import numpy as np import pandas as pd from matplotlib import pyplot as plt from sklearn.cluster import KMeans from sklearn.preprocessing import StandardScaler from sklearn import preprocessing import seaborn as sns from sklearn.decomposition import PCA import scipy.cluster.hierarchy as sch from sklearn.cluster import AgglomerativeClustering from scipy.cluster.hierarchy import linkage df = pd.read\_csv('wine.csv') df1 = df.iloc[:, 1:] In [3]: df1.head() Alcohol Malic Ash Alcalinity Magnesium Phenols Flavanoids Nonflavanoids Proanthocyanins Color Hue Dilution Proline Out[4]: 1.71 2.43 14.23 15.6 127 2.80 3.06 0.28 2.29 5.64 1.04 3.92 1065 13.20 1.78 2.14 11.2 100 2.65 2.76 0.26 1.28 4.38 1.05 3.40 1050 13.16 2.36 2.67 18.6 101 2.80 3.24 0.30 2.81 5.68 1.03 3.17 1185 14.37 1.95 2.50 16.8 113 3.85 3.49 0.24 2.18 7.80 0.86 3.45 1480 13.24 2.59 2.87 21.0 2.80 0.39 118 2.69 1.82 4.32 1.04 2.93 735 df1.describe() In [5]: Out[5]: Alcohol Malic Ash Alcalinity Magnesium Phenols Flavanoids Nonflavanoids Proanthocyanins Color Hue Dilution Proline **count** 178.000000 178.000000 178.000000 178.000000 178.000000 178.000000 178.000000 178.000000 178.000000 178.000000 178.000000 178.000000 178.000000 13.000618 2.336348 2.366517 19.494944 99.741573 2.295112 2.029270 0.361854 1.590899 5.058090 0.957449 2.611685 746.893258 mean 0.811827 0.625851 0.998859 0.228572 1.117146 0.274344 3.339564 14.282484 0.124453 0.572359 2.318286 0.709990 314.907474 std min 11.030000 0.740000 1.360000 10.600000 70.000000 0.980000 0.340000 0.130000 0.410000 1.280000 0.480000 1.270000 278.000000 12.362500 1.602500 2.210000 17.200000 88.000000 1.742500 1.205000 0.270000 1.250000 3.220000 0.782500 1.937500 500.500000 25% 13.050000 1.865000 2.360000 19.500000 98.000000 2.355000 2.135000 0.340000 1.555000 4.690000 0.965000 2.780000 673.500000 985.000000 3.082500 21.500000 107.000000 2.800000 2.875000 0.437500 1.950000 6.200000 1.120000 3.170000 13.677500 2.557500 14.830000 5.800000 3.230000 30.000000 162.000000 3.880000 5.080000 0.660000 3.580000 13.000000 1.710000 4.000000 1680.000000 Correlation cor = df1.corr() cor.style.background\_gradient(cmap='coolwarm') Phenols Flavanoids Nonflavanoids Proanthocyanins Out[7]: Alcohol Malic Ash Alcalinity Magnesium Color Hue Dilution Proline 0.094397 0.211545 -0.310235 0.289101 0.236815 -0.155929 0.546364 -0.071747 0.072343 0.643720 Alcohol 1.000000 0.270798 0.136698 Malic 0.094397 1.000000 0.164045 0.288500 -0.054575 -0.335167 -0.411007 0.292977 -0.220746 0.248985 -0.561296 -0.368710 -0.192011 0.211545 0.164045 0.443367 0.286587 0.128980 0.115077 0.186230 0.258887 -0.074667 0.003911 0.223626 1.000000 0.009652 Ash 0.288500 0.443367 -0.273955 -0.276769 -0.440597 **Alcalinity** -0.310235 1.000000 -0.083333 -0.321113 -0.351370 0.361922 -0.197327 0.018732 0.270798 -0.054575 0.286587 -0.083333 1.000000 0.214401 0.195784 -0.256294 0.199950 0.055398 0.066004 0.393351 0.236441 Magnesium 0.864564 Phenols 0.289101 -0.335167 0.128980 -0.321113 0.214401 1.000000 -0.449935 0.612413 -0.055136 0.433681 0.699949 0.498115 0.864564 Flavanoids 0.236815 -0.411007 0.115077 -0.351370 0.195784 1.000000 -0.537900 0.652692 -0.172379 0.543479 0.787194 0.494193 Nonflavanoids -0.155929 0.292977 0.186230 0.361922 -0.256294 -0.449935 -0.537900 1.000000 -0.365845 0.139057 -0.262640 -0.503270 -0.311385 0.236441 0.612413 0.136698 -0.220746 0.009652 -0.197327 -0.365845 1.000000 -0.025250 0.295544 0.519067 0.330417 **Proanthocyanins** 0.652692 0.546364 0.248985 0.258887 0.018732 0.199950 -0.172379 0.139057 1.000000 -0.521813 -0.428815 0.316100 Color -0.055136 -0.071747 -0.561296 -0.074667 -0.262640 -0.521813 0.236183 Hue -0.273955 0.055398 0.433681 0.543479 0.295544 1.000000 0.565468 Dilution -0.368710 0.003911 -0.276769 0.066004 0.699949 0.787194 -0.503270 -0.428815 0.565468 1.000000 0.312761 0.519067 
