In [1]:
```python
import pandas as pd
import numpy as np
```

In [3]:
```python
df = pd.read_csv('7431_sales_data_sample.csv', encoding='unicode_escape')
```

In [5]:
```python
df.head()
```

Out[5]:

| | ORDERNUMBER | QUANTITYORDERED | PRICEEACH | ORDERLINENUMBER | SALES | OR |
|---|---|---|---|---|---|---|
| **0** | 10107 | 30 | 95.70 | 2 | 2871.00 | |
| **1** | 10121 | 34 | 81.35 | 5 | 2765.90 | |
| **2** | 10134 | 41 | 94.74 | 2 | 3884.34 | |
| **3** | 10145 | 45 | 83.26 | 6 | 3746.70 | |
| **4** | 10159 | 49 | 100.00 | 14 | 5205.27 | |

5 rows × 25 columns

In [7]:
```python
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2823 entries, 0 to 2822
Data columns (total 25 columns):
 #   Column           Non-Null Count  Dtype
---  ------           --------------  -----
 0   ORDERNUMBER      2823 non-null   int64
 1   QUANTITYORDERED  2823 non-null   int64
 2   PRICEEACH        2823 non-null   float64
 3   ORDERLINENUMBER  2823 non-null   int64
 4   SALES            2823 non-null   float64
 5   ORDERDATE        2823 non-null   object
 6   STATUS           2823 non-null   object
 7   QTR_ID           2823 non-null   int64
 8   MONTH_ID         2823 non-null   int64
 9   YEAR_ID          2823 non-null   int64
 10  PRODUCTLINE      2823 non-null   object
 11  MSRP             2823 non-null   int64
 12  PRODUCTCODE      2823 non-null   object
 13  CUSTOMERNAME     2823 non-null   object
 14  PHONE            2823 non-null   object
 15  ADDRESSLINE1     2823 non-null   object
 16  ADDRESSLINE2     302 non-null    object
 17  CITY             2823 non-null   object
 18  STATE            1337 non-null   object
 19  POSTALCODE       2747 non-null   object
 20  COUNTRY          2823 non-null   object
 21  TERRITORY        1749 non-null   object
 22  CONTACTLASTNAME  2823 non-null   object
 23  CONTACTFIRSTNAME 2823 non-null   object
 24  DEALSIZE         2823 non-null   object
dtypes: float64(2), int64(7), object(16)
memory usage: 551.5+ KB
```

In [8]:
```python
#Columns to Remove
to_drop = ['ADDRESSLINE1', 'ADDRESSLINE2', 'STATE', 'POSTALCODE', 'PHONE']
df = df.drop(to_drop, axis=1)
```

In [9]:
```python
#Check for null values
df.isnull().sum()
```

Out[9]:     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 [10]:    ```python
            df.dtypes
            ```

Out[10]:    ORDERNUMBER          int64
            QUANTITYORDERED      int64
            PRICEEACH          float64
            ORDERLINENUMBER      int64
            SALES              float64
            ORDERDATE           object
            STATUS              object
            QTR_ID               int64
            MONTH_ID             int64
            YEAR_ID              int64
            PRODUCTLINE         object
            MSRP                 int64
            PRODUCTCODE         object
            CUSTOMERNAME        object
            CITY                object
            COUNTRY             object
            TERRITORY           object
            CONTACTLASTNAME     object
            CONTACTFIRSTNAME    object
            DEALSIZE            object
            dtype: object

In [12]:    ```python
            #ORDERDATE Should be in date time
            df['ORDERDATE'] = pd.to_datetime(df['ORDERDATE'])
            ```

In [13]:    ```python
            #We need to create some features in order to create cluseters
            #Recency: Number of days between customer's latest order and today's date
            #Frequency : Number of purchases by the customers
            #MonetaryValue : Revenue generated by the customers
            import datetime as dt
            snapshot_date = df['ORDERDATE'].max() + dt.timedelta(days = 1)
            df_RFM = df.groupby(['CUSTOMERNAME']).agg({
                'ORDERDATE' : lambda x : (snapshot_date - x.max()).days,
                'ORDERNUMBER' : 'count',
                'SALES' : 'sum'
            ```

```
})

#Rename the columns
df_RFM.rename(columns = {
    'ORDERDATE' : 'Recency',
    'ORDERNUMBER' : 'Frequency',
    'SALES' : 'MonetaryValue'
}, inplace=True)
```

