

# DWM Programming Assignment

Name - Vishal Kriplani (69)

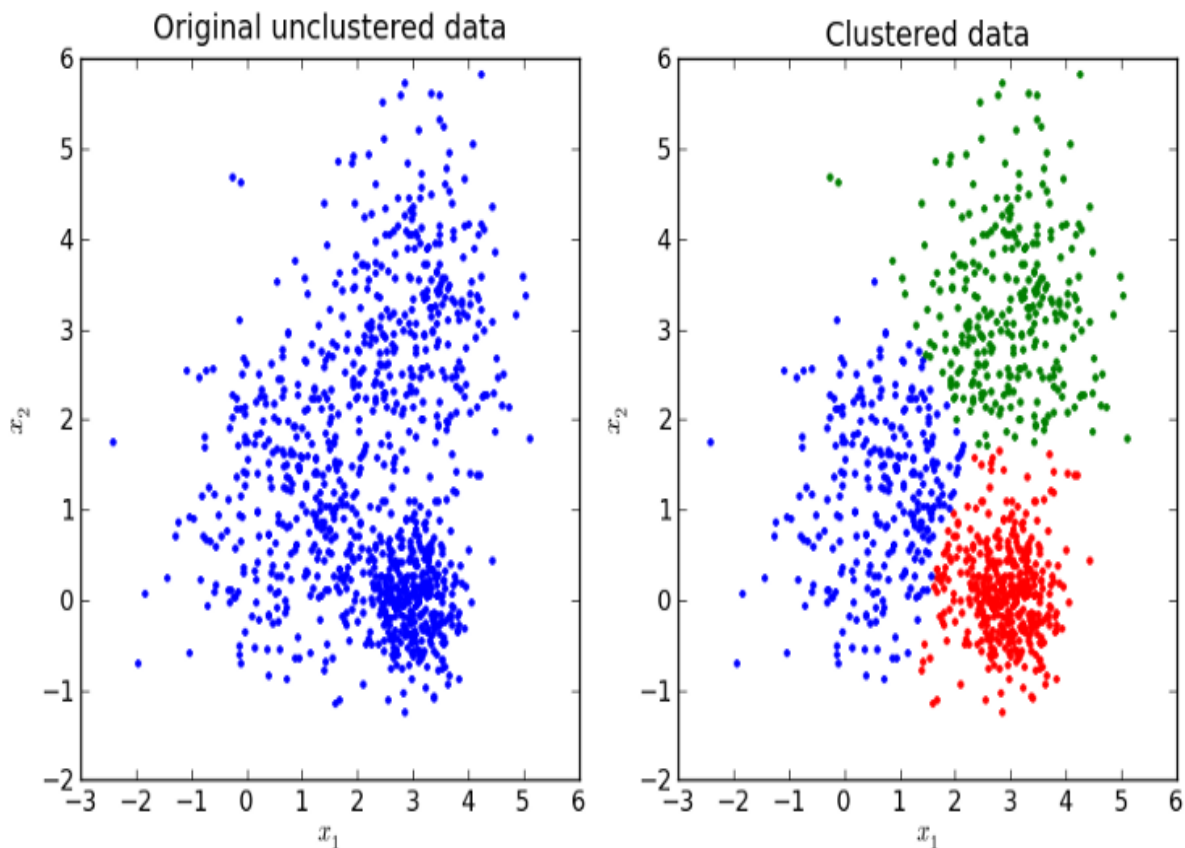
Piyush Chotiya (49)

Ganesh Karwa (39)

## Topic - Single Linkage Hierarchical Clustering

### What is Clustering??

Clustering is basically a technique that groups similar data points such that the points in the same group are more similar to each other than the points in the other groups. The group of similar data points is called a **Cluster**.



### Hierarchical clustering Technique:

Hierarchical clustering is one of the popular and easy to understand clustering technique. This clustering technique is divided into two types:

1. Agglomerative
2. Divisive

**Agglomerative Hierarchical clustering Technique:** In this technique, initially each data point is considered as an individual cluster. At each iteration, the similar clusters merge with other clusters until one cluster or K clusters are formed.

