: [	1
ſ	Data columns (total 13 columns):  # Column Non-Null Count Dtype
: '	DF1  array([[1.423e+01, 1.710e+00, 2.430e+00,, 1.040e+00, 3.920e+00, 1.065e+03], [1.320e+01, 1.780e+00, 2.140e+00,, 1.050e+00, 1.050e+03], [1.316e+01, 2.360e+00, 2.670e+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], [1.317e+01, 2.590e+00, 2.370e+00,, 6.000e-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]])
: [	<pre># Normalization function def norm_func(i):     x = (i-i.min())/(i.max()-i.min())     return (x)  df_norm = norm_func(DF1)  df_norm  array([[8.39350664e-03, 9.40548971e-04, 1.36915357e-03,,     5.41708585e-04, 2.25612696e-03, 6.33900242e-01],     [7.78036396e-03, 9.82218862e-04, 1.19652116e-03,,     5.47661426e-04, 1.94657920e-03, 6.24970980e-01],     [7.5655259e-03, 1.32748367e-03, 1.51202176e-03,,     5.35755743e-04, 1.80966384e-03, 7.05334341e-01],    ,     [7.82203385e-03, 2.47042926e-03, 1.26795526e-03,,     2.73830713e-04, 8.51256347e-04, 4.96984886e-01],</pre>
: [	[7.76250543e-03, 1.46439903e-03, 1.33343652e-03,, 2.79783555e-04, 8.86973397e-04, 4.99961307e-01], [8.33397822e-03, 2.36327811e-03, 1.55369165e-03,, 2.85736396e-04, 8.75067714e-04, 3.33281742e-01]])  pca = PCA(n_components = 13) pca_values = pca.fit_transform(df_norm)  pca_values  array([[ 1.89635495e-01,
: -	-1.32777740e-04, 5.71761518e-05, 2.15343422e-05], [5.56330203e-02, 1.1114442e-02, -1.06460114e-03,, -1.81013009e-05, -3.10079415e-05, 7.90293761e-05], [-1.11284320e-01, -1.26992448e-04, -3.35175331e-03,, 9.56770872e-05, 1.90015753e-05, 1.60351688e-05]])  # The amount of variance that each PCA explains is var = pca.explained_variance_ratio_ var  array([9.98091230e-01, 1.73591562e-03, 9.49589576e-05, 5.02173562e-05, 1.23636847e-05, 8.46213034e-06, 2.80681456e-06, 1.52308053e-06, 1.12783044e-06, 7.21415811e-07, 3.78060267e-07, 2.12013755e-07, 8.25392788e-08])  # Cumulative variance var1 = np.cumsum(np.round(var,decimals = 4)*100) var1
:	array([ 99.81, 99.98, 99.99, 100. , 100. , 100. , 100. , 100. , 100. , 100. , 100. , 100. , 100. ])  # Variance plot for PCA components obtained plt.plot(var1,color="red")  [ <matplotlib.lines.line2d 0x1115a48="" at="">]  100.000</matplotlib.lines.line2d>
: [	finaldf = pd.concat([pd.DataFrame(pca_values[:,0:3],columns=['pc1','pc2','pc3']), df[['Type']]], axis = 1)  finaldf   pc1 pc2 pc3 Type  0 0.189635  0.012794  0.001864
1 •	173 -0.004155 -0.002703 -0.001473 3  174 0.001864 0.001390 -0.002566 3  175 0.052658 0.011177 -0.001332 3  176 0.055633 0.011114 -0.001065 3  177 -0.111284 -0.000127 -0.003352 3  178 rows × 4 columns  Performing K-Means Clustering  Using with Elbow method to find the optimum no of clusters
	<pre>wcss = [] for i in range(1, 11):     kmeans = KMeans(n_clusters = i, init = 'k-means++', random_state = 200)     kmeans.fit(finaldf)     wcss.append(kmeans.inertia_)  plt.figure(figsize=(10, 8))     plt.plot(range(1, 11), wcss)     plt.title('The Elbow Method')     plt.ylabel('Number of clusters')     plt.ylabel('wCSS')     plt.show()</pre> The Elbow Method
	80 - 80 - 40 - 20 -
