

Implement K-Means clustering/ hierarchical clustering on

- sales_data_sample.csv dataset. Determine the number of clusters using the elbow method.

```
1 import pandas as pd
2 import numpy as np
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 #Importing the required libraries.

1 from sklearn.cluster import KMeans, k_means #For clustering
2 from sklearn.decomposition import PCA #Linear Dimensionality reduction.

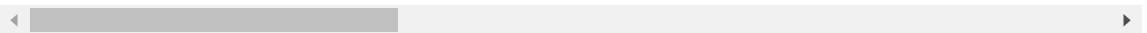
1 df = pd.read_csv('sales_data_sample.csv', encoding = "Latin-1")
2 #df = pd.read_csv("sales_data_sample.csv") #Loading the dataset.
```

Preprocessing

```
1 df.head()
```

	ORDERNUMBER	QUANTITYORDERED	PRICEEACH	ORDERLINENUMBER	SALES	ORDERDATE	
0	10107	30	95.70	2	2871.00	2/24/2003 0:00	S
1	10121	34	81.35	5	2765.90	5/7/2003 0:00	S
2	10134	41	94.74	2	3884.34	7/1/2003 0:00	S
3	10145	45	83.26	6	3746.70	8/25/2003 0:00	S
4	10159	49	100.00	14	5205.27	10/10/2003 0:00	S

5 rows × 25 columns



```
1 df.shape
```

```
(2823, 25)
```

```
1 df.describe()
```

	ORDERNUMBER	QUANTITYORDERED	PRICEEACH	ORDERLINENUMBER	SALES	
count	2823.000000	2823.000000	2823.000000	2823.000000	2823.000000	2
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	

```
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
```

```
1 df.isnull().sum()
```

```
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

1 df.dtypes

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

```
1 df_drop = ['ADDRESSLINE1', 'ADDRESSLINE2', 'STATUS','POSTALCODE', 'CITY', 'TERRITORY', 'PHONE', 'STATE', 'CONTACTLASTNAME', 'CONTACTFIRSTNAME', 'DEALSIZE']
2 df = df.drop(df_drop, axis=1) #Dropping the categorical unnecessary columns along with columns having null values
```

1 df.isnull().sum()

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

1 df.dtypes

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

1 # Checking the categorical columns.

1 df['COUNTRY'].unique()

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

1 df['PRODUCTLINE'].unique()

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

1 df['DEALSIZE'].unique()

array(['Small', 'Medium', 'Large'], dtype=object)

1 productline = pd.get_dummies(df['PRODUCTLINE']) #Converting the categorical columns.
2 Dealsize = pd.get_dummies(df['DEALSIZE'])

1 df = pd.concat([df,productline,Dealsize], axis = 1)

1 df_drop = ['COUNTRY','PRODUCTLINE','DEALSIZE'] #Dropping Country too as there are alot of countries.
2 df = df.drop(df_drop, axis=1)

1 df['PRODUCTCODE'] = pd.Categorical(df['PRODUCTCODE']).codes #Converting the datatype.

1 df.drop('ORDERDATE', axis=1, inplace=True) #Dropping the Orderdate as Month is already included.

1 df.dtypes #All the datatypes are converted into numeric

QUANTITYORDERED    int64
PRICEEACH          float64
ORDERLINENUMBER    int64

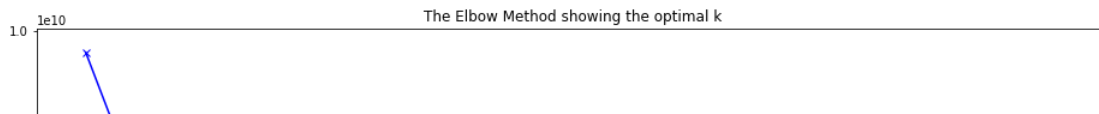
```

```
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
```

▼ Plotting the Elbow Plot to determine the number of clusters.

```
1 distortions = [] # Within Cluster Sum of Squares from the centroid
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_) #Appending the inertia to the Distortions

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()
```



▼ As the number of k increases Inertia decreases.

