

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler
from scipy.cluster.hierarchy import dendrogram, linkage
from sklearn.metrics import silhouette_score
```

```
df = pd.read_csv('sales_data_sample.csv', encoding='ISO-8859-1')
```

```
print(df.head())
print(df.info())
```

```
3      10145      45      83.26      6  3/46.70
4      10159      49     100.00     14  5205.27
```

```
      ORDERDATE  STATUS  QTR_ID  MONTH_ID  YEAR_ID  ...  \
0  2/24/2003 0:00  Shipped      1          2    2003  ...
1  5/7/2003 0:00  Shipped      2          5    2003  ...
2  7/1/2003 0:00  Shipped      3          7    2003  ...
3  8/25/2003 0:00  Shipped      3          8    2003  ...
4 10/10/2003 0:00  Shipped      4         10    2003  ...
```

```
      ADDRESSLINE1  ADDRESSLINE2      CITY STATE  \
0      897 Long Airport Avenue      NaN      NYC  NY
1          59 rue de l'Abbaye      NaN      Reims  NaN
2  27 rue du Colonel Pierre Avia      NaN      Paris  NaN
3      78934 Hillside Dr.      NaN  Pasadena  CA
4      7734 Strong St.      NaN  San Francisco  CA
```

```
      POSTALCODE  COUNTRY  TERRITORY  CONTACTLASTNAME  CONTACTFIRSTNAME  DEALSIZE
0      10022      USA      NaN      Yu      Kwai      Small
1      51100  France      EMEA      Henriot      Paul      Small
2      75508  France      EMEA      Da Cunha      Daniel      Medium
3      90003      USA      NaN      Young      Julie      Medium
4      NaN      USA      NaN      Brown      Julie      Medium
```

```
[5 rows x 25 columns]
<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

```

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
None

```

```

features = df[['QUANTITYORDERED', 'PRICEEACH', 'SALES', 'QTR_ID', 'MONTH_ID', 'YEAR_ID', 'MSRP']]
features = features.dropna()
scaler = StandardScaler()
scaled_features = scaler.fit_transform(features)

```

```

inertia = []
K = range(1, 11)

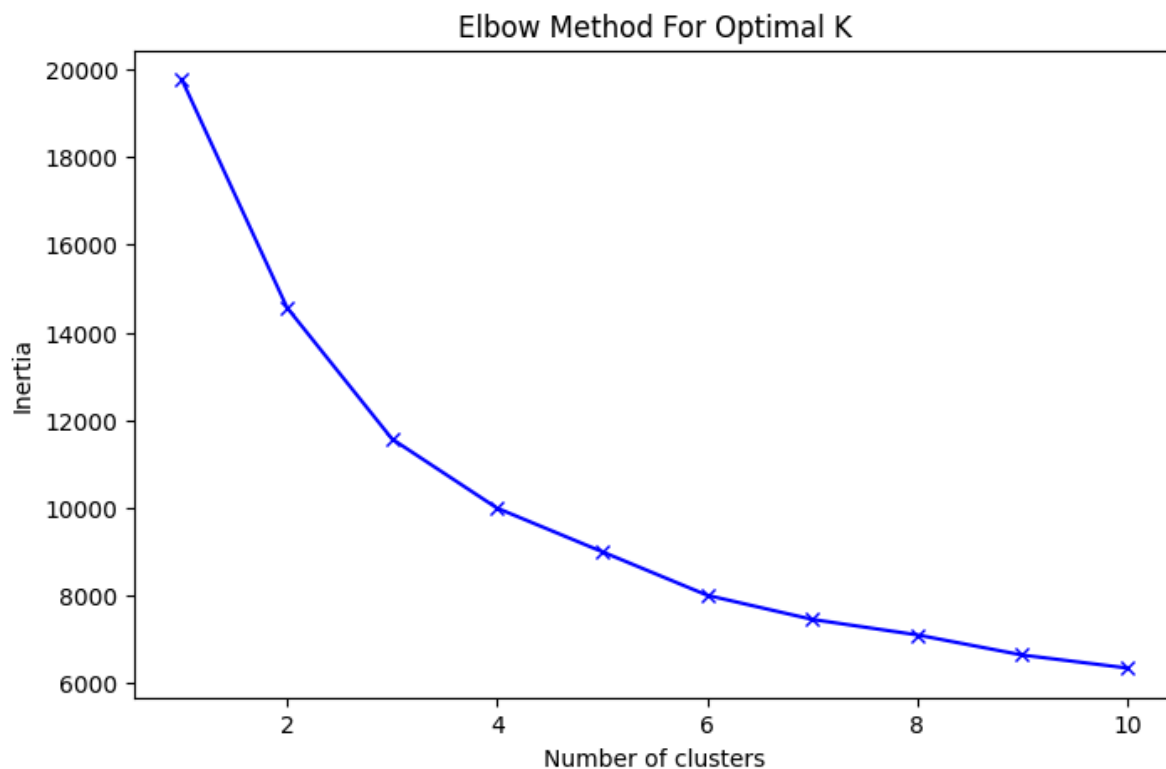
for k in K:
    kmeans = KMeans(n_clusters=k, random_state=42)
    kmeans.fit(scaled_features)
    inertia.append(kmeans.inertia_)

```

```

plt.figure(figsize=(8, 5))
plt.plot(K, inertia, 'bx-')
plt.xlabel('Number of clusters')
plt.ylabel('Inertia')
plt.title('Elbow Method For Optimal K')
plt.show()

```



```

optimal_k = 3
kmeans = KMeans(n_clusters=optimal_k, random_state=42)
kmeans.fit(scaled_features)

df['Cluster'] = kmeans.labels_

```

```
print(df['Cluster'].value_counts())
```

```
Cluster
1    1051
0     923
2     849
Name: count, dtype: int64
```