: [	As per this plot, optimum number of clusters 3  To confirm the same, using Silhouette score's method  for i in range(2,12):     labels=cluster.KMeans(n_clusters=i,init="k-means++",random_state=200).fit(finaldf).labels_     print ("silhouette score for k(clusters) = "+str(i)+" is "
<i>A</i>	Silhouette score for k(clusters) = 4 is 0.7934893991312625 Silhouette score for k(clusters) = 5 is 0.68220630351393922 Silhouette score for k(clusters) = 6 is 0.6914427573862947 Silhouette score for k(clusters) = 7 is 0.5852242596504906 Silhouette score for k(clusters) = 8 is 0.5991039010559479 Silhouette score for k(clusters) = 9 is 0.6002899535946947 Silhouette score for k(clusters) = 10 is 0.580685430101755 Silhouette score for k(clusters) = 11 is 0.5854630977632279  As per the Silhouette score also, K= 3 is the optimum number of clusters  #Hence model=KMeans(n_clusters=3) model.fit(finaldf) model.labels_  array([1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1
: [	1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1
	1 13.744746 2.010678 2.455593 17.037288 106.338983 2.840169 2 13.153750 3.333750 2.437083 21.416667 99.312500 1.678750  plt.figure(figsize=(10,4)) sns.histplot (x='kclust', data=df) plt.xlabel('Cluster') plt.ylabel('Freq') plt.suptitle('Relative comparison of people in respective clusters')  Text(0.5, 0.98, 'Relative comparison of people in respective clusters')  Relative comparison of people in respective clusters
	60 - 50 - 50 - 75 - 70 - 70 - 70 - 70 - 7
:	<pre># create dendrogram plt.figure(figsize=(15, 10)) dendrogram = sch.dendrogram(sch.linkage(finaldf, method='ward'))</pre> 12 10
	8 - 6 - 4 -
	The x-axis contains the samples and y-axis represents the distance between these samples. The vertical line with maximum distance is the blue line. If we decide a threshold of 2 and cut the dendrogram:  plt.figure(figsize=(15, 10)) plt.title("Dendrograms") dend = sch.dendrogram(sch.linkage(finaldf, method='ward')) plt.axhline(y=2, color='r', linestyle='') plt.show()
	Dendrograms  12 -
	4-
: [	Now, testing for optimum no of clusters in the case of original dataset (without PCA)  df_norm1 = norm_func(df1)  Performing K-Means Clustering  Using elbow method:  wcss1 = [] for k in range(1, 11):
:	<pre>kmeans1 = KMeans(n_clusters = k, init = 'k-means++', random_state = 200) kmeans1.fit(df_norm1) wcss1.append(kmeans1.inertia_)  plt.figure(figsize=(10, 8)) plt.plot(range(1, 11), wcss1) plt.vitle('The Elbow Method') plt.xlabel('Number of clusters') plt.ylabel('WcSs') plt.show()</pre> The Elbow Method
	80 - 70 - 85 60 -
	50 -
7 T	
	As per the elbow plot, it seems that the optimum number of clusters = 3  To confirm the optimum no of clusters, using silhouettes  for i in range(2,13):     labels=cluster.KMeans(n_clusters=i,init="k-means++",random_state=200).fit(df_norm1).labels_     print ("Silhouette score for k(clusters) = "+str(i)+" is "
	As per the elbow plot, it seems that the optimum number of clusters = 3  To confirm the optimum no of clusters, using silhouettes  for i in range(2,13):     labels=cluster.kMeans(n_clusters=i,init="k-means++",random_state=260).fit(df_norm1).labels_     print ("Silhouette score for k(clusters) = "*str(1)*" is "