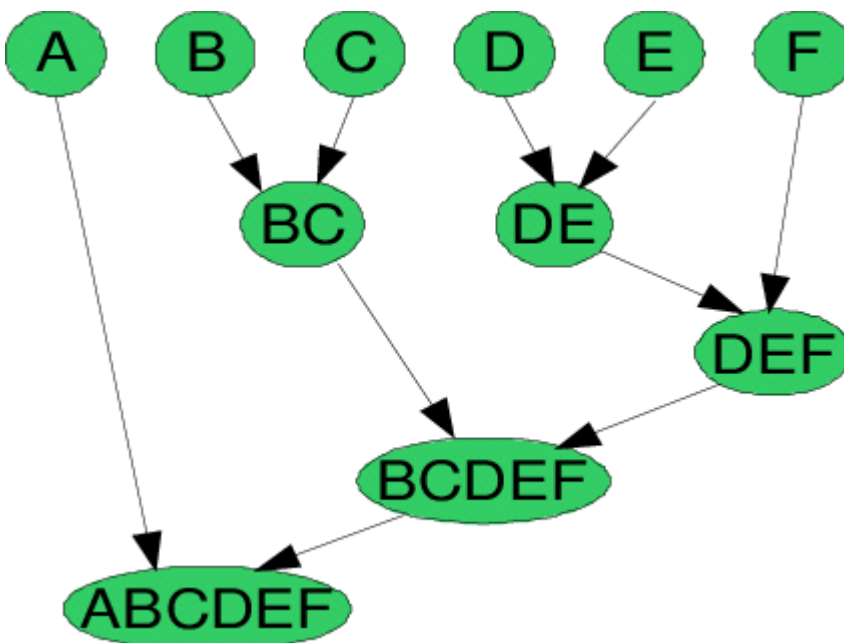
The basic algorithm of Agglomerative is straight forward.

- Compute the proximity matrix
- Let each data point be a cluster
- Repeat: Merge the two closest clusters and update the proximity matrix
- Until only a single cluster remains

Key operation is the computation of the proximity of two clusters

To understand better let's see a pictorial representation of the Agglomerative Hierarchical clustering Technique. Lets say we have six data points {A,B,C,D,E,F}.

- Step- 1: In the initial step, we calculate the proximity of individual points and consider all the six data points as individual clusters as shown in the image below.

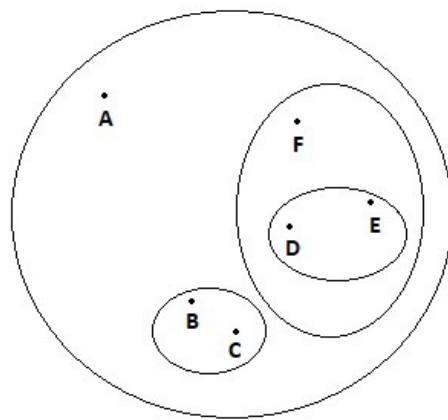
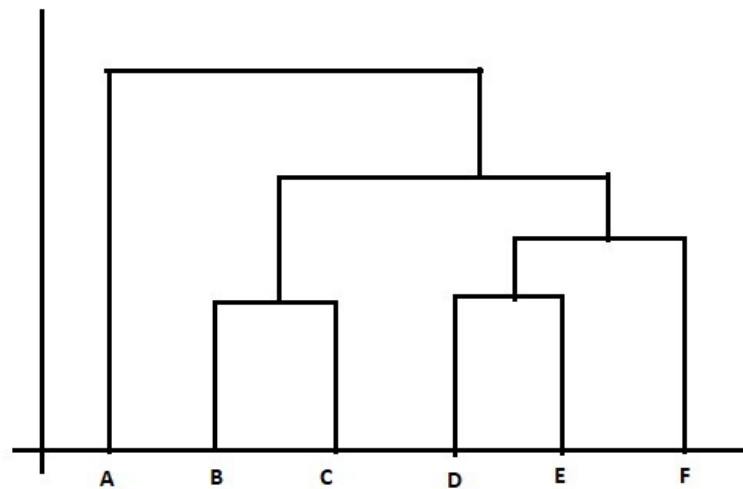


Agglomerative Hierarchical Clustering Technique

- Step- 2: In step two, similar clusters are merged together and formed as a single cluster. Let's consider B,C, and D,E are similar clusters that are merged in step two. Now, we're left with four clusters which are A, BC, DE, F.
- Step- 3: We again calculate the proximity of new clusters and merge the similar clusters to form new clusters A, BC, DEF.
- Step- 4: Calculate the proximity of the new clusters. The clusters DEF and BC are similar and merged together to form a new cluster. We're now left with two clusters A, BCDEF.
- Step- 5: Finally, all the clusters are merged together and form a single cluster.

The Hierarchical clustering Technique can be visualized using a **Dendrogram**.

A **Dendrogram** is a tree-like diagram that records the sequences of merges or splits.



