Name:Gayatri Ghorpade

Roll no:14144

```python
In [11]: import pandas as pd
         import numpy as np
         import seaborn as sns
         import matplotlib.pyplot as plt
```

```python
In [13]: from sklearn.cluster import KMeans
         from sklearn.decomposition import PCA
```
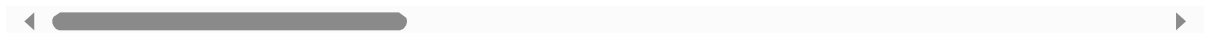
```python
In [20]: df = pd.read_csv("sales_data_sample.csv", encoding ="Latin-1")
```

```python
In [22]: df.head()
```

Out[22]:

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

5 rows × 25 columns

```python
In [24]: df.shape
```

Out[24]: (2823, 25)

```python
In [26]: df.describe()
```

|  | ORDERNUMBER | QUANTITYORDERED | PRICEEACH | ORDERLINENUMBER | SAL |
|---|---|---|---|---|---|
| count | 2823.000000 | 2823.000000 | 2823.000000 | 2823.000000 | 2823.0000 |
| mean | 10258.725115 | 35.092809 | 83.658544 | 6.466171 | 3553.8890 |
| std | 92.085478 | 9.741443 | 20.174277 | 4.225841 | 1841.8651 |
| min | 10100.000000 | 6.000000 | 26.880000 | 1.000000 | 482.1300 |
| 25% | 10180.000000 | 27.000000 | 68.860000 | 3.000000 | 2203.4300 |
| 50% | 10262.000000 | 35.000000 | 95.700000 | 6.000000 | 3184.8000 |
| 75% | 10333.500000 | 43.000000 | 100.000000 | 9.000000 | 4508.0000 |
| max | 10425.000000 | 97.000000 | 100.000000 | 18.000000 | 14082.8000 |

In [28]:
```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 [30]:
```python
df.isnull().sum()
```

```
Out[30]: 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
         PHONE                   0
         ADDRESSLINE1            0
         ADDRESSLINE2         2521
         CITY                    0
         STATE                1486
         POSTALCODE             76
         COUNTRY                 0
         TERRITORY            1074
         CONTACTLASTNAME         0
         CONTACTFIRSTNAME        0
         DEALSIZE                0
         dtype: int64
```

In [32]: `df.dtypes`

```
Out[32]: 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
         PHONE                  object
         ADDRESSLINE1           object
         ADDRESSLINE2           object
         CITY                   object
         STATE                  object
         POSTALCODE             object
         COUNTRY                object
         TERRITORY              object
         CONTACTLASTNAME        object
         CONTACTFIRSTNAME       object
         DEALSIZE               object
         dtype: object
```

```python
In [36]: df_drop = ['ADDRESSLINE1', 'ADDRESSLINE2', 'STATUS', 'POSTALCODE', 'CITY']
```

```python
In [38]: df = df.drop(df_drop, axis=1)
```

```python
In [40]: df.isnull().sum()
```

```
Out[40]: ORDERNUMBER         0
         QUANTITYORDERED     0
         PRICEEACH           0
         ORDERLINENUMBER     0
         SALES               0
         ORDERDATE           0
         QTR_ID              0
         MONTH_ID            0
         YEAR_ID             0
         PRODUCTLINE         0
         MSRP                0
         PRODUCTCODE         0
         CUSTOMERNAME        0
         PHONE               0
         STATE            1486
         COUNTRY             0
         TERRITORY        1074
         CONTACTLASTNAME     0
         CONTACTFIRSTNAME    0
         DEALSIZE            0
         dtype: int64
```

```python
In [42]: df.dtypes
```

```
Out[42]: ORDERNUMBER         int64
         QUANTITYORDERED     int64
         PRICEEACH         float64
         ORDERLINENUMBER     int64
         SALES             float64
         ORDERDATE          object
         QTR_ID              int64
         MONTH_ID            int64
         YEAR_ID             int64
         PRODUCTLINE        object
         MSRP                int64
         PRODUCTCODE        object
         CUSTOMERNAME       object
         PHONE              object
         STATE              object
         COUNTRY            object
         TERRITORY          object
         CONTACTLASTNAME    object
         CONTACTFIRSTNAME   object
         DEALSIZE           object
         dtype: object
```

```python
In [44]: df['COUNTRY'].unique()
```

```
Out[44]:  array(['USA', 'France', 'Norway', 'Australia', 'Finland', 'Austria', 'UK',
                 'Spain', 'Sweden', 'Singapore', 'Canada', 'Japan', 'Italy',
                 'Denmark', 'Belgium', 'Philippines', 'Germany', 'Switzerland',
                 'Ireland'], dtype=object)
```

In [46]:
```python
df['PRODUCTLINE'].unique()
```

```
Out[46]:  array(['Motorcycles', 'Classic Cars', 'Trucks and Buses', 'Vintage Cars',
                 'Planes', 'Ships', 'Trains'], dtype=object)
```

In [48]:
```python
df['DEALSIZE'].unique()
```

