[1]:	PCA Assignment  Data Set: Wine  1. Import Necessary libraries  import pandas as pd import numpy as np import matplotlib.pyplot as plt import seaborn as sns  2. Import Data
[2]:	Type   Alcohol   Male   Ash   Alcalinity   Magnesium   Phenols   Flavanoids   Proanthocyanins   Color   Hue   Dilution   Proline     1   14.23   1.71   2.43   15.6   127   2.80   3.06   0.28   2.29   5.84   1.04   3.92   1065     1   1   13.20   1.78   2.14   11.2   100   2.65   2.76   0.26   2.28   3.86   0.28   2.29   5.84   1.04   3.92   1065     2   1   14.37   1.95   2.65   18.6   101   2.80   3.24   0.30   0.281   5.68   1.03   3.17   1185     3   1   14.37   1.95   2.87   2.10   11.8   2.80   3.24   0.30   0.24   2.18   7.80   0.86   3.45   1480     4   1   13.24   2.59   2.87   2.10   11.8   2.80   2.69   0.39   0.39   1.82   2.31   2.93   7.35     3   13.71   5.65   2.45   2.05   9.5   1.68   0.61   0.52   1.06   7.70   0.64   1.74   7.40     174   3   13.40   3.10   2.48   2.30   1.02   1.80   0.75   0.43   1.41   7.30   0.70   1.56   7.50     175   3   13.27   2.8   2.8   2.0   2.0   1.20   1.59   0.69   0.43   1.41   7.30   0.70   1.56   7.50     176   3   13.17   2.59   2.7   2.00   1.20   1.59   0.69   0.43   1.45   1.35   0.20   0.59   1.56   8.35     176   3   13.17   2.59   2.7   2.00   1.20   1.59   0.69   0.43   1.35   0.20   0.51   1.50   0.50   1.50   0.50     178   Rows × 14 columns:
[3]: [3]: 	3. Data Understanding  Type   Alcohol   Malic   Ash   Alcalinity   Magnesium   Phenols   Flavanoids   Nonflavanoids   Proanthocyanins   Color   Hue   Dilution   Proline     1   14.23   1.71   2.43   15.6   127   2.80   3.06   0.28   2.29   5.64   1.04   3.92   1065     1   1   13.20   1.78   2.14   11.2   100   2.65   2.76   0.26   1.28   4.38   1.05   3.40   1050     2   1   13.16   2.36   2.67   18.6   101   2.80   3.24   0.30   2.81   5.68   1.03   3.17   1185     3   1   14.37   1.95   2.50   16.8   113   3.85   3.49   0.24   2.18   7.80   0.86   3.45   1480     4   1   13.24   2.59   2.87   21.0   118   2.80   2.69   0.39   1.82   4.32   1.04   2.93   7.35      1   1   1   1   1   1   1   1   1
[5]:	Colars   Pandas.core.frame.DataFrame   >
[6]: ,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,	Type
[8]:	25% 1.00000 12.362500 1.602500 2.210000 17.20000 88.00000 1.742500 1.205000 0.270000 1.25000 3.22000 0.782500 1.93750 500.500000  50% 2.000000 13.050000 1.865000 2.360000 19.50000 98.00000 2.355000 2.355000 0.340000 1.555000 4.69000 0.965000 2.78000 673.500000  75% 3.00000 13.677500 3.082500 2.557500 21.50000 107.000000 2.800000 2.875000 0.437500 1.95000 6.20000 1.120000 3.170000 985.000000  max 3.00000 14.830000 5.800000 3.230000 30.00000 162.000000 3.880000 5.080000 0.660000 3.580000 13.000000 1.710000 4.00000 1680.000000  Type
[9]: 10]:	### Converting into numpy array
11]:	liquor_array = data.values  larray([[1.423e+81, 1.710e+00, 2.439e+00,, 1.040e+00, 3.920e+00, 1.050e+03], [1.320e+01, 1.780e+00, 2.140e+00,, 1.030e+00, 3.400e+00, 1.050e+03], [1.310e+01, 2.380e+00, 2.60e+00,, 1.030e+00, 3.170e+00, 1.185e+03],, [1.327e+01, 4.280e+00, 2.260e+00,, 5.900e+01, 1.560e+00, 8.350e+02], [1.337e+01, 2.590e+00, 2.370e+00,, 6.080e+01, 1.620e+00, 8.400e+02], [1.413e+01, 4.100e+00, 2.740e+00,, 6.100e+01, 1.600e+00, 5.600e+02]])  # Normalizing data liquor_norm = scale(liquor_array) liquor_norm = scale(1quor_array) liquor_norm = scale(1quor_array) [[1.51861254, -0.5822498, 0.23205254,, 0.36217728, 1.04709957, 1.0130093], [[2.2402803, -0.49941338, -0.82799632,, 0.40009066, 1.1134493, 0.90524152], [[0.19667903, 0.213125, 1.10933436,, 0.31830389, 0.70857817, 1.9134493, 0.90527817, 1.0933436,, 0.31830389, 0.70857817, 1.7447444, 0.38095541,, -1.61212515, -1.4054488, 0.28057537], [[0.20923168, 0.22769377, 0.01273290,, -1.56825176, -1.4060981, 0.22649777, 0.01273290,, -1.56825176, -1.4060981, 0.2269777, 0.01273290,, -1.52437837, -1.2429777, -0.105106911])
