

```
In [1]: 1 import pandas as pd
2 import numpy as np
3 #viz Libraries
4 import matplotlib.pyplot as plt
5 plt.style.use('ggplot')
6 import seaborn as sns
7 #datetime
8 import datetime as dt
9 #StandardScaler
10 from sklearn.preprocessing import StandardScaler
11 #KMeans
12 from sklearn.cluster import KMeans
```

```
In [11]: 1 df = pd.read_csv(r"C:\Users\kavet\Desktop\BE_Sem_7_Assignments-main\BE_Ser
2
```

```
1
2 df.shape #Dimensions of the data
3
4
```

```
In [12]: 1 df.head() #Glimpse of the data
2
3
```

Out[12]:

	ORDERNUMBER	QUANTITYORDERED	PRICEEACH	ORDERLINENUMBER	SALES	ORDERDA
0	10107	30	95.70	2	2871.00	2/24/2000
1	10121	34	81.35	5	2765.90	5/7/2003
2	10134	41	94.74	2	3884.34	7/1/2003
3	10145	45	83.26	6	3746.70	8/25/2000
4	10159	49	100.00	14	5205.27	10/10/2000

5 rows × 7 columns

```
In [13]: 1 #Removing the variables which dont add significant value for the analysis.
2 to_drop = ['PHONE', 'ADDRESSLINE1', 'ADDRESSLINE2', 'STATE', 'POSTALCODE']
3 df = df.drop(to_drop, axis=1)
4
5 df.isnull().sum()
6
7
```

```
Out[13]: ORDERNUMBER      0
QUANTITYORDERED      0
PRICEEACH            0
ORDERLINENUMBER      0
SALES                0
```

```
In [13]: 1 #Removing the variables which dont add significant value fot the analysis.
2 to_drop = ['PHONE', 'ADDRESSLINE1', 'ADDRESSLINE2', 'STATE', 'POSTALCODE']
3 df = df.drop(to_drop, axis=1)
4
5 df.isnull().sum()
6
7
```

```
Out[13]: ORDERNUMBER      0
QUANTITYORDERED      0
PRICEEACH             0
ORDERLINENUMBER      0
SALES                 0
ORDERDATE             0
STATUS                0
QTR_ID                0
MONTH_ID              0
YEAR_ID               0
PRODUCTLINE           0
MSRP                  0
PRODUCTCODE           0
CUSTOMERNAME          0
CITY                  0
COUNTRY               0
TERRITORY             1074
CONTACTLASTNAME       0
CONTACTFIRSTNAME      0
DEALSIZE              0
dtype: int64
```

```
In [14]: 1 df.dtypes
2
3 df['ORDERDATE'] = pd.to_datetime(df['ORDERDATE'])
4
5 df['ORDERDATE'] = [d.date() for d in df['ORDERDATE']]
6 df.head()
```

```
Out[14]:
```

	ORDERNUMBER	QUANTITYORDERED	PRICEEACH	ORDERLINENUMBER	SALES	ORDERDATE
0	10107	30	95.70	2	2871.00	2003-02-05
1	10121	34	81.35	5	2765.90	2003-05-02
2	10134	41	94.74	2	3884.34	2003-07-08
3	10145	45	83.26	6	3746.70	2003-08-02
4	10159	49	100.00	14	5205.27	2003-10-05

```
In [15]: 1 # Calculate Recency, Frequency and Monetary value for each customer
2 latest_date = df['ORDERDATE'].max() + dt.timedelta(days=1) #Latest date in data
3 df_RFM = df.groupby(['CUSTOMERNAME'])
4
5 df_RFM = df_RFM.agg({
6     'ORDERDATE': lambda x: (latest_date - x.max()).days,
7     'ORDERNUMBER': 'count',
8     'SALES': 'sum'})
```

```
In [16]: 1 #Renaming the columns
In [17]: 2 df_RFM.rename(columns={'ORDERDATE': 'Recency', 'ORDERNUMBER': 'Frequency',
3     'SALES': 'MonetaryValue'})
4 data.head()
```

```
Out[17]:
```

	Recency	Frequency	MonetaryValue
CUSTOMERNAME			
AV Stores, Co.	196	51	157807.81
Alpha Cognac	65	20	70488.44
Amica Models & Co.	265	26	94117.26
Anna's Decorations, Ltd	84	46	153996.13

```

In [16]: 1 #Renaming the columns
In [17]: 2 data.rename(columns={'ORDERNUMBER': 'Frequency', 'MonetaryValue': 'Recency'})
          3 data.head()
          4
Out[17]:

```

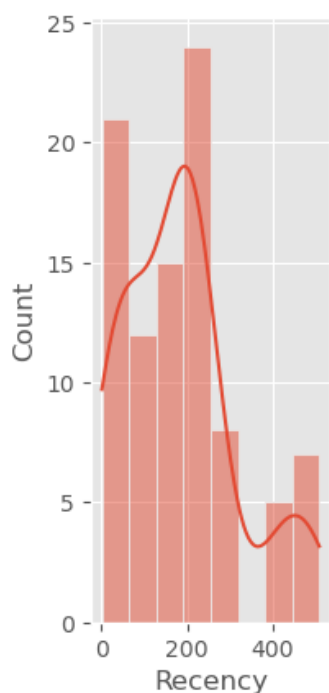
	Recency	Frequency	MonetaryValue
CUSTOMERNAME			
AV Stores, Co.	196	51	157807.81
Alpha Cognac	65	20	70488.44
Amica Models & Co.	265	26	94117.26
Anna's Decorations, Ltd	84	46	153996.13
Atelier graphique	188	7	24179.96

```

In [18]: 1 plt.subplot(1,3,1)
          2 sns.histplot(data['Recency'], kde=True)

```

Out[18]: <Axes: xlabel='Recency', ylabel='Count'>

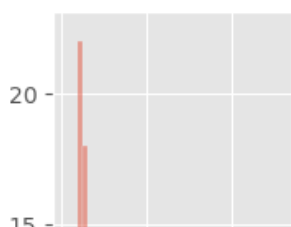


```

In [19]: 1 plt.subplot(1,3,2)
          2 sns.histplot(data['Frequency'], kde=True)

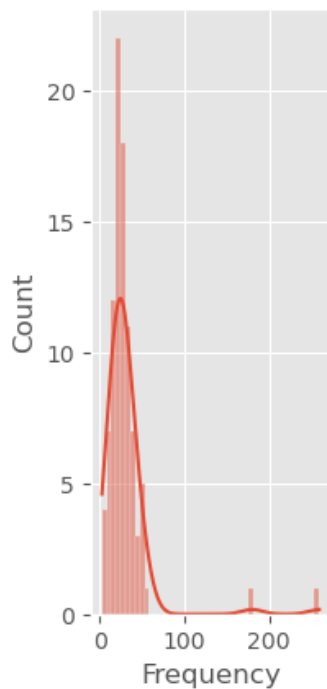
```

Out[19]: <Axes: xlabel='Frequency', ylabel='Count'>



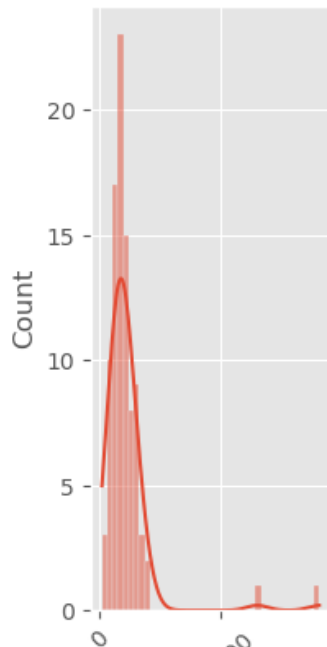
```
In [19]: 1 plt.subplot(1,3,2)
2 sns.histplot(data['Frequency'], kde=True)
```

Out[19]: <Axes: xlabel='Frequency', ylabel='Count'>



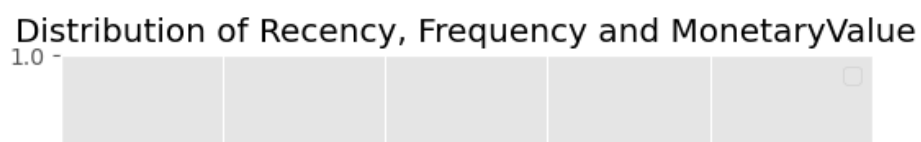
```
In [20]: 1 plt.subplot(1,3,3)
2 plt.xticks(rotation = 45)
3 sns.histplot(data['MonetaryValue'], kde=True)
```

Out[20]: <Axes: xlabel='MonetaryValue', ylabel='Count'>



```
In [21]: 1 plt.title('Distribution of Recency, Frequency and MonetaryValue')
2 plt.legend()
3 plt.show()
```