Observations: A Elbow can be observed at 3 and after that the curve decreases gradually.

```
1 X_train = df.values #Returns a numpy array.
```

```
1 X_train.shape
```

```
(2823, 19)
```

```
1 model = KMeans(n_clusters=3,random_state=2) #Number of cluster = 3
```

```
2 model = model.fit(X_train) #Fitting the values to create a model.
```

```
3 predictions = model.predict(X_train) #Predicting the cluster values (0,1,or 2)
```

```
1 unique,counts = np.unique(predictions,return_counts=True)
```

```
1 counts = counts.reshape(1,3)
```

```
1 counts_df = pd.DataFrame(counts,columns=['Cluster1','Cluster2','Cluster3'])
```

```
1 counts_df.head()
```

	Cluster1	Cluster2	Cluster3
0	1083	1367	373

▼ Visualization

```
1 pca = PCA(n_components=2) #Converting all the features into 2 columns to make it easy to visualize using Princip
```

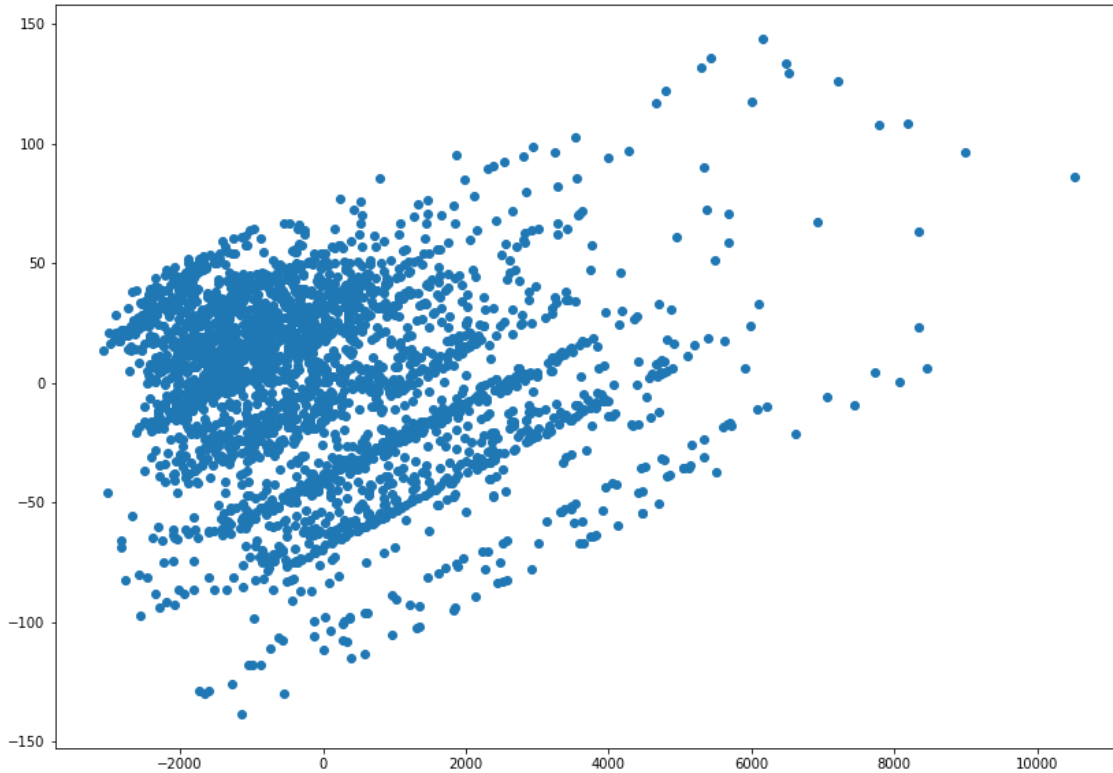
```
1 reduced_X = pd.DataFrame(pca.fit_transform(X_train),columns=['PCA1','PCA2']) #Creating a DataFrame.
```

```
1 reduced_X.head()
```