14]: 15]:	from sklearn.decomposition import PCA from sklearn import preprocessing  PCA = PCA()  pca_liquor = PCA.fit_transform(liquor_norm) pca_liquor  array([[ 3.31675081e+00, -1.44346263e+00, -1.65739045e-01,,
16]: · 16]: ·	5.12492025e-01, 6.98766451e-01, 7.20776948e-02], [-2.38701709e+00, -2.29734668e+00], -5.50696197e-01,, 2.99821968e-01, 3.39820654e-01, -2.18657605e-02], [-3.20875816e+00, -2.76891957e+00, 1.01391366e+00,, -2.29964331e-01, -1.88787963e-01, -3.23964720e-01]])  STEP 2 = Explaining the Variance Using Principal Component:  # The amount of variance that each PCA explains is liquor_var = PCA.explained_variance_ratio_ liquor_var  array([0.3619848, 0.1920749 , 0.11123631, 0.0706903 , 0.06563294,
17]: <sup>1</sup>	liquor_cum_var  array([ 36.2 , 55.41, 66.53, 73.6 , 80.16, 85.1 , 89.34, 92.02, 94.24, 96.17, 97.91, 99.21, 100.01])  # Variance plot for PCA components obtained plt.plot(liquor_cum_var, color = "red") plt.show()  100
19]:	STEP 3 = Combine the Target and the Principal Components:  #create data frame from pca1 and pca2  df_liquor = pd.concat([pd.DataFrame(pca_liquor[:, 0:3], columns = ['PC1', 'PC2', 'PC3']), liquor[['Type']]], axis = 1)
1	2 2.516740 -1.031151 0.982819 1 3 3.757066 -2.756372 -0.176192 1 4 1.008908 -0.869831 2.026688 1 173 -3.370524 -2.216289 -0.342570 3 174 -2.601956 -1.757229 0.207581 3 175 -2.677839 -2.760899 -0.940942 3 176 -2.387017 -2.297347 -0.550696 3 177 -3.208758 -2.768920 1.013914 3
20]:	<pre>df_liquor['Type'].value_counts() 2     71 3     59 3     48 Name: Type, dtype: int64  STEP 4 = Plot the Principal Components on 2D : plt.suptitle('PC1 vs PC2 - Type', size = 20) sns.scatterplot(data = df_liquor, x = 'PC1', y = 'PC2', hue = 'Type') plt.show()  PC1 vs PC2 - Type  4</pre>
22]:	5. Clustering  Clustering a Data Frame with PC1, PC2 and PC3  (a) Hierarchical Clustering  STEP 1 = Data Pre-processing:  import pandas as pd import numpy as np impo
24]:	1 1 1 2 20 1.78 2.14 112 100 2.65 2.76 0.26 1.28 4.38 1.05 3.40 1050 2 1 1 1 1 1 2 2 3 2.67 18.0 101 2.80 3.24 0.30 2.81 5.68 1.03 3.17 1185 3 1 1 4.37 1.95 2.50 16.8 113 3.85 3.49 0.24 2.18 7.80 0.85 3.45 1.480 11quor_x_1 = pd.DataFrame(pca_liquor[:, 9:3], columns = ['PC1', 'PC2', 'PC3']) 11quor_x_1 = pd.DataFrame(pca_liquor[:, 9:3], columns = ['PC1', 'PC2', 'PC3']) 1 2 209465 0.333393 -2.02645 1 2 2.51740 1.031151 0.832819 3 3.757066 2.758379 0.898831 2.026688 4 10 2 2.95889 0.898831 2.02689 1 2 2.518740 2.031851 0.48289 0.489831 0.058989 1 2 2.518740 1.031151 0.832819 3 3.757066 2.758379 0.898831 0.056688 1 2 2.518740 1.031151 0.832819 3 3.757066 2.758379 0.898831 0.056688 1 2 2.518740 1.031151 0.832819 3 3.757066 2.758379 0.898831 0.056688 1 2 2.518740 0.031151 0.832819 1 2 2.518740 0.031151
25]:	STEP 2 = Finding the optimal number of clusters using the Dendrogram:  from scipy.cluster.hierarchy import linkage import scipy.cluster.hierarchy as sch  # create dendrogram  x = linkage(liquor_x_1, method = 'complete',metric = 'euclidean')  plt.figure(figsize = (13, 8)) plt.title('Hierarchical Dendrogram Plot') plt.xlabel('Index') plt.ylabel('Euclidean Distances')  sch.dendrogram(x) plt.show()
	Hierarchical Dendrogram Plot    STEP 3 = Training the hierarchical clustering model:
27]: 28]: 28]: 4 28]: 4	STEP 3 = Training the hierarchical clustering model:  from sklearn.cluster import AgglomerativeClustering  hc = AgglomerativeClustering(n_clusters = 3, affinity = 'euclidean', linkage = 'ward')  hc  AgglomerativeClustering(n_clusters=3)  x_hc = hc.fit_predict(liquor_x_1) x_hc  array([2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2
