3.2	7	0

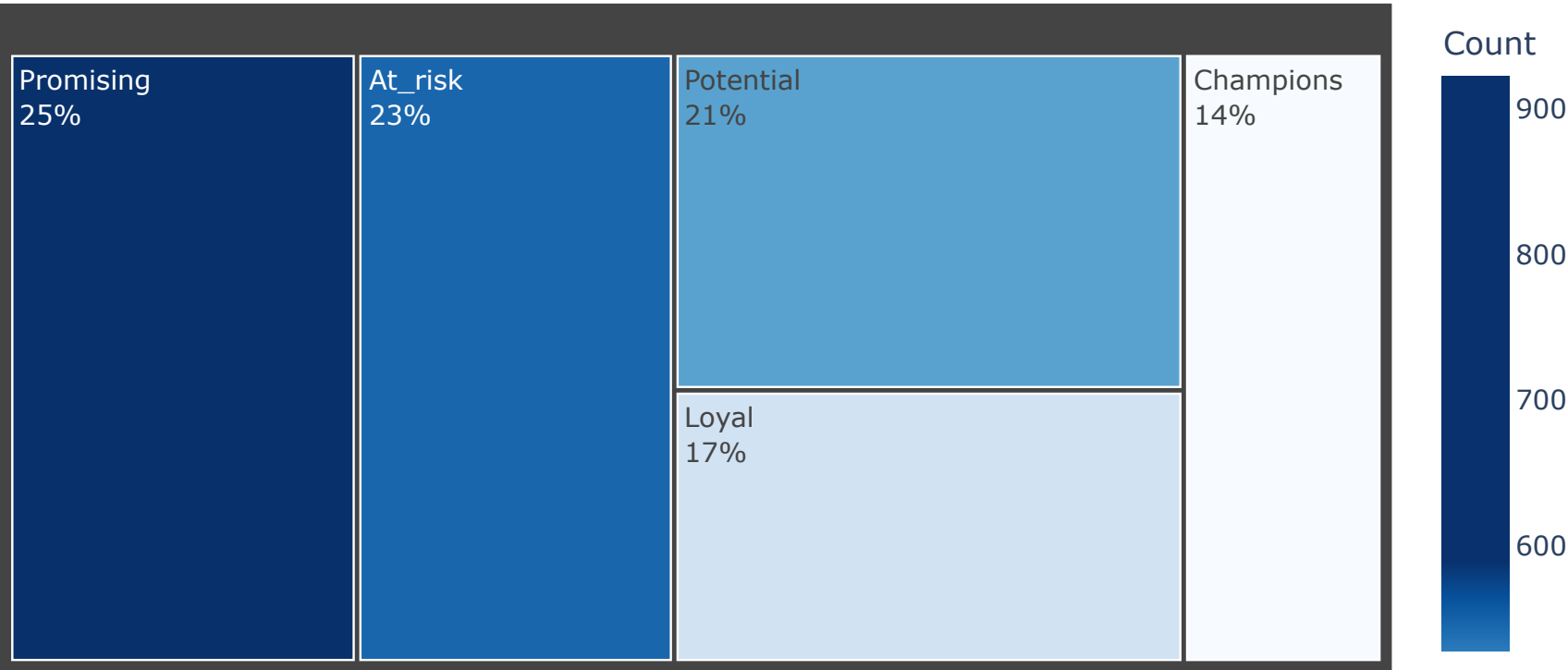
3643 rows × 9 columns

```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4 # Count the number of churned and non-churned customers
5 churn_counts = datamart['Churn'].value_counts()
6
7 # Create a bar plot
8 plt.figure(figsize=(6, 4))
9 sns.barplot(x=churn_counts.index, y=churn_counts.values, palette='pastel')
10 plt.xlabel('Churn')
11 plt.ylabel('Number of Customers')
12 plt.title('Churn vs. Non-Churn Customers')
13 plt.xticks([1, 0], ['Churn', 'Non-Churn_loyal'])
14 plt.show()
```



```
1 def map_churn_segment(x):
2     if x in [0, 1]:
3         return 'At_risk'
4     elif x in [2, 3]:
5         return 'Potential'
6     elif x in [4, 5]:
7         return 'Loyal'
8     elif x in [6, 7]:
9         return 'Promising'
10    elif x == 8:
11        return 'Champions'
12    else:
13        return 'other'
14
15 datamart['Churn_Segment'] = datamart['RFM_Decile'].apply(map_churn_segment)
16
17 import plotly.express as px
18
19 churn_segment_counts = datamart['Churn_Segment'].value_counts().reset_index()
20 churn_segment_counts.columns = ['Churn_Segment', 'Count']
21
22 # Calculate percentage distribution
23 churn_segment_counts['Percentage'] = (churn_segment_counts['Count'] / churn_segment_counts['Count
24 churn_segment_counts['Percentage'] = churn_segment_counts['Percentage'].round(2).astype(str) + '%'
25
26 fig = px.treemap(churn_segment_counts,
27                  path=['Churn_Segment'],
28                  values='Count',
29                  title='Churn Segment Distribution',
30                  color='Count',
31                  color_continuous_scale='Blues',
32                  hover_data=['Percentage'])
33
34 # Update tooltip text to display label and percentage
35 fig.update_traces(textinfo='label+percent root')
36
37 # Adjust figsize
38 fig.update_layout(width=800, height=480)
39
40 fig.show()
41
```

Churn Segment Distribution



The above plot indicates that a significant portion of the customer base consists of loyal and engaged customers ("Champions" and "Loyal"). However, there is a need to focus on the "At_risk" segment and explore strategies to retain these customers. Additionally, efforts can be

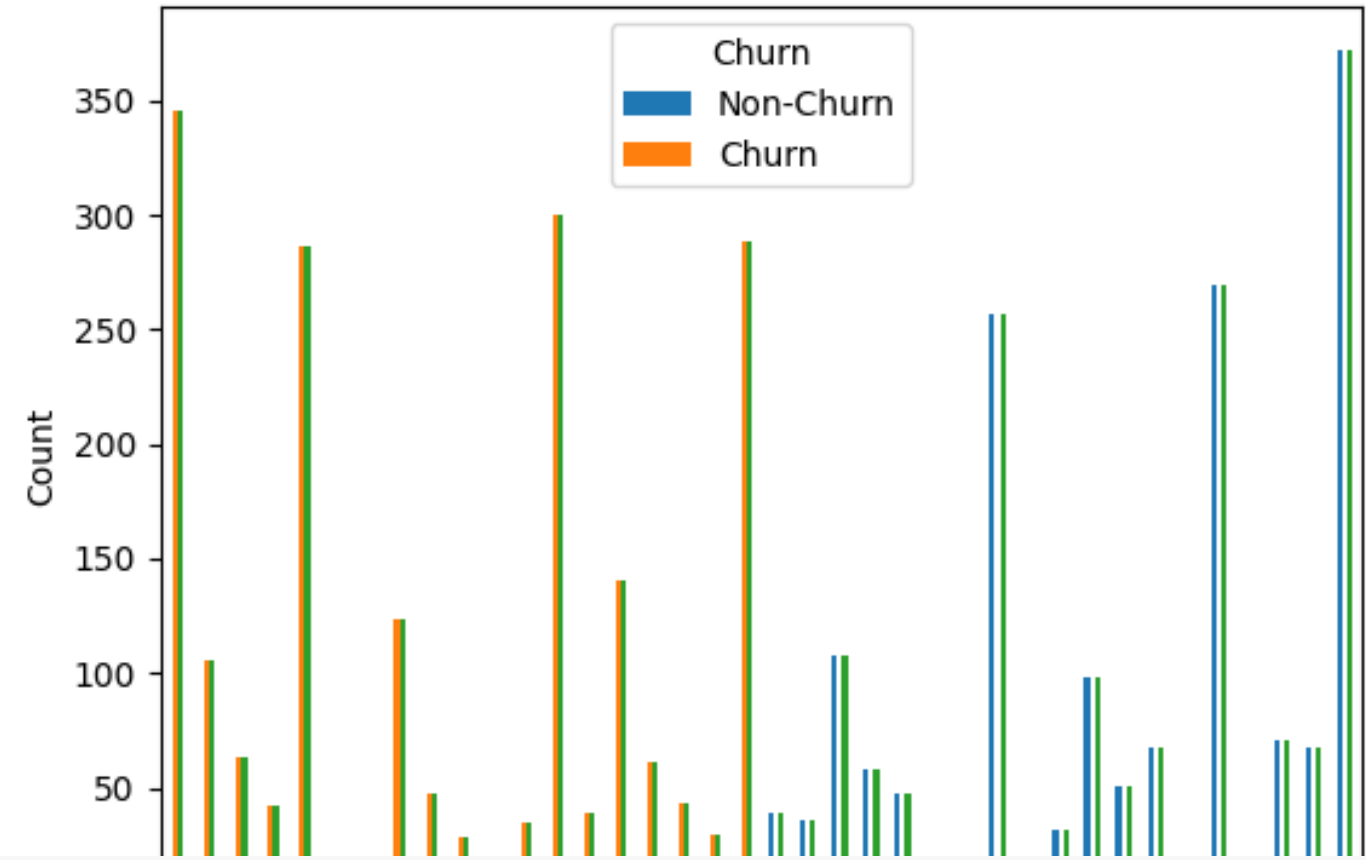
directed towards converting the "Potential" customers into more loyal ones, as they show promising engagement levels. Understanding and targeting these segments appropriately can help in optimizing marketing and retention strategies.

```
1 # Group customers by 'RFM_Score' and 'Churn', and count churn and non-churn customers
2 score_wise_churn_counts = datamart.groupby(['RFM_Score', 'Churn']).size().unstack()
3
4 # Add a column for the total count of customers in each RFM_Score category
5 score_wise_churn_counts['Total'] = score_wise_churn_counts.sum(axis=1)
6
7 # Rename the columns for better readability
8 score_wise_churn_counts.rename(columns={0: 'Non-Churn', 1: 'Churn'}, inplace=True)
9
10 # Display the count of churn and non-churn customers score-wise based on RFM_Score
11 print(score_wise_churn_counts)
12
```