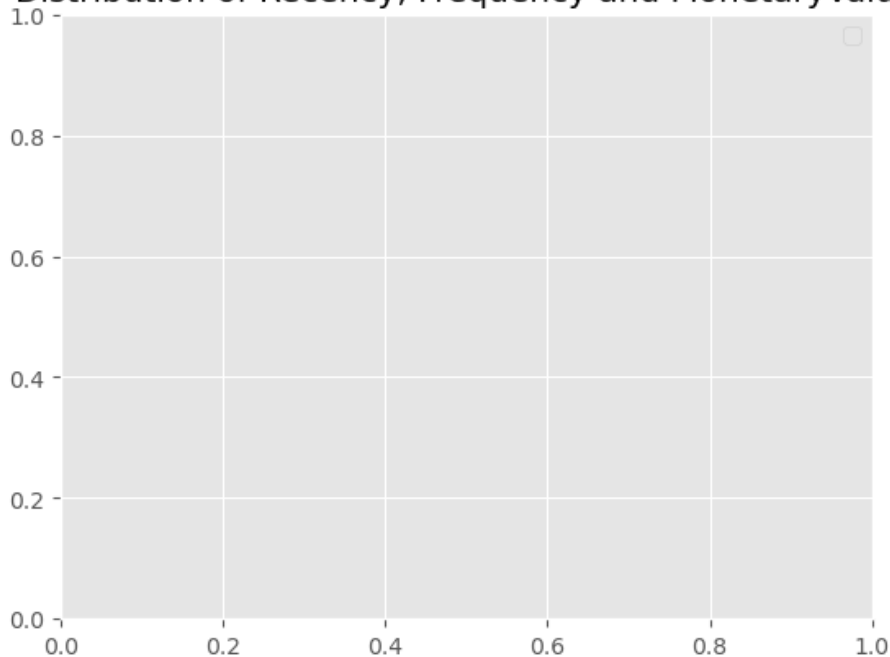
No artists with labels found to put in legend. Note that artists whose label start with an underscore are ignored when legend() is called with no argument.



```
In [21]: 1 plt.title('Distribution of Recency, Frequency and MonetaryValue')
2         plt.legend()
3         plt.show()
```

No artists with labels found to put in legend. Note that artists whose label start with an underscore are ignored when legend() is called with no argument.

Distribution of Recency, Frequency and MonetaryValue



```
In [22]: 1 data_log = np.log(data)
```

```
In [23]: 1 data_log.head()
```

Out[23]:

	Recency	Frequency	MonetaryValue
CUSTOMERNAME			
AV Stores, Co.	5.278115	3.931826	11.969133
Alpha Cognac	4.174387	2.995732	11.163204
Amica Models & Co.	5.579730	3.258097	11.452297
Anna's Decorations, Ltd	4.430817	3.828641	11.944683
Atelier graphique	5.236442	1.945910	10.093279

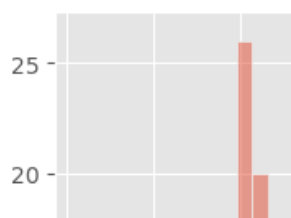
```
In [24]: 1 plt.figure(figsize=(10,6))
```

Out[24]: <Figure size 1000x600 with 0 Axes>

<Figure size 1000x600 with 0 Axes>

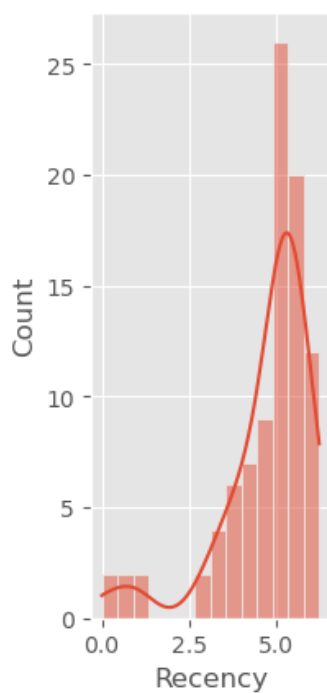
```
In [25]: 1 plt.subplot(1,3,1)
2         sns.histplot(data_log['Recency'], kde=True)
```

Out[25]: <Axes: xlabel='Recency', ylabel='Count'>



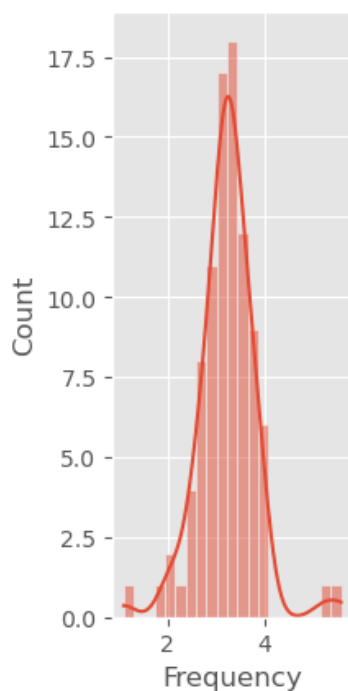
```
In [25]: 1 plt.subplot(1,3,1)
        2 sns.histplot(data_log['Recency'], kde=True)
```

Out[25]: <Axes: xlabel='Recency', ylabel='Count'>



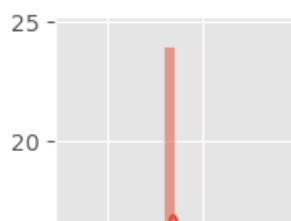
```
In [26]: 1 plt.subplot(1,3,2)
        2 sns.histplot(data_log['Frequency'], kde=True)
```