30]:	1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
31]:	
32]:	liquor_x_1['Clusters_1'].value_counts()  0 66 2 65 1 47 Name: Clusters_1, dtype: int64  STEP 5 = Plot the Principal Components of Hierarchical in 2D:  plt.suptitle('PC1 vs PC2 - Hierarchical cluster', size = 20)  sns.scatterplot(data = liquor_x_1, x = 'PC1', y = 'PC2', hue = 'Clusters_1')  PC1 vs PC2 - Hierarchical cluster  PC1 vs PC2 - Hierarchical cluster
	Custers 1 0 1 2 2 1 1 2 2 1 1 2 2 1 1 2 2 1 1 2 2 1 1 2 2 1 1 2 2 1 1 2 2 1 1 2 2 1 1 2 2 1 1 2 2 1 1 2 2 1 1 2 2 1 1 2 2 1 1 1 2 2 1 1 1 2 2 1 1 1 2 2 1 1 1 2 2 1 1 1 2 2 1 1 1 2 2 1 1 1 2 2 1
34]:	(b) K-MEANS  STEP 1 = Data Pre-processing:  import pandas as pd import numpy as np import matplotlib.pyplot as plt import seaborn as sns  import warnings warnings.filterwarnings('ignore')  liquor = pd.read_csv('wine.csv') liquor.head()
36]:	Type         Alcoho         Malic         Ash         Alcalinity         Magnesium         Pleanus         Flavancia         Proanthocyanins         Color         Hue         Dilution         Proline           0         1         14.23         1.71         2.43         15.6         127         2.80         3.06         0.28         2.29         5.64         1.04         3.92         1050           1         1         13.20         1.78         2.14         11.2         100         2.65         2.76         0.26         1.28         4.38         1.56         1.56         1.56         1.28         4.38         1.56         1.56         1.56         1.56         1.28         4.38         1.56         3.45         1.56         1.28         4.38         1.56         3.45         1.185         1.28         2.81         5.68         1.56         3.45         1.480         <
	2 2.516740 -1.031151 0.982819 3 3.757066 -2.756372 -0.176192 4 1.008908 -0.869831 2.026688 173 -3.370524 -2.216289 -0.342570 174 -2.601956 -1.757229 0.207581 175 -2.677839 -2.760899 -0.940942 176 -2.387017 -2.297347 -0.550696 177 -3.208758 -2.768920 1.013914
37]:	STEP 2 = Finding the optimal number of clusters using the Elbow Method :  from matplotlib import pyplot as plt from sklearn.cluster import KMeans  wcss = [] for i in range(1, 11):     kmeans = KMeans(n_clusters = i, random_state = 0)     kmeans.fit(liquor_x_1)     wcss.append(kmeans.inertia_) wcss.  [1539.5034801883066, 886.1611364823499, 512.9995067661521, 429.8030732473654, 371.02394449456744,
	326.2768904959661, 291.52701349085964, 261.6843759885969, 239.04461949060345, 213.23994073321188]  plt.plot(range(1, 11), wcss) plt.scatter(range(1, 11), wcss, color = 'red') plt.title('ELBOW METHOD') plt.xlabel('Number of clusters') plt.ylabel('wcss') plt.show()  ELBOW METHOD  1600  ELBOW METHOD
40]:	STEP 3 = Training the K - means clustering model :  km = KMeans(3, random_state = 42) km.fit(liquor_x_1)
40]:	<pre>km.fit(liquor_x_1)  km.labels_ array([1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1</pre>
42]:	liquor_x_1['Clusters_2'] = km.labels_ liquor_x_1  PC1 PC2 PC3 Clusters_2  0 3.316751 -1.443463 -0.165739
1 43]:	176 -2.387017 -2.297347 -0.550696 0 177 -3.208758 -2.768920 1.013914 0 78 rows × 4 columns  STEP 4 = Clusters with Number of Records :  liquor_x_1['Clusters_2'].value_counts() 2 65 1 62 0 51 Name: Clusters_2, dtype: int64
44]:	STEP 5 = Plot the Principal Components of KMEANS in 2D :¶  plt.suptitle('PC1 vs PC2 - K means clusters', size = 20)  sns.scatterplot(data = liquor_x_1, x = 'PC1', y = 'PC2', hue = 'Clusters_2')  plt.show()  PC1 vs PC2 - K means clusters