In [14]: `df_RFM.head()`

Out[14]:

| CUSTOMERNAME | Recency | Frequency | MonetaryValue |
|---|---|---|---|
| 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 [15]:
```
# Divide into segments
# We create 4 quartile ranges
df_RFM['M'] = pd.qcut(df_RFM['MonetaryValue'], q = 4, labels = range(1,5))
df_RFM['R'] = pd.qcut(df_RFM['Recency'], q = 4, labels = list(range(4,0,-1)))
df_RFM['F'] = pd.qcut(df_RFM['Frequency'], q = 4, labels = range(1,5))

df_RFM.head()
```

Out[15]:

| CUSTOMERNAME | Recency | Frequency | MonetaryValue | M | R | F |
|---|---|---|---|---|---|---|
| AV Stores, Co. | 196 | 51 | 157807.81 | 4 | 2 | 4 |
| Alpha Cognac | 65 | 20 | 70488.44 | 2 | 4 | 2 |
| Amica Models & Co. | 265 | 26 | 94117.26 | 3 | 1 | 2 |
| Anna's Decorations, Ltd | 84 | 46 | 153996.13 | 4 | 3 | 4 |
| Atelier graphique | 188 | 7 | 24179.96 | 1 | 2 | 1 |

In [16]:
```
#Create another column for RFM score
df_RFM['RFM_Score'] = df_RFM[['R', 'M', 'F']].sum(axis=1)
df_RFM.head()
```

Out[16]:

| CUSTOMERNAME | Recency | Frequency | MonetaryValue | M | R | F | RFM_Score |
|---|---|---|---|---|---|---|---|
| AV Stores, Co. | 196 | 51 | 157807.81 | 4 | 2 | 4 | 10 |
| Alpha Cognac | 65 | 20 | 70488.44 | 2 | 4 | 2 | 8 |
| Amica Models & Co. | 265 | 26 | 94117.26 | 3 | 1 | 2 | 6 |
| Anna's Decorations, Ltd | 84 | 46 | 153996.13 | 4 | 3 | 4 | 11 |
| Atelier graphique | 188 | 7 | 24179.96 | 1 | 2 | 1 | 4 |

In [17]:
```python
def rfm_level(df):
    if bool(df['RFM_Score'] >= 10):
        return 'High Value Customer'

    elif bool(df['RFM_Score'] < 10) and bool(df['RFM_Score'] >= 6):
        return 'Mid Value Customer'
    else:
        return 'Low Value Customer'
df_RFM['RFM_Level'] = df_RFM.apply(rfm_level, axis = 1)
df_RFM.head()
```

Out[17]:

| CUSTOMERNAME | Recency | Frequency | MonetaryValue | M | R | F | RFM_Score | RFM_Leve |
|---|---|---|---|---|---|---|---|---|
| AV Stores, Co. | 196 | 51 | 157807.81 | 4 | 2 | 4 | 10 | High Valu Custome |
| Alpha Cognac | 65 | 20 | 70488.44 | 2 | 4 | 2 | 8 | Mid Valu Custome |
| Amica Models & Co. | 265 | 26 | 94117.26 | 3 | 1 | 2 | 6 | Mid Valu Custome |
| Anna's Decorations, Ltd | 84 | 46 | 153996.13 | 4 | 3 | 4 | 11 | High Valu Custome |
| Atelier graphique | 188 | 7 | 24179.96 | 1 | 2 | 1 | 4 | Low Valu Custome |

In [18]:
```python
# Time to perform KMeans
data = df_RFM[['Recency', 'Frequency', 'MonetaryValue']]
data.head()
```

Out[18]:

| CUSTOMERNAME | Recency | Frequency | MonetaryValue |
|---|---|---|---|
| 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 [19]:
```python
# Our data is skewed we must remove it by performing log transformation
data_log = np.log(data)
data_log.head()
```

Out[19]:

| CUSTOMERNAME | Recency | Frequency | MonetaryValue |
|---|---|---|---|
| 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 [20]:
```python
#Standardization
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
scaler.fit(data_log)
data_normalized = scaler.transform(data_log)
data_normalized = pd.DataFrame(data_normalized, index = data_log.index, columns=
data_normalized.describe().round(2)
```