PCA1 PCA2

```
1 #Plotting the normal Scatter Plot
2 plt.figure(figsize=(14,10))
3 plt.scatter(reduced_X['PCA1'],reduced_X['PCA2'])
```

<matplotlib.collections.PathCollection at 0x7f0fe6fad810>



```
1 model.cluster_centers_ #Finding the centriods. (3 Centriods in total. Each Array contains a centriods for partic
```

```
array([[ 3.72031394e+01,  9.52120960e+01,  6.44967682e+00,
         4.13868425e+03,  2.72022161e+00,  7.09879963e+00,
         2.00379409e+03,  1.13248384e+02,  5.04469067e+01,
         3.74884580e-01,  1.15420129e-01,  9.41828255e-02,
         8.21791320e-02,  1.84672207e-02,  1.16343490e-01,
         1.98522622e-01,  2.08166817e-17,  1.00000000e+00,
        -3.38618023e-15],
       [ 3.08302853e+01,  7.00755230e+01,  6.67300658e+00,
         2.12409474e+03,  2.71762985e+00,  7.09509876e+00,
         2.00381127e+03,  7.84784199e+01,  6.24871982e+01,
         2.64813460e-01,  1.21433797e-01,  1.29480614e-01,
         1.00219459e-01,  3.87710315e-02,  9.21726408e-02,
         2.53108998e-01,  6.93889390e-18,  6.21799561e-02,
         9.37820044e-01],
       [ 4.45871314e+01,  9.98931099e+01,  5.75603217e+00,
         7.09596863e+03,  2.71045576e+00,  7.06434316e+00,
         2.00389008e+03,  1.45823056e+02,  3.14959786e+01,
         5.33512064e-01,  1.07238606e-01,  7.23860590e-02,
         2.14477212e-02,  1.07238606e-02,  1.31367292e-01,
```

```
1.23324397e-01, 4.20911528e-01, 5.79088472e-01,  
-1.99840144e-15]])
```

```
1 reduced_centers = pca.transform(model.cluster_centers_) #Transforming the centroids into 3 in x and y coordinate
```