Dendrogram representation

**2. Divisive Hierarchical clustering Technique:** Since the Divisive Hierarchical clustering Technique is not much used in the real world, I'll give a brief of the Divisive Hierarchical clustering Technique.

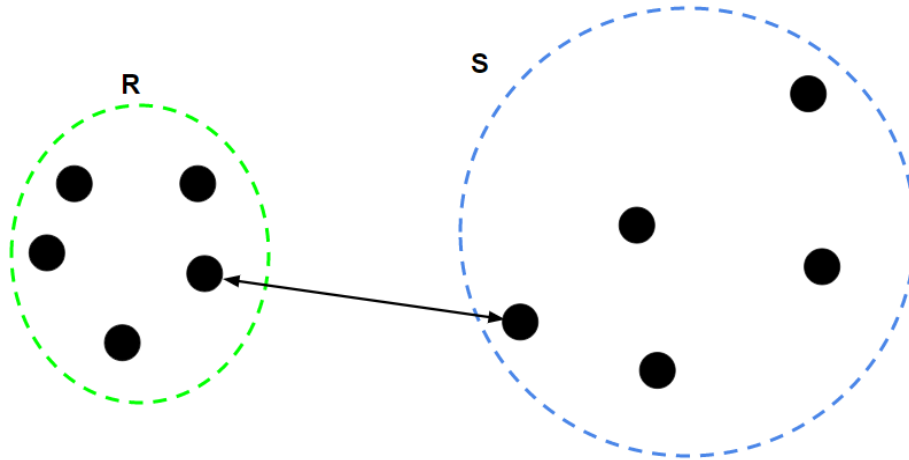
In simple words, we can say that the Divisive Hierarchical clustering is exactly the opposite of the **Agglomerative Hierarchical clustering**. In Divisive Hierarchical clustering, we consider all the data points as a single cluster and in each iteration, we separate the data points from the cluster which

are not similar. Each data point which is separated is considered as an individual cluster. In the end, we'll be left with n clusters.

As we're dividing the single clusters into n clusters, it is named as **Divisive Hierarchical clustering**.

**Single Linkage:** For two clusters R and S, the single linkage returns the minimum distance between two points i and j such that i belongs to R and j belongs to S.

$$L(R, S) = \min(D(i, j)), i \in R, j \in S$$



## ▼ Problem statement

---

You are owning a supermarket mall and through membership cards, you have some basic data about your customers like Customer ID, age, gender, annual income and spending score. Problem Statement You own the mall and want to understand the customers like who can be easily converge [Target Customers] so that the sense can be given to marketing team and plan the strategy accordingly.

## ▼ All Imports

```
1 %matplotlib inline
2 import matplotlib.pyplot as plt
3 import numpy as np
4 import random
5 import seaborn as sns
```

```

6 import scipy as sp
7 from sklearn import datasets
8 from numpy import linalg as LA
9 import scipy.cluster.hierarchy as shc
10 from sklearn.metrics.pairwise import pairwise_distances
11 import sys

```

## ▼ Reading Raw Mall Customer Data from csv

```

1 import pandas as pd
2
3 url = 'https://raw.githubusercontent.com/Vk76/Single-Linkage-Hierarchical-Clust
4 df = pd.read_csv(url)
5 data = df.iloc[:, [3,4]].values

```

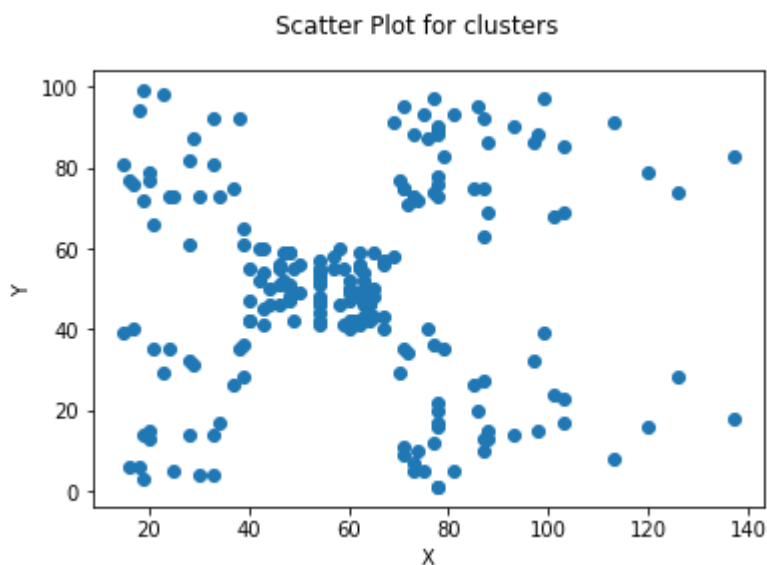
## ▼ Plotting all the data set Points

```

1 fig = plt.figure()
2 fig.suptitle('Scatter Plot for clusters')
3 ax = fig.add_subplot(1,1,1)
4 ax.set_xlabel('X')
5 ax.set_ylabel('Y')
6 ax.scatter(data[:,0],data[:,1])