Out[26]: <Axes: xlabel='Frequency', ylabel='Count'>



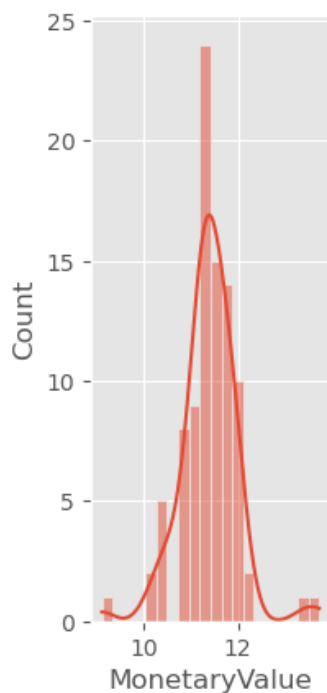
```
In [27]: 1 plt.subplot(1,3,3)
        2 sns.histplot(data_log['MonetaryValue'], kde=True)
```

Out[27]: <Axes: xlabel='MonetaryValue', ylabel='Count'>



```
frequency
In [27]: 1 plt.subplot(1,3,3)
        2 sns.histplot(data_log['MonetaryValue'], kde=True)
```

Out[27]: <Axes: xlabel='MonetaryValue', ylabel='Count'>

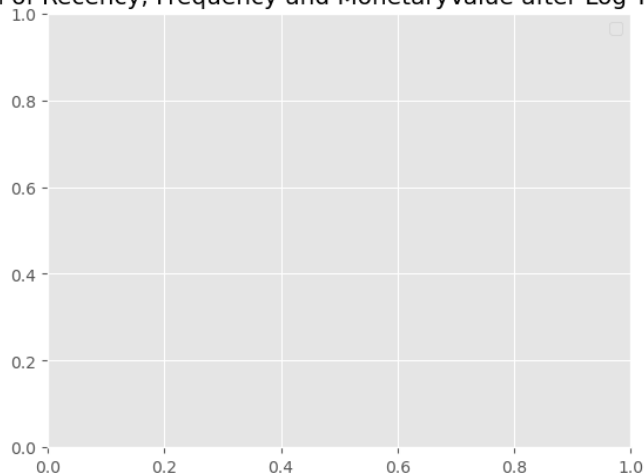


```
In [29]: 1 plt.title('Distribution of Recency, Frequency and MonetaryValue after Log
        2 plt.legend()
        3
```

No artists with labels found to put in legend. Note that artists whose label start with an underscore are ignored when legend() is called with no argument.

Out[29]: <matplotlib.legend.Legend at 0x1e3e66fce90>

Distribution of Recency, Frequency and MonetaryValue after Log Transformation



```
In [30]: 1 plt.show()
        2 # Fit the scaler
        3 scaler.fit(data_log)
In [31]: 1 # Initialize a scaler
Out[32]: StandardScaler()
```

StandardScaler()

```
In [33]: 1 # Scale and center the data
        2 data_normalized = scaler.transform(data_log)
```

```
In [34]: 1 # Create a pandas DataFrame
```

```
In [30]: 1 plt.show()
          # Fit the scaler
          2 scaler.fit(data_log)

In [31]: 1 # Initialize a scaler
Out[32]: 2 scaler = StandardScaler()
          StandardScaler()
```

```
In [33]: 1 # Scale and center the data
          2 data_normalized = scaler.transform(data_log)
```

```
In [34]: 1 # Create a pandas DataFrame
          2 data_normalized = pd.DataFrame(data_normalized, index=data_log.index, columns=
```

```
In [35]: 1 # Print summary statistics
          2 data_normalized.describe().round(2)
```

```
Out[35]:
```