[+ Code](#)[+ Text](#)

```
1 reduced_centers
```

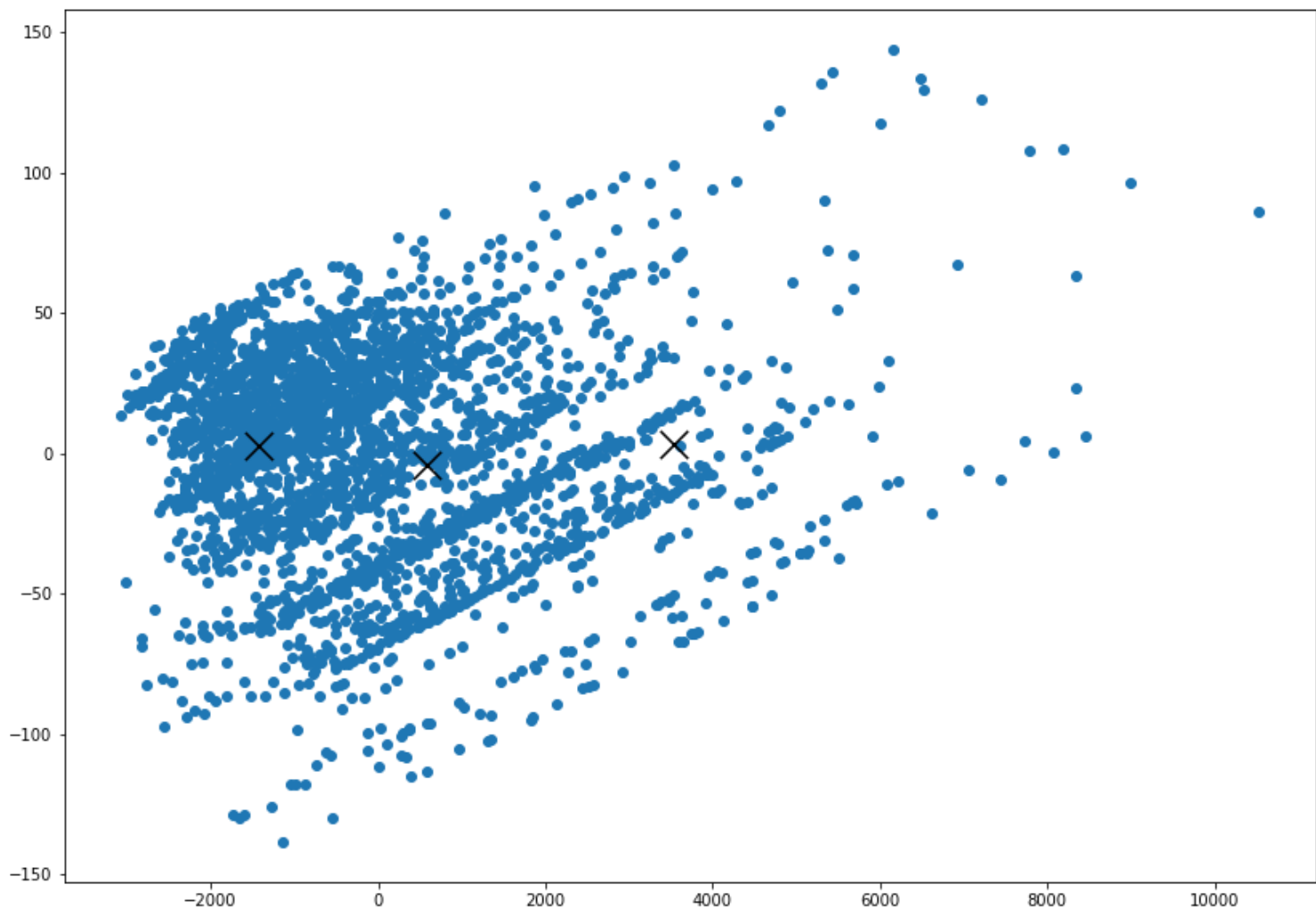
```
array([[ 5.84994044e+02, -4.36786931e+00],  
       [-1.43005891e+03,  2.60041009e+00],  
       [ 3.54247180e+03,  3.15185487e+00]])
```

```
1 plt.figure(figsize=(14,10))
```

```
2 plt.scatter(reduced_X['PCA1'],reduced_X['PCA2'])
```

```
3 plt.scatter(reduced_centers[:,0],reduced_centers[:,1],color='black',marker='x',s=300) #Plotting the centroids
```

<matplotlib.collections.PathCollection at 0x7f0fe6f2e410>



```
1 reduced_X['Clusters'] = predictions #Adding the Clusters to the reduced dataframe.
```

```
1 reduced_X.head()
```