 0.643720
 -0.192011
 0.223626
 -0.440597 0.393351 0.498115 0.494193 -0.311385 1.000000 Proline There are some quite correlation between variables. For example the correlation between flavanoids and dilution is pretty high (78%). Thus we can remove that variable from our dataset. However this method is long and tedious. Hence we PCA method for **Dimensionality Reduction** Dimensionality Reduction with PCA # normalizing the data df\_norm = StandardScaler().fit\_transform(df1)  $pca = PCA(n\_components=13)$ principalComponents = pca.fit\_transform(df\_norm) PC = range(1, pca.n\_components\_+1) plt.bar(PC, pca.explained\_variance\_ratio\_, color='blue') plt.xlabel('Principal Components') plt.ylabel('Variance %') plt.xticks(PC) Out[11]: ([<matplotlib.axis.XTick at 0x1ed9ab34220>, <matplotlib.axis.XTick at 0x1ed9ab341f0>, <matplotlib.axis.XTick at 0x1ed9a863160>, <matplotlib.axis.XTick at 0x1ed9ab70f10>, <matplotlib.axis.XTick at 0x1ed9ab82460>, <matplotlib.axis.XTick at 0x1ed9ab82970>, <matplotlib.axis.XTick at 0x1ed9ab82e80>, <matplotlib.axis.XTick at 0x1ed9ab843d0>, <matplotlib.axis.XTick at 0x1ed9ab848e0>, <matplotlib.axis.XTick at 0x1ed9ab84df0>, <matplotlib.axis.XTick at 0x1ed9ab89340>, <matplotlib.axis.XTick at 0x1ed9ab89850>, <matplotlib.axis.XTick at 0x1ed9ab84460>], [Text(0, 0, ''), Text(0, 0, ''),
Text(0, 0, ''), Text(0, 0, ''), Text(0, 0, ''), Text(0, 0, Text(0, 0, ''), Text(0, 0, '')]) 0.35 0.30 0.25 0.20 0.15 0.10 0.05 0.00 1 2 3 4 5 6 7 8 9 10 11 12 13 Principal Components PCA\_components = pd.DataFrame(principalComponents) plt.scatter(PCA\_components[0], PCA\_components[1], alpha=.3, color='blue') plt.xlabel('PCA 1') plt.ylabel('PCA 2') plt.show() PCA 2 PCA 1 As shown in the bar graph, the most of variance is put in the first 2 components. Since there is not much variance present from 3rd component, lets just the first 2 componets in our analysis. The scatter plot given an indication that there may be 3 clusters present Finding out the optimal number of clusters wcss = []In [14]: **for** i **in** range(1, 15): kmeans = KMeans(n\_clusters = i, init = 'k-means++', random\_state = 42) kmeans.fit(PCA\_components.iloc[:,:3]) wcss.append(kmeans.inertia\_) plt.plot(range(1, 15), wcss) plt.title('The Elbow Method') plt.xlabel('Number of clusters') plt.ylabel('WCSS') plt.show() The Elbow Method 1600 1400 1200 1000 800 600 400 200 The scree plot levels off at k=3 and let's use it to determine the clusters K clusters model = KMeans(n\_clusters=3) model.fit(PCA\_components.iloc[:,:2]) Out[18]: KMeans(n\_clusters=3) labels = model.predict(PCA\_components.iloc[:,:2]) plt.scatter(PCA\_components[0], PCA\_components[1], c=labels) -2 k\_new\_df=pd.DataFrame(principalComponents[:,0:2]) In [23]: model\_k = KMeans(n\_clusters=3) model\_k.fit(k\_new\_df) Out[23]: KMeans(n\_clusters=3) model\_k.labels\_ Out[24]: array([0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 2, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 2, 2, In [25]: md=pd.Series(model\_k.labels\_) df1['clust']=md k\_new\_df.head() Out[27]: 1 **0** 3.316751 -1.443463 **1** 2.209465 0.333393 **2** 2.516740 -1.031151 **3** 3.757066 -2.756372 4 1.008908 -0.869831 df1.groupby(df1.clust).mean() In [28]: Out[28]: Alcohol Malic Ash Alcalinity Magnesium Phenols Flavanoids Nonflavanoids Proanthocyanins Color **Hue Dilution Proline** clust **0** 13.659219 1.975781 2.463750 17.596875 107.312500 2.859688 3.012656 0.290000 1.921719 5.406250 1.069688 3.157188 1082.562500 **1** 12.238308 1.931385 2.219385 19.898462 0.365538 92.830769 2.204308 1.989231 1.587692 2.992615 1.051631 2.769231 506.353846 **2** 13.151633 3.344490 2.434694 21.438776 99.020408 1.678163 0.797959 0.450816 1.163061 7.343265 0.685918 1.690204 H clusters model2 = AgglomerativeClustering(n\_clusters=3, affinity='euclidean', linkage='ward') h\_cluster = model2.fit(PCA\_components.iloc[:,:2]) labels2 = model2.labels\_ X = PCA\_components.iloc[:,:1] Y = PCA\_components.iloc[:,1:2] plt.figure(figsize=(10, 7)) In [33]: plt.scatter(X, Y, c=labels2) Out[33]: <matplotlib.collections.PathCollection at 0x1ed9b5c5c10> -1 -2 -3 h\_new\_df=pd.DataFrame(principalComponents[:,0:2]) h\_new\_df.head() 1 Out[35]: **0** 3.316751 -1.443463 **1** 2.209465 0.333393 **2** 2.516740 -1.031151 **3** 3.757066 -2.756372 4 1.008908 -0.869831 hcf = linkage(h\_new\_df, method="complete", metric="euclidean") plt.figure(figsize=(15, 5));plt.title('Hierarchical Clustering Dendrogram');plt.xlabel('Index');plt.ylabel('Distance') In [37]: sch.dendrogram( hcf, leaf\_rotation=0., leaf\_font\_size=8., plt.show() Hierarchical Clustering Dendrogram Index h\_complete = AgglomerativeClustering(n\_clusters=5,linkage='complete',affinity = "euclidean").fit(h\_new\_df) h\_complete.labels\_ In [39]: 3, 3, 3, 3, 4, 3, 3, 4, 3, 3, 3, 3, 4, 2, 0, 0, 2, 2, 2, 2, 0, 2, 0, 2, 0, 3, 2, 2, 2, 0, 2, 2, 2, 2, 2, 0, 2, 2, 2, 0, 0, 0, 2, 2, 3, 0, 2, 2, 2, 2, 2, 2, 2, 2, 0, 2, 0, 2, 0, 2, 2, 2, 2, 0, 2, 2, 3, 0, 0, 2, 2, 2, 0, 2, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 1, 1, 1, 1, 1, 0, 1, 1, 0, 1, 1, 0, 1, 0, 0, 1, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 1], dtype=int64) cluster\_labels=pd.Series(h\_complete.labels\_) df1['clust']=cluster\_labels In [41]: df1.head() Alcohol Malic Ash Alcalinity Magnesium Phenols Flavanoids Nonflavanoids Proanthocyanins Color Hue Dilution Proline clust Out[42]: 1.71 2.43 2.80 0.28 15.6 127 3.06 2.29 5.64 1.04 3.92 1065 13.20 1.78 2.14 11.2 100 2.65 2.76 0.26 1050 1.28 4.38 1.05 3.40 2.36 2.67 101 2.80 0.30 13.16 18.6 3.24 2.81 5.68 1.03 3.17 1185 1.95 2.50 16.8 113 1480 14.37 3.85 3.49 7.80 0.86 3.45 2.59 2.87 118 2.80 735 13.24 2.69 4.32 1.04 2.93 df1.groupby(df1.clust).mean() Out[43]: Malic Alcalinity Magnesium Phenols Flavanoids Nonflavanoids Proanthocyanins Color Dilution **Proline** Alcohol Hue clust **0** 12.686222 2.924000 2.380667 20.966667 95.000000 1.691556 1.088222 0.469556 4.613556 0.824578 2.021111 21.847826 0.897826 0.428696 0.640000 **1** 13.420435 3.390435 2.486087 103.130435 1.777391 1.406522 9.203043 1.640870 654.782609 **2** 12.238163 1.713061 2.165714 19.528571 92.061224 2.395510 2.204694 0.318163 2.992653 1.096327 2.913265 509.755102 **3** 13.573265 2.016531 2.481429 18.069388 107.795918 2.758571 0.296122 5.175510 1.065510 3.166939 2.887551 **4** 14.150000 1.963333 2.435000 15.150000 109.500000 3.248333 3.505833 0.276667 2.185000 6.735000 1.055833 3.188333 1289.333333 Conclusion Using PCA we reduced the variables to only 2 from 13 and use clustering classification, we can safely assume that there exists 3 cluster in the wine data sets