	Recency	Frequency	MonetaryValue
count	92.00	92.00	92.00
mean	0.00	-0.00	0.00
std	1.01	1.01	1.01
min	-3.51	-3.67	-3.82
25%	-0.24	-0.41	-0.39
50%	0.37	0.06	-0.04
75%	0.53	0.45	0.52
max	1.12	4.03	3.92

```
In [36]: 1 # Fit KMeans and calculate SSE for each k
          2 sse={}
          3 for k in range(1, 21):
          4     kmeans = KMeans(n_clusters=k, random_state=1)
          5     kmeans.fit(data_normalized)
          6     sse[k] = kmeans.inertia_

ans.py:1412: FutureWarning: The default value of `n_init` will change from
10 to 'auto' in 1.4. Set the value of `n_init` explicitly to suppress the
warning
  super()._check_params_vs_input(X, default_n_init=10)
C:\Users\kavet\anaconda3\New folder\Lib\site-packages\sklearn\cluster\_kme
ans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows
with MKL, when there are less chunks than available threads. You can avoid
it by setting the environment variable OMP_NUM_THREADS=1.
  warnings.warn(
C:\Users\kavet\anaconda3\New folder\Lib\site-packages\sklearn\cluster\_kme
ans.py:1412: FutureWarning: The default value of `n_init` will change from
10 to 'auto' in 1.4. Set the value of `n_init` explicitly to suppress the
warning
  super()._check_params_vs_input(X, default_n_init=10)
C:\Users\kavet\anaconda3\New folder\Lib\site-packages\sklearn\cluster\_kme
ans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows
with MKL, when there are less chunks than available threads. You can avoid
it by setting the environment variable OMP_NUM_THREADS=1.
  warnings.warn(
C:\Users\kavet\anaconda3\New folder\Lib\site-packages\sklearn\cluster\_kme
```

```
In [37]: 1 plt.figure(figsize=(10,6))
          2 plt.title('The Elbow Method')
```

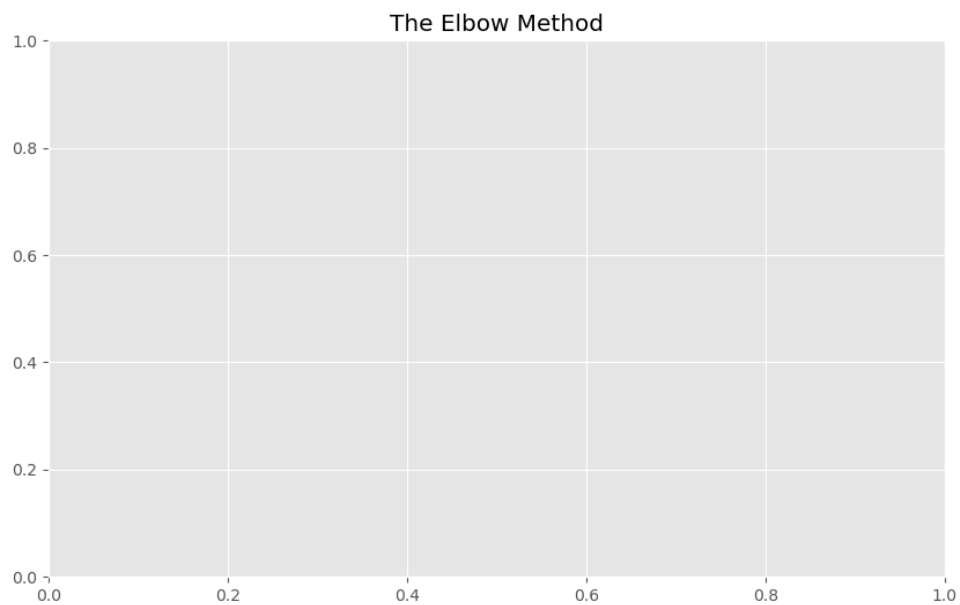
```
Out[37]: Text(0.5, 1.0, 'The Elbow Method')
```





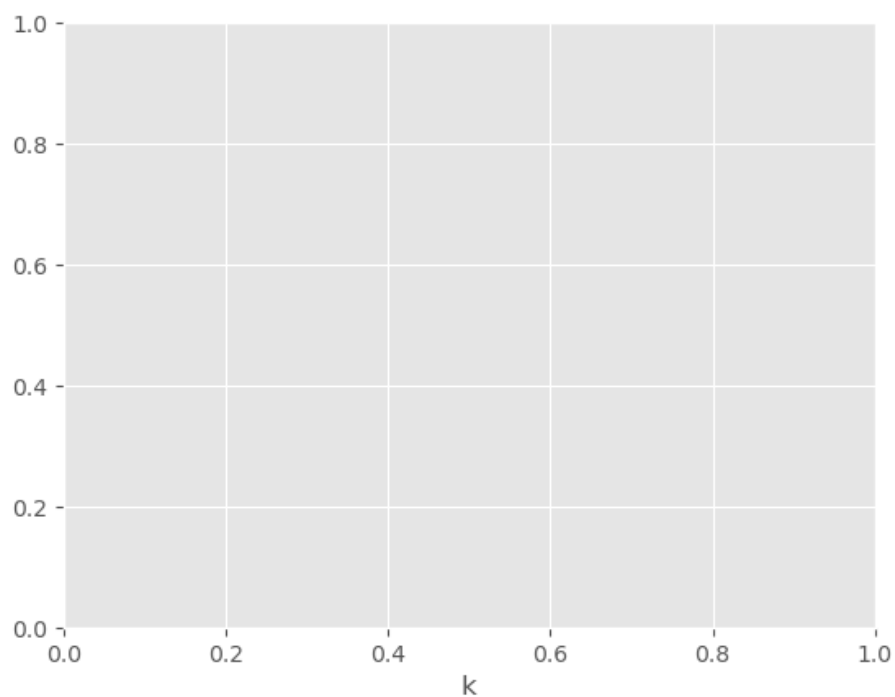
```
In [37]: 1 plt.figure(figsize=(10,6))
        2 plt.title('The Elbow Method')
```