```
Out[48]:  array(['Small', 'Medium', 'Large'], dtype=object)
```

In [50]:
```python
productline = pd.get_dummies(df['PRODUCTLINE'])
Dealsize = pd.get_dummies(df['DEALSIZE'])
```

In [52]:
```python
df = pd.concat([df, productline,Dealsize],axis=1)
```

In [54]:
```python
df_drop = ['COUNTRY', 'PRODUCTLINE', 'DEALSIZE']
df = df.drop(df_drop, axis =1 )
```

In [56]:
```python
df['PRODUCTCODE'] = pd.Categorical(df[ 'PRODUCTCODE']).codes
```
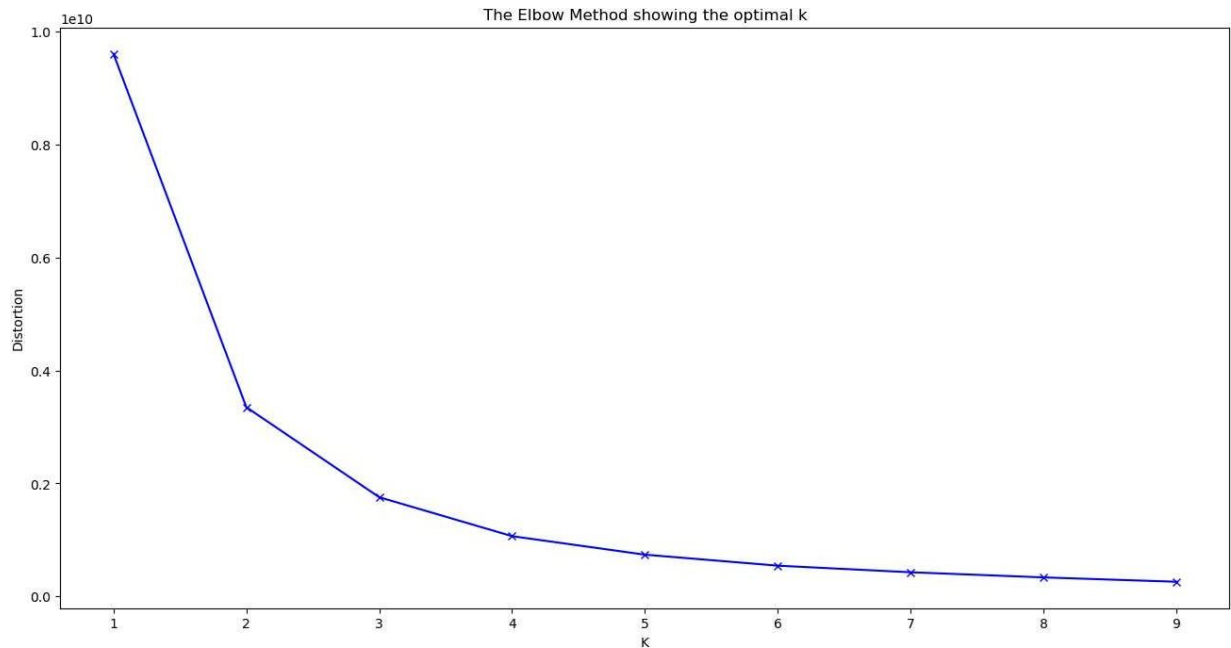
In [58]:
```python
df.drop('ORDERDATE', axis = 1, inplace=True)
```

In [60]:
```python
df.dtypes
```

```
Out[60]:  ORDERNUMBER              int64
          QUANTITYORDERED          int64
          PRICEEACH              float64
          ORDERLINENUMBER          int64
          SALES                  float64
          QTR_ID                   int64
          MONTH_ID                 int64
          YEAR_ID                  int64
          MSRP                     int64
          PRODUCTCODE               int8
          CUSTOMERNAME            object
          PHONE                  object
          STATE                  object
          TERRITORY              object
          CONTACTLASTNAME        object
          CONTACTFIRSTNAME       object
          Classic Cars             bool
          Motorcycles              bool
          Planes                   bool
          Ships                    bool
          Trains                   bool
          Trucks and Buses         bool
          Vintage Cars             bool
          Large                    bool
          Medium                   bool
          Small                    bool
          dtype: object
```

```python
In [72]:  distortions = []
          K = range(1, 10)
          for k in K:
              kmeanModel = KMeans(n_clusters=k, random_state=42)
              kmeanModel.fit(df_numeric)
              distortions.append(kmeanModel.inertia_)
```

```python
In [74]:  plt.figure(figsize=(16,8))
          plt.plot(K, distortions, 'bx-')
          plt.xlabel('K')
          plt.ylabel('Distortion')
          plt.title('The Elbow Method showing the optimal k')
          plt.show()
```

The Elbow Method showing the optimal k

```
In [76]: x_train = df.values
```

```
In [78]: x_train.shape
```

```
Out[78]: (2823, 26)
```

```
In [112... model = KMeans(n_clusters=3, random_state=2)
         model.fit(x_train)
         predictions = model.predict(x_train)
```

```
In [114... unique, counts = np.unique(predictions, return_counts=True)
```

```
In [116... counts = counts.reshape(1,3)
```

```
In [129... counts_df = pd.DataFrame(counts, columns=['Cluster', 'Cluster2', 'Cluster3'])
```

```
In [131... counts_df.head()
```

Out[131...

|   | Cluster | Cluster2 | Cluster3 |
|---|---------|----------|----------|
| **0** | 1344 | 398 | 1081 |

```
In [139... pca = PCA(n_components=2)
         reduced_X = pd.DataFrame(pca.fit_transform(x_train), columns=['PCA1', 'PEA2'])
         reduced_X.head()
```

| | PCA1 | PEA2 |
|---|---|---|
| 0 | -683.111504 | -150.512921 |
| 1 | -788.262223 | -136.226986 |
| 2 | 330.177513 | -125.205410 |
| 3 | 192.505559 | -113.881926 |
| 4 | 1651.047570 | -102.509854 |

In [ ]: