

```
In [1]: 1 import pandas as pd
        2 import numpy as np
        3 import seaborn as sns
        4 import matplotlib.pyplot as plt
        5
```

```
In [2]: 1 from sklearn.cluster import KMeans
        2 from sklearn.decomposition import PCA
```

```
In [8]: 1 df = pd.read_csv("sales_data_sample.csv", encoding = "Latin-1")
```

```
In [11]: 1 df.head()
```

Out[11]:

| | ORDERNUMBER | QUANTITYORDERED | PRICEEACH | ORDERLINENUMBER | SALES | ORDERDATE |
|---|-------------|-----------------|-----------|-----------------|---------|------------|
| 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 [12]: 1 df.shape
```

Out[12]: (2823, 7)

```
In [13]: 1 df.describe()
```

Out[13]:

| | ORDERNUMBER | QUANTITYORDERED | PRICEEACH | ORDERLINENUMBER | SALES |
|-------|--------------|-----------------|-------------|-----------------|--------------|
| count | 2823.000000 | 2823.000000 | 2823.000000 | 2823.000000 | 2823.000000 |
| mean | 10258.725115 | 35.092809 | 83.658544 | 6.466171 | 3553.889072 |
| std | 92.085478 | 9.741443 | 20.174277 | 4.225841 | 1841.865106 |
| min | 10100.000000 | 6.000000 | 26.880000 | 1.000000 | 482.130000 |
| 25% | 10180.000000 | 27.000000 | 68.860000 | 3.000000 | 2203.430000 |
| 50% | 10262.000000 | 35.000000 | 95.700000 | 6.000000 | 3184.800000 |
| 75% | 10333.500000 | 43.000000 | 100.000000 | 9.000000 | 4508.000000 |
| max | 10425.000000 | 97.000000 | 100.000000 | 18.000000 | 14082.800000 |



```
In [15]: 1 df.info()
```

In [15]: 1 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 [16]: 1 df.isnull().sum()

In [16]: 1 df.isnull().sum()

```
Out[16]: 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 [17]: 1 df.dtypes

```
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
```

In [18]: 1 df_drop = ['ADDRESSLINE1', 'ADDRESSLINE2', 'STATUS', 'POSTALCODE', 'CITY']

In [19]: 1 df = df.drop(df_drop, axis=1)

In [21]: 1 df.isnull().sum()

```
In [21]: 1 df.isnull().sum()
```

```
Out[21]: 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
          COUNTRY          0
          DEALSIZE         0
          dtype: int64
```

```
In [22]: 1 df.dtypes
```

```
Out[22]: 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
          COUNTRY          object
          DEALSIZE         object
          dtype: object
```

```
In [26]: 1 df['COUNTRY'].unique()
         2
```

```
Out[26]: array(['USA', 'France', 'Norway', 'Australia', 'Finland', 'Austria', 'UK',
                'Spain', 'Sweden', 'Singapore', 'Canada', 'Japan', 'Italy',
                'Denmark', 'Belgium', 'Philippines', 'Germany', 'Switzerland',
                'Ireland'], dtype=object)
```

```
In [27]: 1 df['PRODUCTLINE'].unique()
```

```
Out[27]: array(['Motorcycles', 'Classic Cars', 'Trucks and Buses', 'Vintage Cars',
                'Planes', 'Ships', 'Trains'], dtype=object)
```

```
In [28]: 1 df['DEALSIZE'].unique()
```

```
Out[28]: array(['Small', 'Medium', 'Large'], dtype=object)
```

```
In [30]: 1 productline = pd.get_dummies(df['PRODUCTLINE'])
         2 Dealsize = pd.get_dummies(df['DEALSIZE'])
```

```
In [35]: 1 df = pd.concat([df, productline, Dealsize], axis=1)
```

```
In [35]: 1 df = pd.concat([df, productline,Dealsize],axis=1)
```

```
In [36]: 1 df_drop = ['COUNTRY', 'PRODUCTLINE', 'DEALSIZE']  
2 df = df.drop(df_drop, axis =1 )
```

```
In [37]: 1 df['PRODUCTCODE'] = pd.Categorical(df['PRODUCTCODE']).codes
```

```
In [40]: 1 df.drop('ORDERDATE', axis = 1, inplace=True)
```

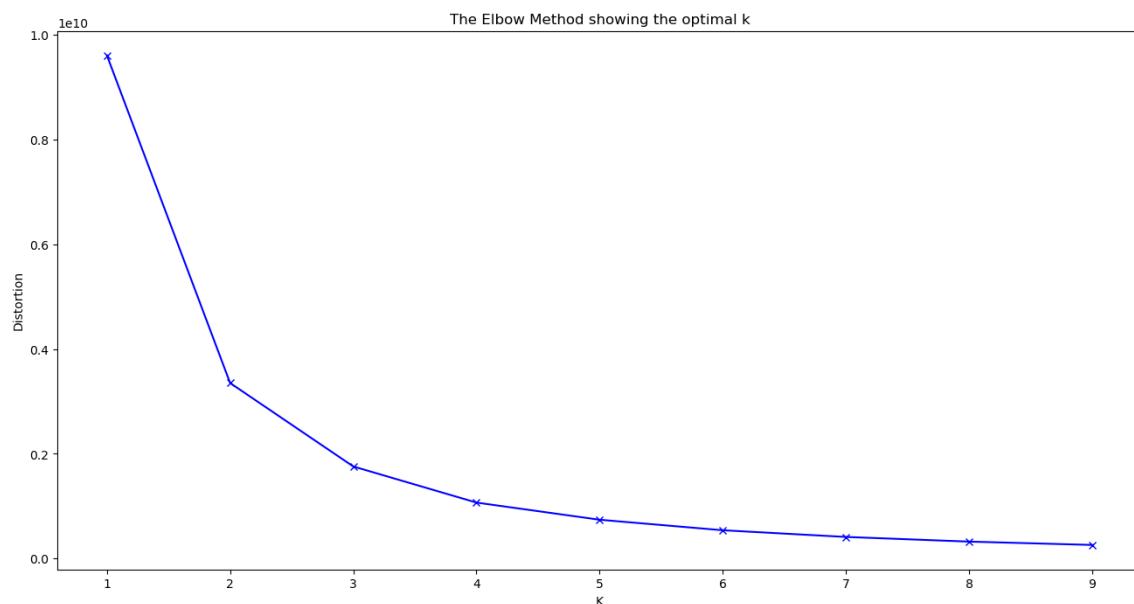
```
In [41]: 1 df.dtypes
```

```
Out[41]: ORDERNUMBER          int64  
QUANTITYORDERED          int64  
PRICEEACH                float64  
ORDERLINENUMBER          int64  
SALES                    float64  
QTR_ID                   int64  
MONTH_ID                 int64  
YEAR_ID                  int64  
MSRP                     int64  
PRODUCTCODE              int8  
Classic Cars              uint8  
Motorcycles               uint8  
Planes                    uint8  
Ships                     uint8  
Trains                    uint8  
Trucks and Buses          uint8  
Vintage Cars              uint8  
Large                     uint8  
Medium                    uint8  
Small                     uint8  
dtype: object
```

```
In [43]: 1 distortions = []  
2 K = range(1,10)  
3 for k in K:  
4     kmeanModel = KMeans(n_clusters=k)  
5     kmeanModel.fit(df)  
6     distortions.append(kmeanModel.inertia_)
```

```
In [44]: 1 plt.figure(figsize=(16,8))
```

```
In [44]: 1 plt.figure(figsize=(16,8))
2 plt.plot(K, distortions, 'bx-')
3 plt.xlabel('K')
4 plt.ylabel('Distortion')
5 plt.title('The Elbow Method showing the optimal k')
6 plt.show()
```



```
In [45]: 1 x_train = df.values
```

```
In [46]: 1 x_train.shape
```

```
Out[46]: (2823, 20)
```

```
In [49]: 1 model = KMeans (n_clusters=3, random_state=2)
2 model = model.fit(x_train)
3 predictions = model.predict(x_train)
```

```
In [51]: 1 unique, counts = np.unique(predictions, return_counts=True)
```

```
In [52]: 1 counts = counts.reshape(1,3)
```

```
In [55]: 1 counts_df = pd.DataFrame(counts, columns=['Cluster', 'Cluster2', 'Cluster3'])
```

```
In [57]: 1 counts_df.head()
```

```
Out[57]:
```

| | Cluster | Cluster2 | Cluster3 |
|---|---------|----------|----------|
| 0 | 1083 | 1367 | 373 |

```
In [59]: 1 pca = PCA(n_components=2)
```

```
In [59]: 1 pca = PCA(n_components=2)
          2 reduced_X = pd.DataFrame(pca.fit_transform(x_train), columns=['PCA1', 'PCA2'])
          3 reduced_X.head()
```

Out[59]:

| | PCA1 | PCA2 |
|---|-------------|-------------|
| 0 | -682.790370 | -151.271539 |
| 1 | -787.939342 | -136.994834 |
| 2 | 330.482091 | -125.876905 |
| 3 | 192.812426 | -114.565402 |
| 4 | 1651.330150 | -103.067424 |

In []:

1