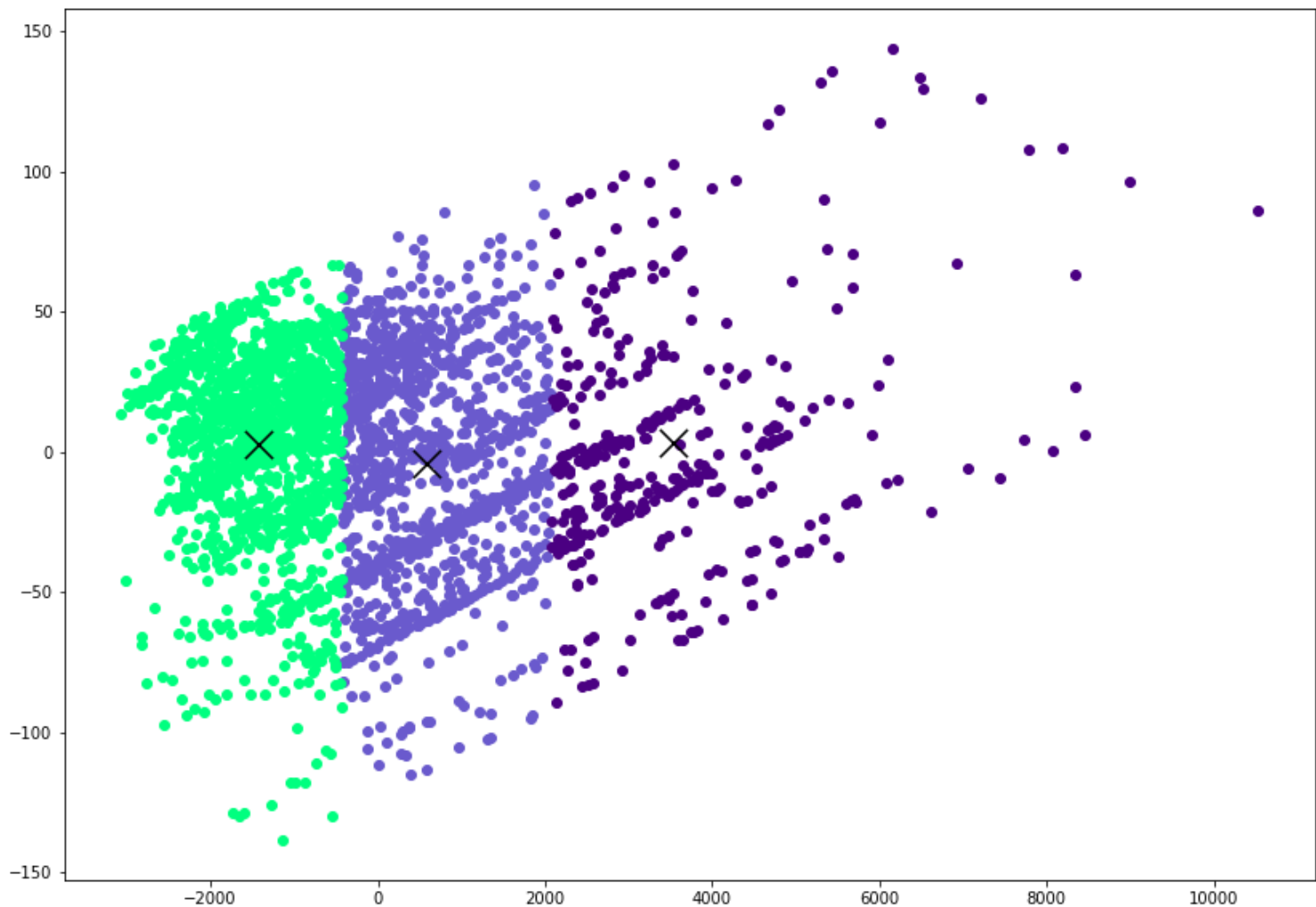

	PCA1	PCA2	Clusters
0	-682.488323	-42.819535	1
1	-787.665502	-41.694991	1
2	330.732170	-26.481208	0

```

1 #Plotting the clusters
2 plt.figure(figsize=(14,10))
3 #                                taking the cluster number and first column                                taking the same cluster number and se
4 plt.scatter(reduced_X[reduced_X['Clusters'] == 0].loc[:, 'PCA1'],reduced_X[reduced_X['Clusters'] == 0].loc[:, 'PCA2'],color='red',s=100)
5 plt.scatter(reduced_X[reduced_X['Clusters'] == 1].loc[:, 'PCA1'],reduced_X[reduced_X['Clusters'] == 1].loc[:, 'PCA2'],color='blue',s=100)
6 plt.scatter(reduced_X[reduced_X['Clusters'] == 2].loc[:, 'PCA1'],reduced_X[reduced_X['Clusters'] == 2].loc[:, 'PCA2'],color='green',s=100)
7
8
9 plt.scatter(reduced_centers[:,0],reduced_centers[:,1],color='black',marker='x',s=300)

```

<matplotlib.collections.PathCollection at 0x7f0fe700c1d0>



[Colab paid products](#) - [Cancel contracts here](#)