	Churn	Non-Churn	Churn	Total
RFM_Score				
1.0		NaN	345.0	345.0
1.2		NaN	105.0	105.0
1.3		NaN	63.0	63.0
1.4		NaN	42.0	42.0
1.5		NaN	286.0	286.0
1.6		NaN	17.0	17.0
1.6		NaN	16.0	16.0
1.7		NaN	124.0	124.0
1.8		NaN	47.0	47.0
1.8		NaN	28.0	28.0
1.9		NaN	1.0	1.0
1.9		NaN	35.0	35.0
2.0		NaN	300.0	300.0
2.1		NaN	39.0	39.0
2.2		NaN	140.0	140.0
2.3		NaN	61.0	61.0
2.3		NaN	43.0	43.0
2.4		NaN	29.0	29.0
2.5		NaN	288.0	288.0
2.6		39.0	NaN	39.0
2.7		36.0	NaN	36.0
2.7		108.0	NaN	108.0
2.8		58.0	NaN	58.0
2.8		47.0	NaN	47.0
2.9		9.0	NaN	9.0
2.9		14.0	NaN	14.0
3.0		257.0	NaN	257.0
3.1		5.0	NaN	5.0
3.1		32.0	NaN	32.0
3.2		98.0	NaN	98.0
3.3		51.0	NaN	51.0
3.3		67.0	NaN	67.0
3.4		16.0	NaN	16.0
3.5		269.0	NaN	269.0
3.6		18.0	NaN	18.0
3.7		71.0	NaN	71.0
3.8		67.0	NaN	67.0
4.0		372.0	NaN	372.0

```
1 import matplotlib.pyplot as plt
2
3 # Plot the stacked bar chart
4 ax = score_wise_churn_counts.plot(kind='bar')
5
6 # Customize the x-axis labels
7 x_labels = [f'{float(score):.2f}' for score in score_wise_churn_counts.index]
8 ax.set_xticklabels(x_labels, rotation=0)
9
10 plt.xlabel('RFM Score')
11 plt.ylabel('Count')
12 plt.title('Churn and Non-Churn Customers Score-wise')
13 plt.legend(title='Churn', labels=['Non-Churn','Churn'])
14 plt.xticks(rotation=90)
15 plt.show()
16
```

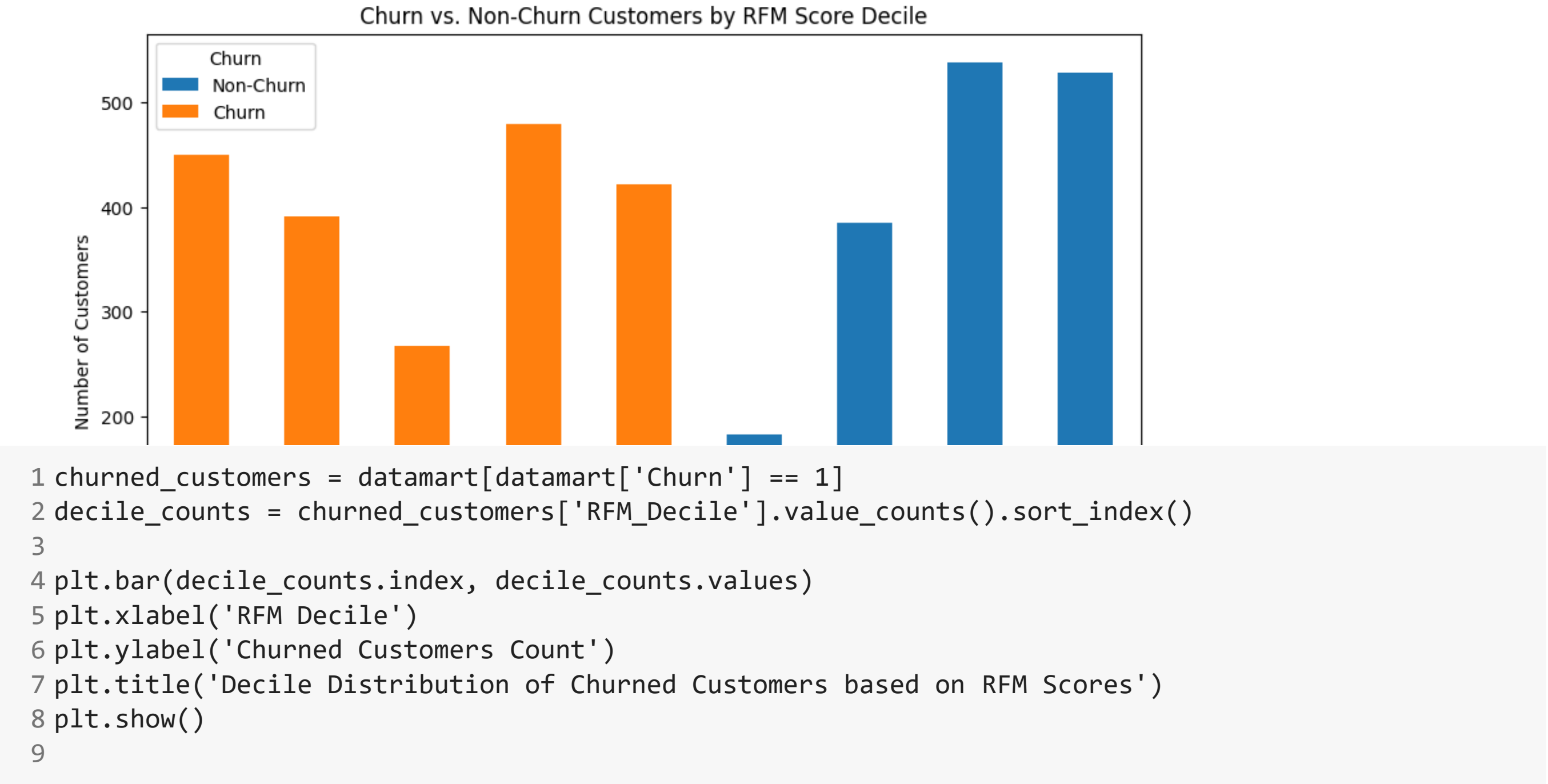
Churn and Non-Churn Customers Score-wise



```
1 # Group customers by their RFM score deciles and count churn and non-churn customers
2 score_wise_counts = datamart.groupby(['RFM_Decile', 'Churn']).size().unstack()
3
4 # Add a column for the total count of customers in each score decile
5 score_wise_counts['Total'] = score_wise_counts.sum(axis=1)
6
7 # Rename the columns for better readability
8 score_wise_counts.rename(columns={0: 'Non-Churn', 1: 'Churn'}, inplace=True)
9
10 # Display the count of churn and non-churn customers score-wise
11 print(score_wise_counts)
```

Churn	Non-Churn	Churn	Total
RFM_Decile			
0	NaN	450.0	450.0
1	NaN	391.0	391.0
2	NaN	268.0	268.0
3	NaN	479.0	479.0
4	NaN	421.0	421.0
5	183.0	NaN	183.0
6	385.0	NaN	385.0
7	538.0	NaN	538.0
8	528.0	NaN	528.0

```
1 import matplotlib.pyplot as plt
2
3 # Group customers by their RFM score deciles and count churn and non-churn customers
4 score_wise_counts = datamart.groupby(['RFM_Decile', 'Churn']).size().unstack()
5
6 # Add a column for the total count of customers in each score decile
7 score_wise_counts['Total'] = score_wise_counts.sum(axis=1)
8
9 # Rename the columns for better readability
10 score_wise_counts.rename(columns={0: 'Non-Churn', 1: 'Churn'}, inplace=True)
11
12 # Plot the count of churn and non-churn customers score-wise
13 score_wise_counts[['Non-Churn', 'Churn']].plot(kind='bar', stacked=True, figsize=(10, 6))
14 plt.xlabel('RFM Score Decile')
15 plt.ylabel('Number of Customers')
16 plt.title('Churn vs. Non-Churn Customers by RFM Score Decile')
17 plt.legend(title='Churn', loc='upper left', labels=['Non-Churn', 'Churn'])
18 plt.xticks(rotation=0)
19 plt.show()
```



```
1 print(datamart.columns)

Index(['Recency', 'Frequency', 'MonetaryValue', 'R', 'F', 'M', 'RFM_Score',
      'RFM_Decile', 'Churn', 'Churn_Segment'],
      dtype='object')

1 # Convert 'R', 'F', and 'M' columns to string data types
2 datamart['R'] = datamart['R'].astype(str)
3 datamart['F'] = datamart['F'].astype(str)
4 datamart['M'] = datamart['M'].astype(str)
5
6 # Concatenate 'R', 'F', and 'M' columns to create the 'RFM_Segment' column
7 datamart['RFM_Segment'] = datamart['R'] + datamart['F'] + datamart['M']
8
9 # Now, let's check the DataFrame to see if the 'RFM_Segment' column is created
10 print(datamart.head())
11
```

CustomerID	Recency	Frequency	MonetaryValue	R	F	M	RFM_Score	RFM_Decile	\
12747	3	25	948.70	4	4	4	4.0	8	
12748	1	888	7046.16	4	4	4	4.0	8	
12749	4	37	813.45	4	4	4	4.0	8	
12820	4	17	268.02	4	3	3	3.5	7	
12822	71	9	146.15	2	2	3	2.2	3	

CustomerID	Churn	Churn_Segment	RFM_Segment
12747	0	Champions	444
12748	0	Champions	444
12749	0	Champions	444
12820	0	Promising	433
12822	1	Potential	223

```
1 # Datamart RFM_Segment analysis
2 datamart.groupby('RFM_Segment').size().sort_values(ascending=False)[:10]
```

RFM_Segment
444 372
111 345
211 169
344 156
233 129
222 128
333 120
122 117
311 114
433 113
dtype: int64

```
1 datamart[datamart['RFM_Segment']=='444']
```

	Recency	Frequency	MonetaryValue	R	F	M	RFM_Score	RFM_Decile	Churn	Churn_Segment	R
CustomerID											
12747	3	25	948.70	4	4	4	4.0	8	0	Champions	
12748	1	888	7046.16	4	4	4	4.0	8	0	Champions	
12749	4	37	813.45	4	4	4	4.0	8	0	Champions	
12839	3	54	947.63	4	4	4	4.0	8	0	Champions	
12841	5	70	630.95	4	4	4	4.0	8	0	Champions	
...
18225	4	56	1175.02	4	4	4	4.0	8	0	Champions	
18229	12	30	1626.54	4	4	4	4.0	8	0	Champions	
18245	8	40	673.61	4	4	4	4.0	8	0	Champions	
18272	3	38	642.44	4	4	4	4.0	8	0	Champions	
18283	4	152	432.93	4	4	4	4.0	8	0	Champions	

372 rows × 11 columns

```
1 # Analyze statistical summaries of RFM parameters for each RFM_Segment level
2 datamart.groupby('RFM_Segment').agg({
3     'Recency': 'mean',  # Mean Recency for each RFM_Segment
4     'Frequency': 'mean',  # Mean Frequency for each RFM_Segment
5     'MonetaryValue': ['mean', 'count']  # Mean MonetaryValue and count of customers for each RFM_Segment
6 }).round(1)
7
```

```
1 # Analyze statistical summaries of RFM parameters for each RFM_Score level
2 datamart.groupby('RFM_Score').agg({
3     'Recency': 'mean',  # Mean Recency for each RFM_Score level
4     'Frequency': 'mean',  # Mean Frequency for each RFM_Score level
5     'MonetaryValue': ['mean', 'count']  # Mean MonetaryValue and count of customers for each RFM_Score level
6 }).round(1)
7
```

	Recency	Frequency	MonetaryValue	
	mean	mean	mean	count
RFM_Score				
1.0	246.9	2.1	28.4	345
1.2	234.5	2.9	82.4	105
1.3	246.5	6.5	38.4	63
1.4	254.1	2.3	202.6	42
1.5	147.1	3.7	56.2	286
1.6	226.0	13.1	43.9	17
1.6	225.9	2.2	1434.6	16
1.7	138.4	4.8	127.7	124
1.8	229.7	13.3	92.2	47
1.8	85.7	6.2	33.1	28
1.9	162.0	22.0	52.8	1
1.9	101.1	3.2	293.2	35
2.0	89.7	6.1	96.1	300
2.1	93.2	10.6	235.3	39
2.2	75.6	6.5	183.6	140
2.3	121.9	19.0	122.0	61
2.3	32.7	6.8	37.8	43
2.4	56.0	4.4	394.7	29
2.5	59.5	10.4	165.4	288
2.6	49.6	16.1	149.6	39
2.7	9.5	3.1	79.2	36
2.7	64.3	12.2	468.0	108
2.8	33.7	14.9	96.3	58
2.8	51.7	23.1	157.0	47
2.9	7.8	3.6	218.1	9
2.9	32.1	8.3	578.8	14
3.0	40.3	19.2	334.0	257
3.1	7.0	2.2	2026.2	5
3.1	23.7	20.4	74.3	32
3.2	25.5	13.2	565.0	98
3.3	9.6	15.9	97.7	51
3.3	31.6	35.6	241.8	67
3.4	9.5	7.6	658.6	16
3.5	22.4	34.6	604.7	269
3.6	9.4	27.1	101.3	18
3.7	10.5	16.7	776.4	71
3.8	10.3	38.6	231.1	67
4.0	8.0	75.6	1653.9	372

```
1 rfm_score_stats = datamart['RFM_Score'].describe()
2
3 print(rfm_score_stats)
```