Out[20]:

| | 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 [21]:
```python
#Fit KMeans and use elbow method to choose the number of clusters
import matplotlib.pyplot as plt
import seaborn as sns
```

```python
from sklearn.cluster import KMeans

sse = {}

for k in range(1, 21):
    kmeans = KMeans(n_clusters = k, random_state = 1)
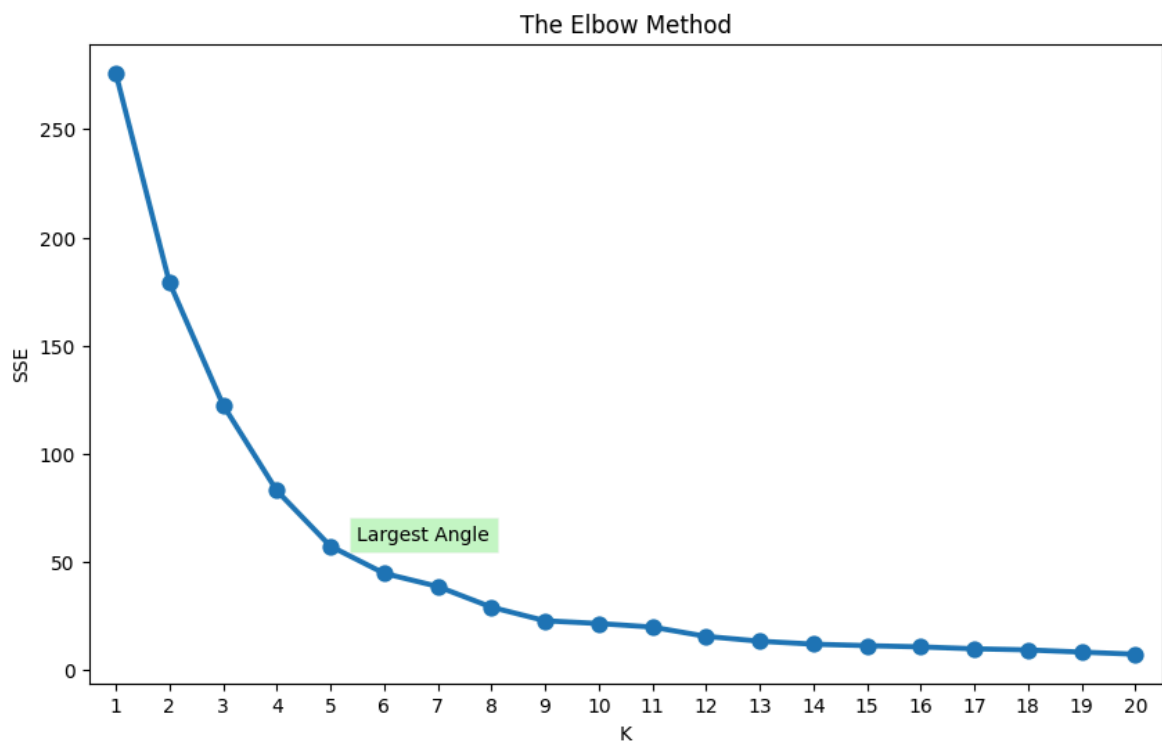    kmeans.fit(data_normalized)
    sse[k] = kmeans.inertia_
```

In [22]:
```python
plt.figure(figsize=(10,6))
plt.title('The Elbow Method')

plt.xlabel('K')
plt.ylabel('SSE')
plt.style.use('ggplot')

sns.pointplot(x=list(sse.keys()), y = list(sse.values()))
plt.text(4.5, 60, "Largest Angle", bbox = dict(facecolor = 'lightgreen', alpha =
plt.show()
```



In [23]:
```python
# 5 number of clusters seems good
kmeans = KMeans(n_clusters=5, random_state=1)
kmeans.fit(data_normalized)
cluster_labels = kmeans.labels_

data_rfm = data.assign(Cluster = cluster_labels)
data_rfm.head()
```

Out[23]:

| CUSTOMERNAME | Recency | Frequency | MonetaryValue | Cluster |
|---|---|---|---|---|
| AV Stores, Co. | 196 | 51 | 157807.81 | 4 |
| Alpha Cognac | 65 | 20 | 70488.44 | 2 |
| Amica Models & Co. | 265 | 26 | 94117.26 | 2 |
| Anna's Decorations, Ltd | 84 | 46 | 153996.13 | 4 |
| Atelier graphique | 188 | 7 | 24179.96 | 1 |