Out[37]: Text(0.5, 1.0, 'The Elbow Method')



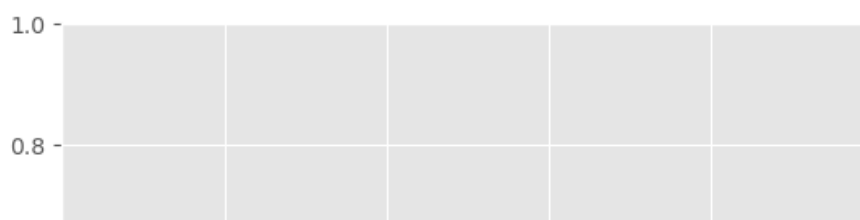
```
In [38]: 1 # Add X-axis Label "k"
        2 plt.xlabel('k')
        3
```

Out[38]: Text(0.5, 0, 'k')



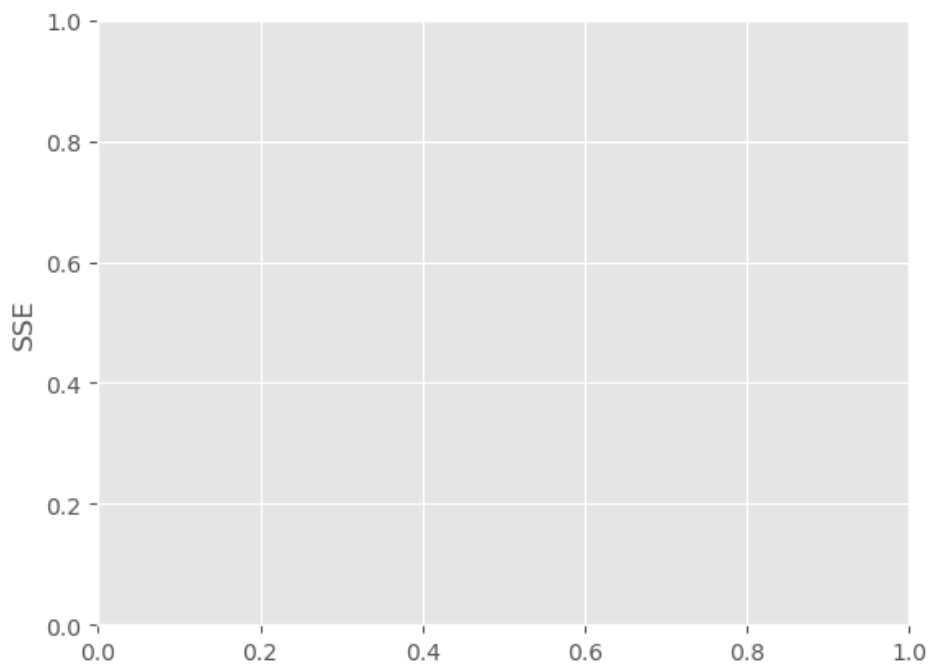
```
In [39]: 1 # Add Y-axis Label "SSE"
        2 plt.ylabel('SSE')
```

Out[39]: Text(0, 0.5, 'SSE')

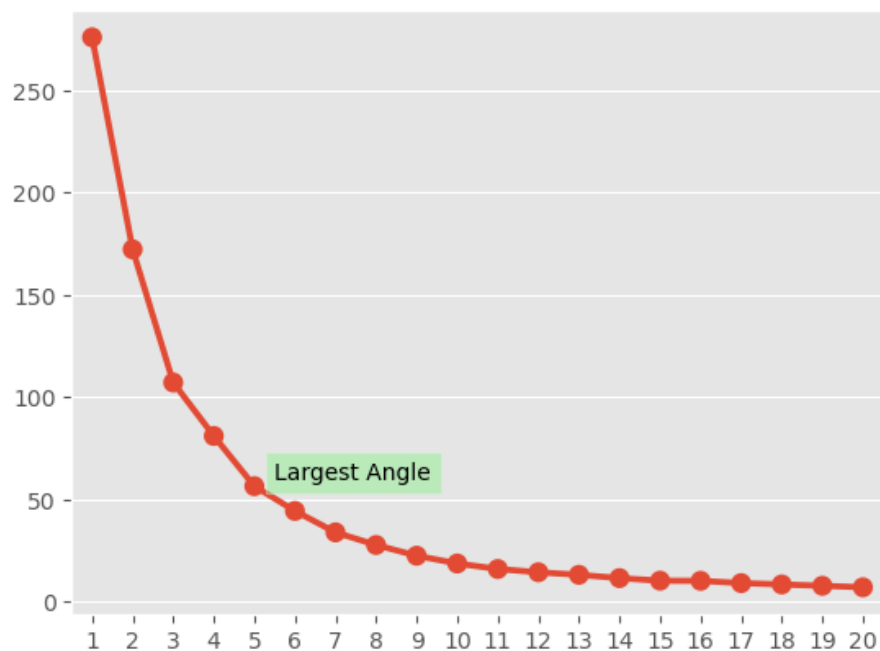


```
In [39]: 1 # Add Y-axis Label "SSE"
          2 plt.ylabel('SSE')
```

Out[39]: Text(0, 0.5, 'SSE')



```
In [40]: 1 # Plot SSE values for each key in the dictionary
          2 sns.pointplot(x=list(sse.keys()), y=list(sse.values()))
          3 plt.text(4.5,60,"Largest Angle",bbox=dict(facecolor='lightgreen', alpha=0
          4 plt.show())
```



```
In [41]: 1 # Running KMeans with 5 clusters on the normalized data set
          2 # Initializing KMeans
          3 kmeans = KMeans(n_clusters=5, random_state=1)
```

C:\Users\kavet\anaconda3\New folder\Lib\site-packages\sklearn\cluster\\_kmean  
s.py:1412: FutureWarning: The default value of `n\_init` will change from 10 t  
o 'auto' in 1.4. Set the value of `n\_init` explicitly to suppress the warning  
super().\_check\_params\_vs\_input(X, default\_n\_init=10)

C:\Users\kavet\anaconda3\New folder\Lib\site-packages\sklearn\cluster\\_kmean  
s.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with  
MKL, when there are less chunks than available threads. You can avoid it by s  
etting the environment variable OMP\_NUM\_THREADS=1.  
warnings.warn(

```
In [42]: 1 # Run KMeans with 5 clusters on the normalized data set
2 # Fit the KMeans model
3 kmeans = KMeans(n_clusters=5, random_state=1)

C:\Users\kavet\anaconda3\New folder\Lib\site-packages\sklearn\cluster\_kmean
s.py:1412: FutureWarning: The default value of `n_init` will change from 10 t
o 'auto' in 1.4. Set the value of `n_init` explicitly to suppress the warning
  super()._check_params_vs_input(X, default_n_init=10)
C:\Users\kavet\anaconda3\New folder\Lib\site-packages\sklearn\cluster\_kmean
s.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with
MKL, when there are less chunks than available threads. You can avoid it by s
etting the environment variable OMP_NUM_THREADS=1.
  warnings.warn(
```

```
Out[42]: KMeans
KMeans(n_clusters=5, random_state=1)
```

```
In [43]: 1 # Extract cluster labels
2 cluster_labels = kmeans.labels_
```

```
In [44]: 1 # Assigning Cluster Labels to Raw Data
2 # Create a DataFrame by adding a new cluster label column
3 data_rfm = data.assign(Cluster=cluster_labels)
4 data_rfm.head()
```

```
Out[44]:
```

	Recency	Frequency	MonetaryValue	Cluster
CUSTOMERNAME				
AV Stores, Co.	196	51	157807.81	1
Alpha Cognac	65	20	70488.44	2
Amica Models & Co.	265	26	94117.26	2
Anna's Decorations, Ltd	84	46	153996.13	1
Atelier graphique	188	7	24179.96	0

```
In [45]: 1 # Group the data by cluster
2 grouped = data_rfm.groupby(['Cluster'])
```

```
In [46]: 1 # Calculate average RFM values and segment sizes per cluster value
2 grouped.agg({
3     'Recency': 'mean',
4     'Frequency': 'mean',
5     'MonetaryValue': ['mean', 'count']
6 }).round(1)
```

```
Out[46]:
```

	Recency	Frequency	MonetaryValue	
Cluster	mean	mean	mean	count
0	324.2	10.7	35628.7	12
1	126.5	37.1	133158.0	31
2	209.2	22.1	78633.2	43
3	2.0	219.5	783576.1	2

```
In [47]: 1 #4-----2.0-----38.75-----132201.6-----4---#
2 # Calculating relative importance of each attribute
3 # Calculate average RFM values for each cluster
4 cluster_avg = data_rfm.groupby(['Cluster']).mean()
5 print(cluster_avg)
```

Cluster	Recency	Frequency	MonetaryValue
0	324.250000	10.666667	35628.653333
1	126.548387	37.129032	133158.014516
2	209.162791	22.093023	78633.205814
3	2.000000	219.500000	783576.085000
4	2.000000	38.750000	132201.635000

```
In [47]: 1 #4-----2.0-----35628.653333-----4---#
2 # # Calculating relative importance of each attribute
3 # Calculate average RFM values for each cluster
4 cluster_avg = data_rfm.groupby(['Cluster']).mean()
5 print(cluster_avg)
```

	Recency	Frequency	MonetaryValue
Cluster			
0	324.250000	10.666667	35628.653333
1	126.548387	37.129032	133158.014516
2	209.162791	22.093023	78633.205814
3	2.000000	219.500000	783576.085000
4	2.000000	38.750000	132201.635000

```
In [48]: 1 # Calculate average RFM values for the total customer population
2 population_avg = data.mean()
3 print(population_avg)
```

```
Recency      182.826087
Frequency     30.684783
MonetaryValue 109050.313587
dtype: float64
```

```
In [49]: 1 # Calculate relative importance of cluster's attribute value compared to p
2 relative_imp = cluster_avg / population_avg - 1
```

```
In [50]: 1 # Print relative importance score rounded to 2 decimals
2 print(relative_imp.round(2))
```

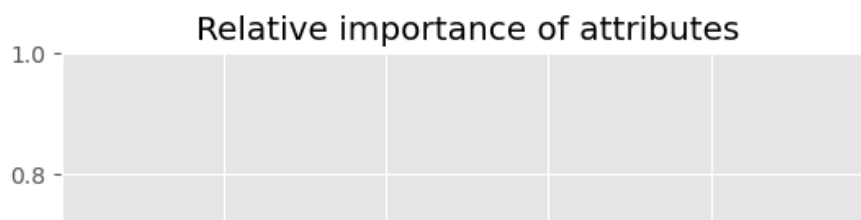
	Recency	Frequency	MonetaryValue
Cluster			
0	0.77	-0.65	-0.67
1	-0.31	0.21	0.22
2	0.14	-0.28	-0.28
3	-0.99	6.15	6.19
4	-0.99	0.26	0.21

```
In [51]: 1 #Plot Relative Importance
2 # Initialize a plot with a figure size of 8 by 2 inches
3 plt.figure(figsize=(8, 2))
```

```
Out[51]: <Figure size 800x200 with 0 Axes>
<Figure size 800x200 with 0 Axes>
```

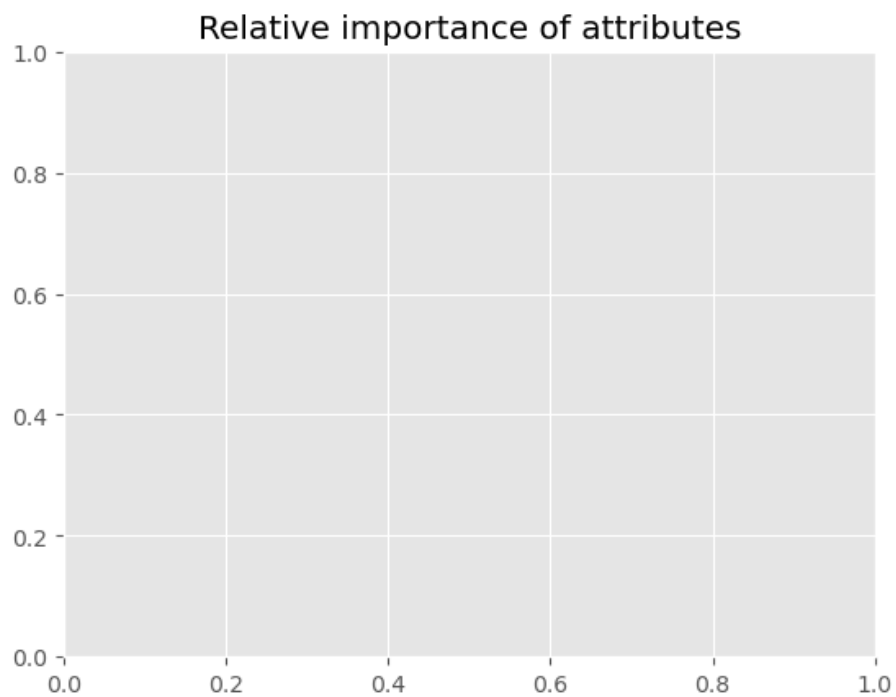
```
In [52]: 1 # Add the plot title
2 plt.title('Relative importance of attributes')
```

```
Out[52]: Text(0.5, 1.0, 'Relative importance of attributes')
```



```
In [52]: 1 # Add the plot title
         2 plt.title('Relative importance of attributes')
```

```
Out[52]: Text(0.5, 1.0, 'Relative importance of attributes')
```



```
In [53]: 1 # Plot the heatmap
         2 sns.heatmap(data=relative_imp, annot=True, fmt='.2f', cmap='RdYlGn')
         3 plt.show()
```