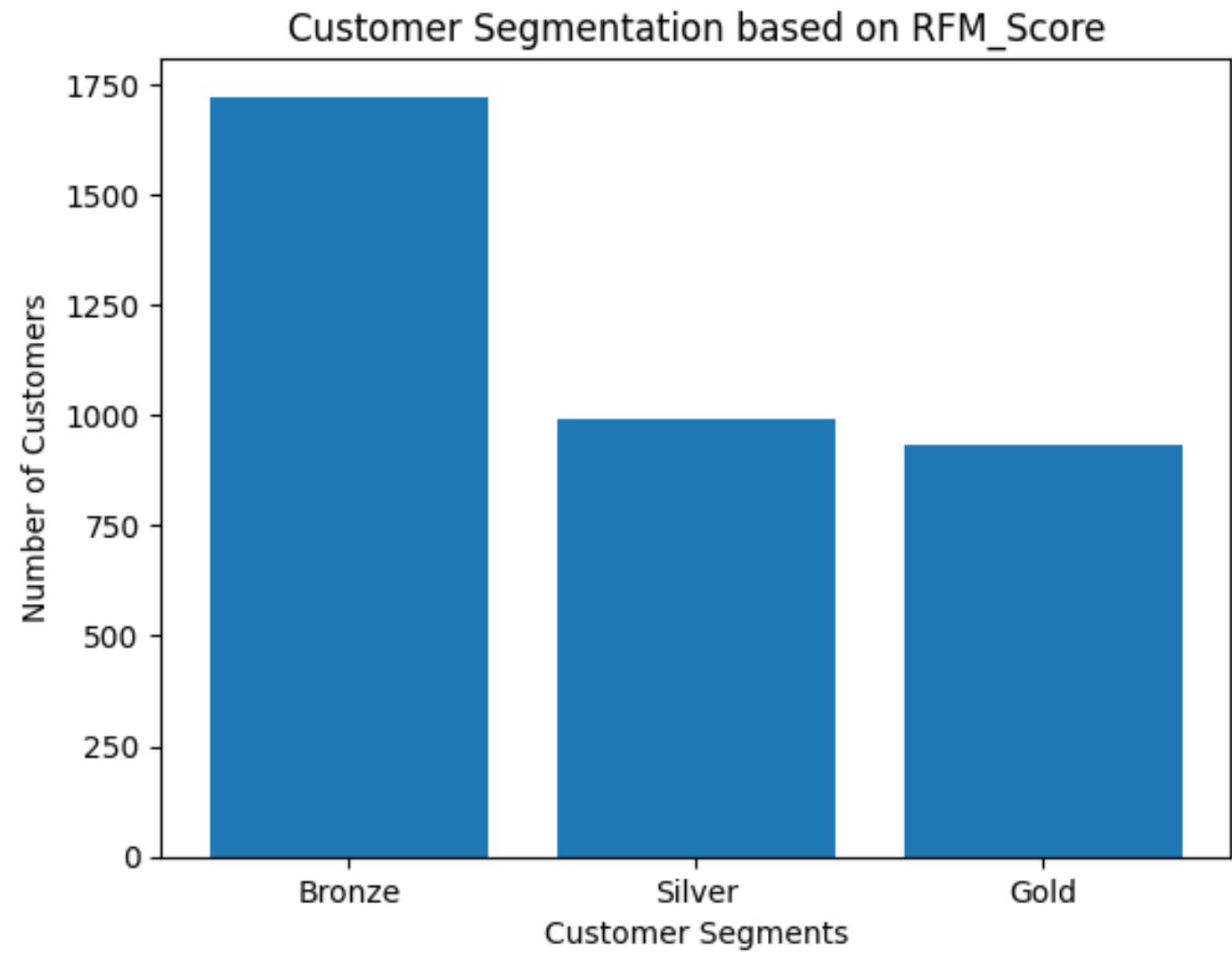


```
count    3643.000000
mean      2.479275
std       0.938874
min       1.000000
25%       1.700000
50%       2.500000
75%       3.300000
max       4.000000
Name: RFM_Score, dtype: float64
```

```
1 # Set cutoff points based on the number of RFM_Score levels
2 def segment_me(df):
3     if df['RFM_Score'] >= 3.3:
4         return 'Gold'
5     elif (df['RFM_Score'] >= 2.5) and (df['RFM_Score'] < 3.3):
6         return 'Silver'
7     else:
8         return 'Bronze'
9
10 # Apply the segmentation function to create the 'General_Segment' column
11 datamart['General_Segment'] = datamart.apply(segment_me, axis=1)
12
13 datamart.head()
```

	Recency	Frequency	MonetaryValue	R	F	M	RFM_Score	RFM_Decile	Churn	Churn_Segment	RFM
CustomerID											
12747	3	25	948.70	4	4	4	4.0	8	0	Champions	
12748	1	888	7046.16	4	4	4	4.0	8	0	Champions	
12749	4	37	813.45	4	4	4	4.0	8	0	Champions	
12820	4	17	268.02	4	3	3	3.5	7	0	Promising	
12822	71	9	146.15	2	2	3	2.2	3	1	Potential	

```
1 import matplotlib.pyplot as plt
2
3 # Count the number of customers in each 'General_Segment'
4 segment_counts = datamart['General_Segment'].value_counts()
5
6 # Plot the bar graph
7 plt.bar(segment_counts.index, segment_counts.values)
8 plt.xlabel('Customer Segments')
9 plt.ylabel('Number of Customers')
10 plt.title('Customer Segmentation based on RFM_Score')
11 plt.show()
12
```



```
1 # Analyze statistical summaries of RFM parameters for each General_Segment level
2 datamart.groupby('General_Segment').agg({
3     'Recency': 'mean',  # Mean Recency for each General_Segment level
4     'Frequency': 'mean',  # Mean Frequency for each General_Segment level
5     'MonetaryValue': ['mean', 'count']  # Mean MonetaryValue and count of customers for each Gener
6 }).round(1)
7
```

	Recency	Frequency	MonetaryValue	
	mean	mean	mean	count
General_Segment				
Bronze	157.6	5.3	107.4	1721
Gold	14.3	48.4	947.4	931
Silver	45.3	14.2	286.2	991

In some cases, taking multiple trial-and-error attempts to find the right segmentation points can be cumbersome, low-precision, and result in poor business decisions.

K-Means clustering is an ideal machine learning model for customer segmentation without arbitrarily chosen thresholds (as we have been doing so far).

K-Means seeks a balance between having an appropriate commercial customer segmentation and minimizing the error of prediction (SSE).

K-means clustering is one of the most popular, simple, and fast unsupervised machine learning methods.

Key assumptions of k-means clustering regarding variables (R/F/M):

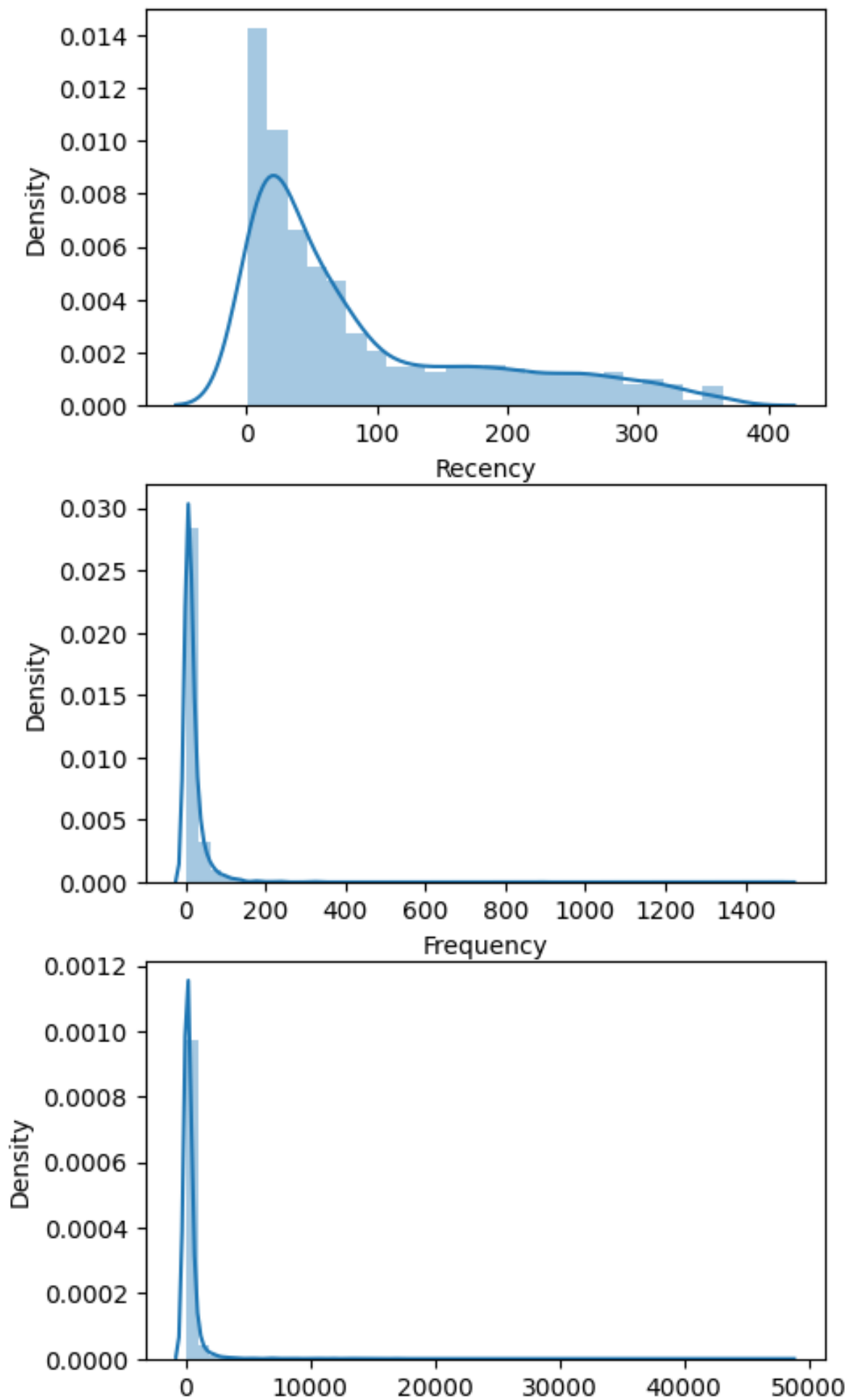
Variables should have symmetric distributions (skewed variables can be managed with logarithmic transformation). Variables should have the same mean (to ensure equal importance is assigned to each variable). Variables should have the same variance (to ensure equal importance is assigned to each variable). Exploring data for k-means clustering Steps:

Explore variables with asymmetric distributions - apply logarithmic transformation to them. Normalize/Standardize variables to have the same mean. Normalize/Standardize variables to have the same variance. Store the processed variables as a separate 'array' to be used later for clustering.

```
1 # Filter columns 'Recency', 'Frequency', and 'MonetaryValue' into a new DataFrame
2 datamart_rfm = datamart[['Recency', 'Frequency', 'MonetaryValue']]
3
4 # Perform statistical summary on the 'datamart_rfm' DataFrame
5 datamart_rfm.describe()
6
```

	Recency	Frequency	MonetaryValue
count	3643.00000	3643.000000	3643.000000
mean	90.43563	18.714247	370.694387
std	94.44651	43.754468	1347.443451
min	1.00000	1.000000	0.650000
25%	19.00000	4.000000	58.705000
50%	51.00000	9.000000	136.370000
75%	139.00000	21.000000	334.350000
max	365.00000	1497.000000	48060.350000

```
1 # Create a vertical layout of three subplots
2 plt.figure(figsize=(5, 10))
3 plt.subplot(3, 1, 1)
4 sns.distplot(datamart_rfm['Recency'])  # Plot distribution of 'Recency'
5 plt.subplot(3, 1, 2)
6 sns.distplot(datamart_rfm['Frequency'])  # Plot distribution of 'Frequency'
7 plt.subplot(3, 1, 3)
8 sns.distplot(datamart_rfm['MonetaryValue'])  # Plot distribution of 'MonetaryValue'
9
10 plt.show()
```

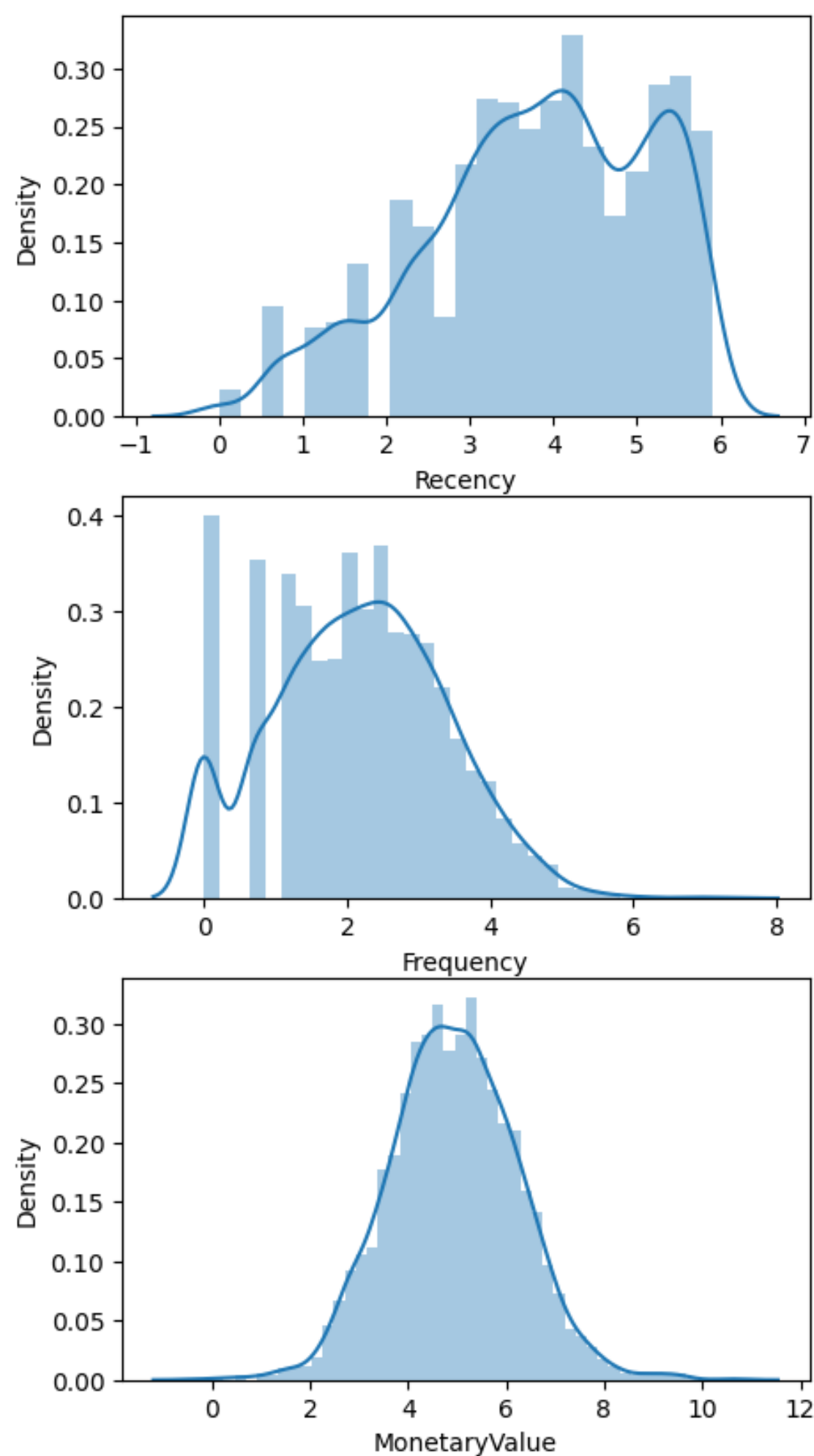


The output will show three separate plots, each representing the distribution of 'Recency', 'Frequency', and 'MonetaryValue' respectively. These plots help visualize the spread of data for each metric and identify any skewness or patterns in the data.

```
1 # Apply logarithmic transformation to 'datamart_rfm' and store the result in 'datamart_rfm_log'
2 datamart_rfm_log = np.log(datamart_rfm)
3
4 # Perform statistical summary on the 'datamart_rfm_log' DataFrame
5 datamart_rfm_log.describe()
```

	Recency	Frequency	MonetaryValue
count	3643.000000	3643.000000	3643.000000
mean	3.806481	2.171902	4.934900
std	1.352631	1.210321	1.310945
min	0.000000	0.000000	-0.430783
25%	2.944439	1.386294	4.072524
50%	3.931826	2.197225	4.915372
75%	4.934474	3.044522	5.812188
max	5.899897	7.311218	10.780213

```
1 # Create a vertical layout of three subplots
2 plt.figure(figsize=(5, 10))
3 plt.subplot(3, 1, 1)
4 sns.distplot(datamart_rfm_log['Recency']) # Plot distribution of log-transformed 'Recency'
5 plt.subplot(3, 1, 2)
6 sns.distplot(datamart_rfm_log['Frequency']) # Plot distribution of log-transformed 'Frequency'
7 plt.subplot(3, 1, 3)
8 sns.distplot(datamart_rfm_log['MonetaryValue']) # Plot distribution of log-transformed 'MonetaryValue'
9
10 # Show the plot
11 plt.show()
```



```
1 from sklearn.preprocessing import StandardScaler
2 # Initialize a StandardScaler
3 scaler = StandardScaler()
4
5 # Fit the StandardScaler to the log-transformed data
6 scaler.fit(datamart_rfm_log)
7
8 # Normalize and center the data using the fitted scaler
9 datamart_rfm_normalized = scaler.transform(datamart_rfm_log)
10
11 # Create a DataFrame with normalized values
12 datamart_rfm_normalized = pd.DataFrame(datamart_rfm_normalized,
13                                       index=datamart_rfm_log.index,
14                                       columns=datamart_rfm_log.columns)
15
16 # Statistical summary of the new DataFrame with normalized variables
17 datamart_rfm_normalized.describe().round(2)
18
```

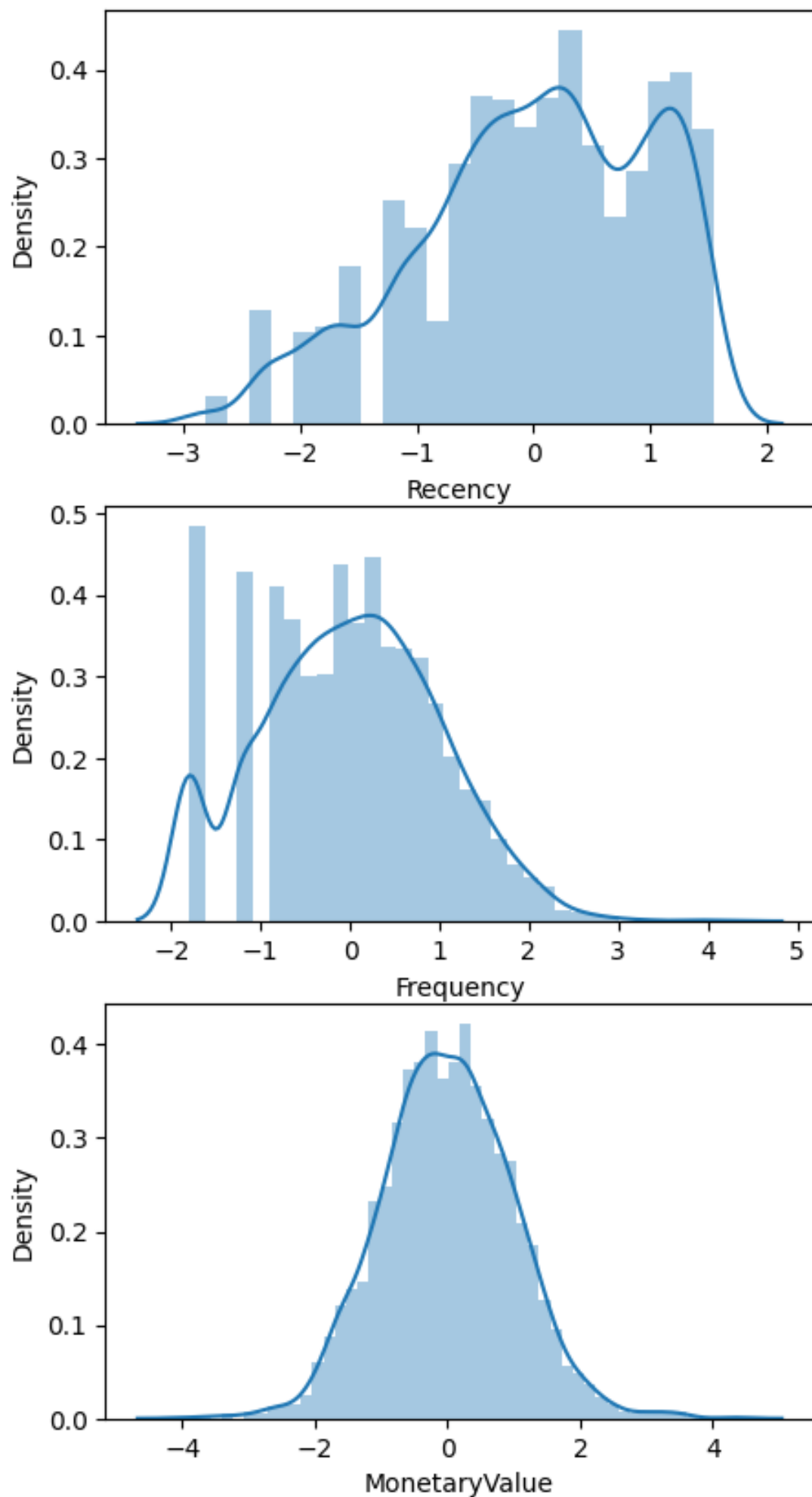
	Recency	Frequency	MonetaryValue
count	3643.00	3643.00	3643.00
mean	-0.00	0.00	0.00
std	1.00	1.00	1.00
min	-2.81	-1.79	-4.09
25%	-0.64	-0.65	-0.66
50%	0.09	0.02	-0.01
75%	0.83	0.72	0.67
max	1.55	4.25	4.46



```

1 # Create a vertical layout of three subplots
2 plt.figure(figsize=(5, 10))
3 plt.subplot(3, 1, 1)
4 sns.distplot(datamart_rfm_normalized['Recency']) # Plot distribution of normalized 'Recency'
5 plt.subplot(3, 1, 2)
6 sns.distplot(datamart_rfm_normalized['Frequency']) # Plot distribution of normalized 'Frequency'
7 plt.subplot(3, 1, 3)
8 sns.distplot(datamart_rfm_normalized['MonetaryValue']) # Plot distribution of normalized 'MonetaryValue'
9
10 # Show the plot
11 plt.show()

```



▼ Model Selection and Training

```

1 from sklearn.model_selection import train_test_split
2 # Add the 'Churn' column back to the DataFrame
3 datamart_rfm_normalized['Churn'] = datamart['Churn']
4
5 # Separate the target variable (Churn) from the features
6 X = datamart_rfm_normalized.drop('Churn', axis=1)
7 y = datamart_rfm_normalized['Churn']
8
9 # Split the data
10 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
11

```

```

1 from sklearn.linear_model import LogisticRegression
2 from sklearn.tree import DecisionTreeClassifier
3 from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier

```



```

4 from sklearn.model_selection import GridSearchCV
5
6 # Define the models and their respective hyperparameter grids for grid search
7 models = {
8     'Logistic Regression': (LogisticRegression(), {'C': [0.01, 0.1, 1, 10]}),
9     'Decision Tree': (DecisionTreeClassifier(), {'max_depth': [None, 5, 10, 20]}),
10    'Random Forest': (RandomForestClassifier(), {'n_estimators': [50, 100, 200]}),
11    'Gradient Boosting': (GradientBoostingClassifier(), {'n_estimators': [50, 100, 200]})
12 }
13
14 # Train and tune the models using grid search
15 best_models = {}
16 for model_name, (model, param_grid) in models.items():
17     grid_search = GridSearchCV(model, param_grid, cv=5, scoring='accuracy')
18     grid_search.fit(X_train, y_train)
19     best_models[model_name] = grid_search.best_estimator_
20
21 # Print the best hyperparameters for each model
22 for model_name, best_model in best_models.items():
23     print(f'Best hyperparameters in {model_name}: {best_model.get_params()}')
24

```

```

Best hyperparameters in Logistic Regression: {'C': 0.01, 'class_weight': None, 'dual': False, 'fit_intercept': True, 'intercept_scaling': 1,
Best hyperparameters in Decision Tree: {'ccp_alpha': 0.0, 'class_weight': None, 'criterion': 'gini', 'max_depth': 10, 'max_features': None,
Best hyperparameters in Random Forest: {'bootstrap': True, 'ccp_alpha': 0.0, 'class_weight': None, 'criterion': 'gini', 'max_depth': None,
Best hyperparameters in Gradient Boosting: {'ccp_alpha': 0.0, 'criterion': 'friedman_mse', 'init': None, 'learning_rate': 0.1, 'loss': 'log

```

```

1 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, roc_auc_score
2
3 # Evaluate each model's performance on the testing set
4 for model_name, model in best_models.items():
5     y_pred = model.predict(X_test)
6     accuracy = accuracy_score(y_test, y_pred)
7     precision = precision_score(y_test, y_pred)
8     recall = recall_score(y_test, y_pred)
9     f1 = f1_score(y_test, y_pred)
10    roc_auc = roc_auc_score(y_test, model.predict_proba(X_test)[:, 1])
11
12    print(f'{model_name} Performance:')
13    print(f'Accuracy: {accuracy:.2f}, Precision: {precision:.2f}, Recall: {recall:.2f}, F1-score:

```

```

Logistic Regression Performance:
Accuracy: 0.95, Precision: 0.95, Recall: 0.96, F1-score: 0.95, ROC-AUC: 0.99
Decision Tree Performance:
Accuracy: 1.00, Precision: 1.00, Recall: 1.00, F1-score: 1.00, ROC-AUC: 1.00
Random Forest Performance:
Accuracy: 1.00, Precision: 1.00, Recall: 1.00, F1-score: 1.00, ROC-AUC: 1.00
Gradient Boosting Performance:
Accuracy: 1.00, Precision: 1.00, Recall: 1.00, F1-score: 1.00, ROC-AUC: 1.00

```

```

1 # Create a DataFrame with the RFM values for the new customer
2 new_customer_rfm = pd.DataFrame({
3     'Recency': [2.5],
4     'Frequency': [1.0],
5     'MonetaryValue': [-0.5]
6 })
7
8 # Scale the RFM values using the fitted scaler
9 new_customer_rfm_scaled = scaler.transform(new_customer_rfm)
10
11 # Create a DataFrame with the scaled RFM values
12 new_customer_rfm_scaled = pd.DataFrame(new_customer_rfm_scaled, columns=new_customer_rfm.columns)
13
14 # Use the selected model to make the churn prediction
15 prediction = best_models['Random Forest'].predict(new_customer_rfm_scaled)
16
17 # Display the prediction (0: Not churn, 1: Churn)
18 if prediction[0] == 0:

```

```
19     print("Prediction: Not Churn")
20 else:
21     print("Prediction: Churn")
22
Prediction: Churn
```

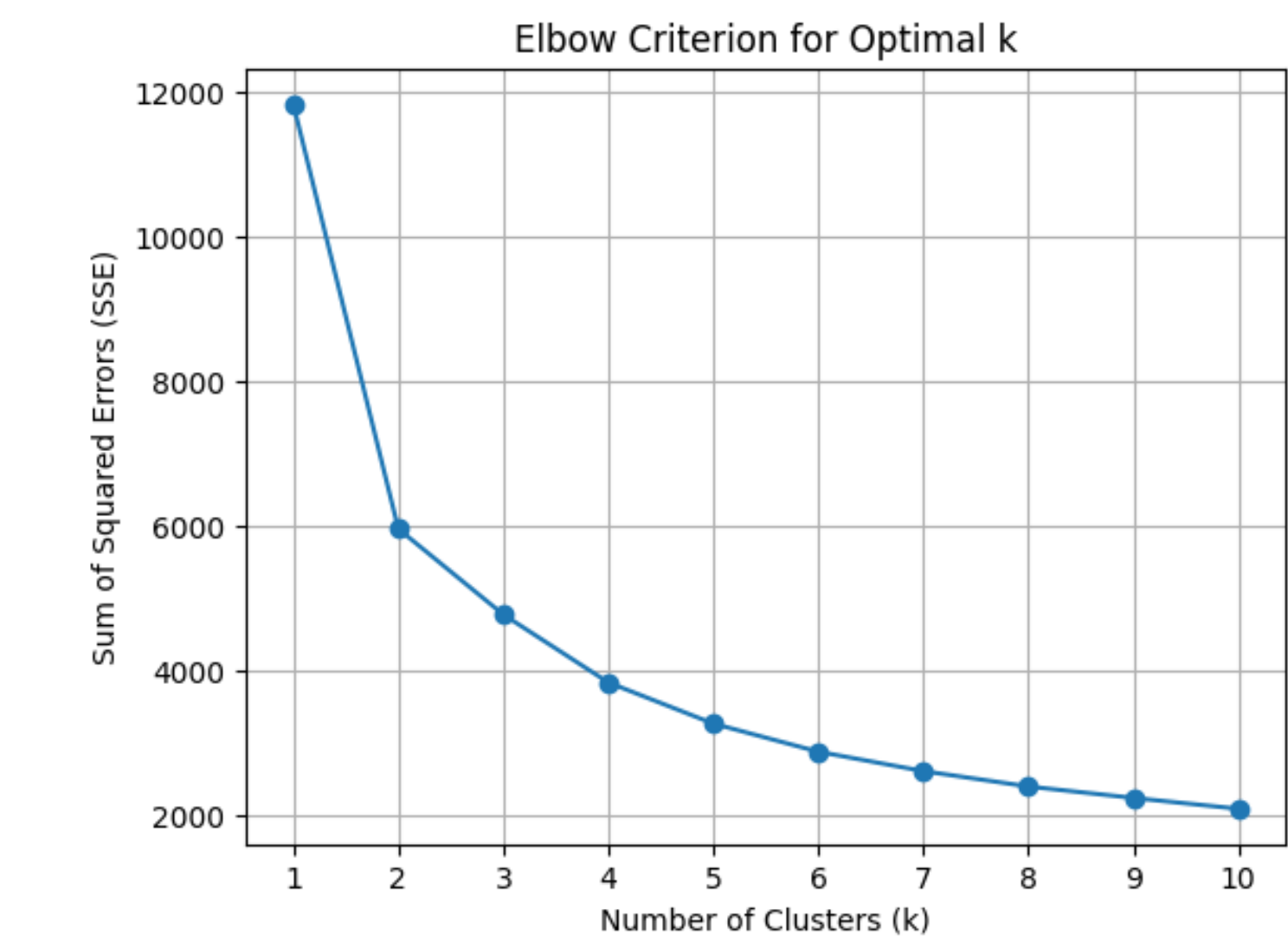
K_MEANS

The Elbow criterion is a widely used method to determine the optimal number of clusters (k) in K-Means clustering. It involves plotting the number of clusters against the Sum of Squared Errors (SSE) within the cluster.

Steps to find the optimal number of clusters using the Elbow criterion:

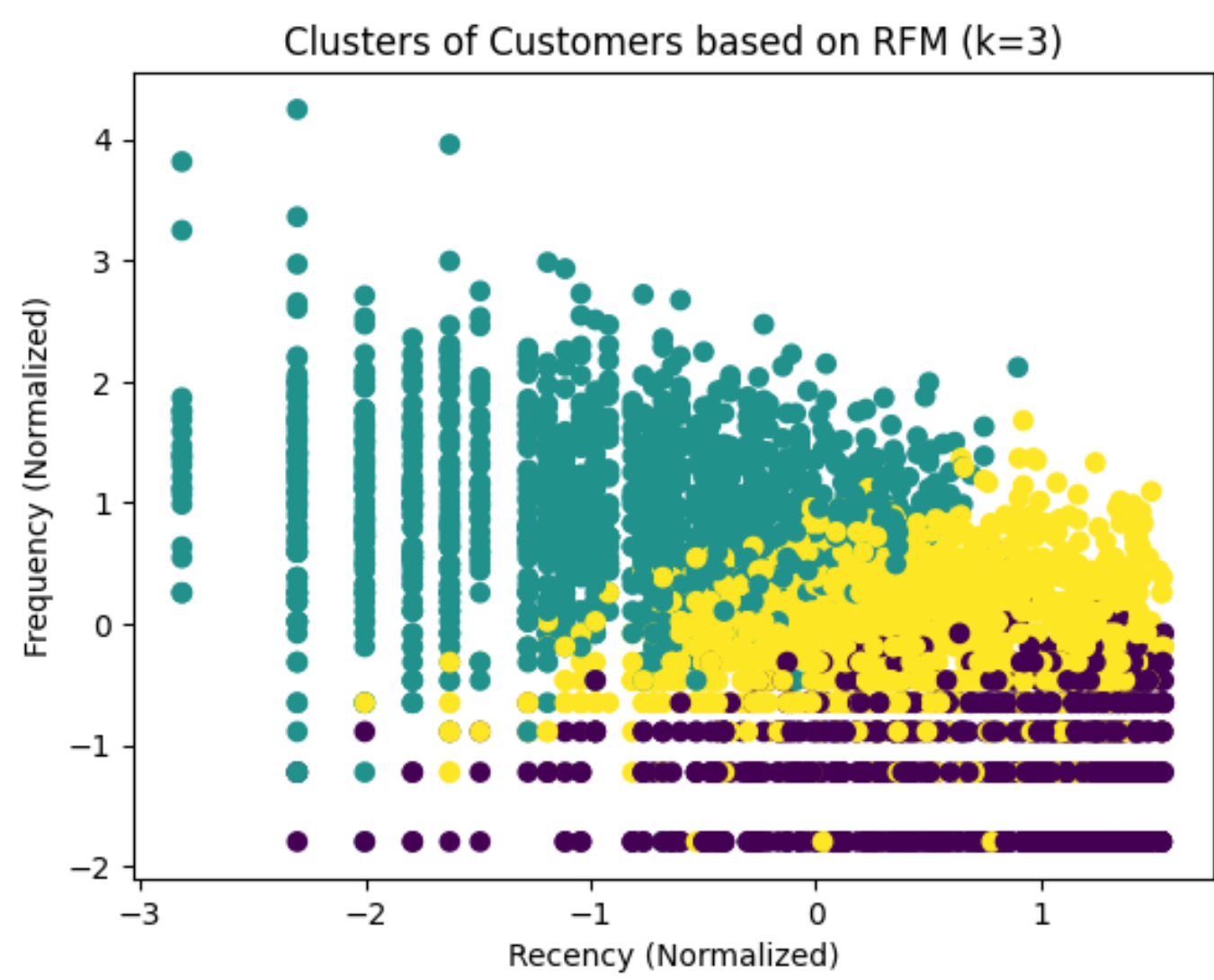
Calculate the SSE for different values of k (number of clusters). Plot the number of clusters (k) against the SSE. Identify the "elbow point" on the plot where the SSE decreases at a slower rate, indicating diminishing returns in reducing SSE with more clusters. Define the "elbow point" as the optimal value of k for K-Means.

```
1 from sklearn.cluster import KMeans
2
3 # Create an empty list to store SSE values for different k values
4 sse = []
5
6 # Try different values of k from 1 to 10
7 for k in range(1, 11):
8     # Create a KMeans instance with k clusters
9     kmeans = KMeans(n_clusters=k, random_state=42)
10
11     # Fit the KMeans model to the normalized data
12     kmeans.fit(datamart_rfm_normalized)
13
14     # Append the SSE value for the current k to the list
15     sse.append(kmeans.inertia_)
16
17 # Plot the number of clusters (k) against the SSE
18 plt.plot(range(1, 11), sse, marker='o')
19 plt.xlabel('Number of Clusters (k)')
20 plt.ylabel('Sum of Squared Errors (SSE)')
21 plt.title('Elbow Criterion for Optimal k')
22 plt.xticks(range(1, 11))
23 plt.grid(True)
24 plt.show()
25
```



```
1 # Importing required libraries
2 import matplotlib.pyplot as plt
3
4 # Create a KMeans instance with k=3 clusters (as per the elbow method result)
```

```
5 kmeans = KMeans(n_clusters=3, random_state=42)
6
7 # Fit the KMeans model to the normalized data
8 kmeans.fit(datamart_rfm_normalized)
9
10 # Get the cluster labels for each data point
11 cluster_labels = kmeans.labels_
12
13 # Plot the data points using Recency and Frequency as x and y axes, respectively
14 plt.scatter(datamart_rfm_normalized.iloc[:, 0], datamart_rfm_normalized.iloc[:, 1], c=cluster_labels)
15 plt.xlabel('Recency (Normalized)')
16 plt.ylabel('Frequency (Normalized)')
17 plt.title('Clusters of Customers based on RFM (k=3)')
18 plt.show()
19
```

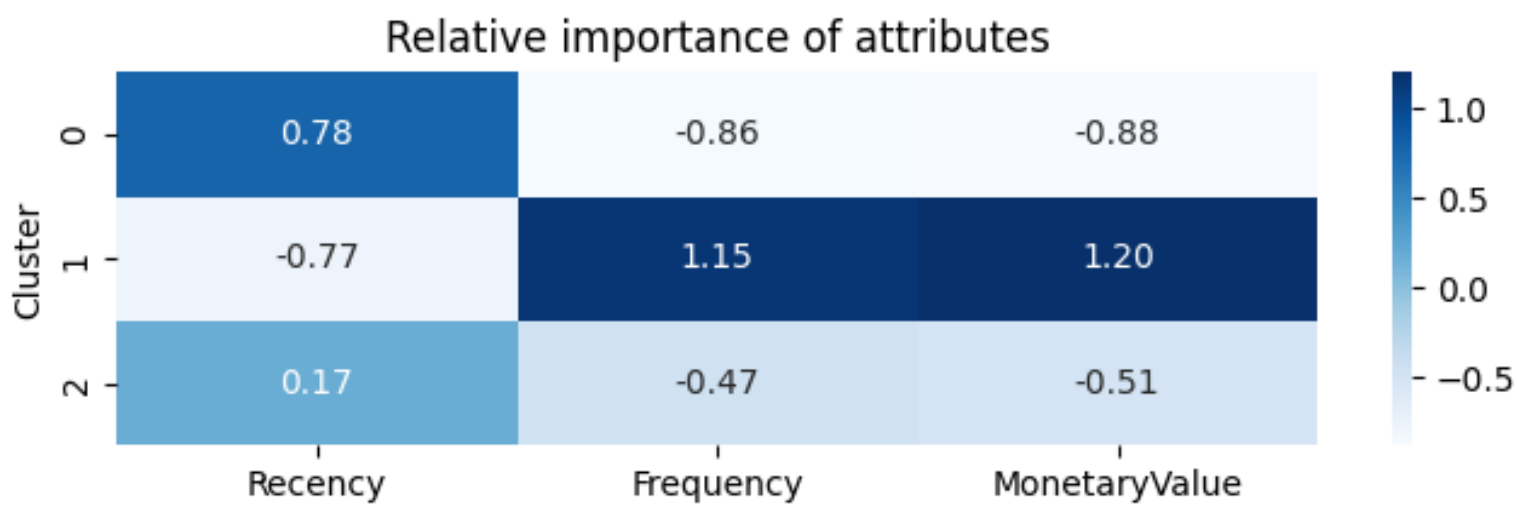


```
1 # Create a DataFrame 'datamart_rfm_k3' with the added 'Cluster' column
2 datamart_rfm_k3 = datamart_rfm.assign(Cluster=cluster_labels)
3
4 # Calculate the average RFM values and the number of customers for each cluster
5 cluster_summary = datamart_rfm_k3.groupby(['Cluster']).agg({
6     'Recency': 'mean',
7     'Frequency': 'mean',
8     'MonetaryValue': ['mean', 'count']
9 }).round(0)
10
11 # Display the summary statistics for each cluster
12 cluster_summary
13
```

	Recency	Frequency	MonetaryValue	
	mean	mean	mean	count
Cluster				
0	161.0	3.0	44.0	970
1	21.0	40.0	817.0	1296
2	106.0	10.0	181.0	1377

```
1 # Calculate the average RFM values for each cluster
2 cluster_avg = datamart_rfm_k3.groupby(['Cluster']).mean()
3
4 # Calculate the average RFM values for the total population of customers
5 population_avg = datamart_rfm.mean()
6
7 # Calculate the relative importance of each attribute's value in the cluster compared to the population
8 relative_imp = cluster_avg / population_avg - 1
```

```
9
10 # Round the relative importance scores to 2 decimal places
11 relative_imp_rounded = relative_imp.round(2)
12
13 # Plot the heatmap
14 plt.figure(figsize=(8, 2))
15 plt.title('Relative importance of attributes')
16 sns.heatmap(data=relative_imp_rounded, annot=True, fmt='.2f', cmap='Blues')
17
18 # Show the plot
19 plt.show()
20
```



► **check and concat**

[] ↳ 6 cells hidden