```

<matplotlib.collections.PathCollection at 0x7f4391c0d0f0>



## ▼ Function for plotting the clusters with different color

---

input -

**data** : two dimensional data frame

**linkage** : single

**no\_of\_clusters** : clusters to make

```
1 def hierarchical_clustering(data,linkage,no_of_clusters):
2     #first step is to calculate the initial distance matrix
3     #it consists distances from all the point to all the point
4     color = ['r','b','g','c','m','y','k','w']
5
6     initial_distances = pairwise_distances(data,metric='euclidean')
7     #making all the diagonal elements infinity
8     np.fill_diagonal(initial_distances,sys.maxsize)
9     clusters = find_clusters(initial_distances,linkage)
10
11     #plotting the clusters
12     iteration_number = initial_distances.shape[0] - no_of_clusters
13     clusters_to_plot = clusters[iteration_number]
14     arr = np.unique(clusters_to_plot)
15
16     indices_to_plot = []
17     fig = plt.figure()
18     fig.suptitle('Scatter Plot for clusters')
19     ax = fig.add_subplot(1,1,1)
20     ax.set_xlabel('X')
21     ax.set_ylabel('Y')
22     for x in np.nditer(arr):
23         indices_to_plot.append(np.where(clusters_to_plot==x))
24     p=0
25
26     clusters.pop(5)
27     for i in range(0,len(indices_to_plot)):
28         for j in np.nditer(indices_to_plot[i]):
29             ax.scatter(data[j,0],data[j,1], c= color[p])
30             p = p + 1
31     plt.title('Clusters of customers using Hierarchical Clustering')
32     plt.xlabel('Annual Income (K$)')
33     plt.ylabel('Spending Score (1-100)')
34     plt.legend()
35     plt.show()
36
37
```

## ▼ Function for finding the clusters in sets

---

input -

**data** : two dimensional data frame

**linkage** : single

output -

**clusters** : two dimensional list of clusters

```
1 def find_clusters(input,linkage):
2     clusters = {}
3     row_index = -1
4     col_index = -1
5     array = []
6
7
8     for n in range(input.shape[0]):
9         array.append(n)
10
11     clusters[0] = array.copy()
12
13     #finding minimum value from the distance matrix
14     #note that this loop will always return minimum value from bottom triangle
15     for k in range(1, input.shape[0]):
16         min_val = sys.maxsize
17
18         for i in range(0, input.shape[0]):
19             for j in range(0, input.shape[1]):
20                 if(input[i][j]<=min_val):
21                     min_val = input[i][j]
22                     row_index = i
23                     col_index = j
24
25         #once we find the minimum value, we need to update the distance matrix
26         #updating the matrix by calculating the new distances from the cluster
27
28         #for Single Linkage
29         if(linkage == "single" or linkage == "Single"):
30             for i in range(0,input.shape[0]):
31                 if(i != col_index):
32                     #we calculate the distance of every data point from newly f
33                     temp = min(input[col index][i],input[row index][i])
```

```

34         #we update the matrix symmetrically as our distance matrix
35         input[col_index][i] = temp
36         input[i][col_index] = temp
37
38     #set the rows and columns for the cluster with higher index i.e. the row
39     #Set input[row_index][for_all_i] = infinity
40     #set input[for_all_i][row_index] = infinity
41     for i in range (0,input.shape[0]):
42         input[row_index][i] = sys.maxsize
43         input[i][row_index] = sys.maxsize
44
45     #Manipulating the dictionary to keep track of cluster formation in each
46     #if k=0,then all datapoints are clusters
47
48     minimum = min(row_index,col_index)
49     maximum = max(row_index,col_index)
50     for n in range(len(array)):
51         if(array[n]==maximum):
52             array[n] = minimum
53     clusters[k] = array.copy()
54
55     return clusters

```

## Plotting the graph for 8 different cluster values using single linkage hierarchical clustering

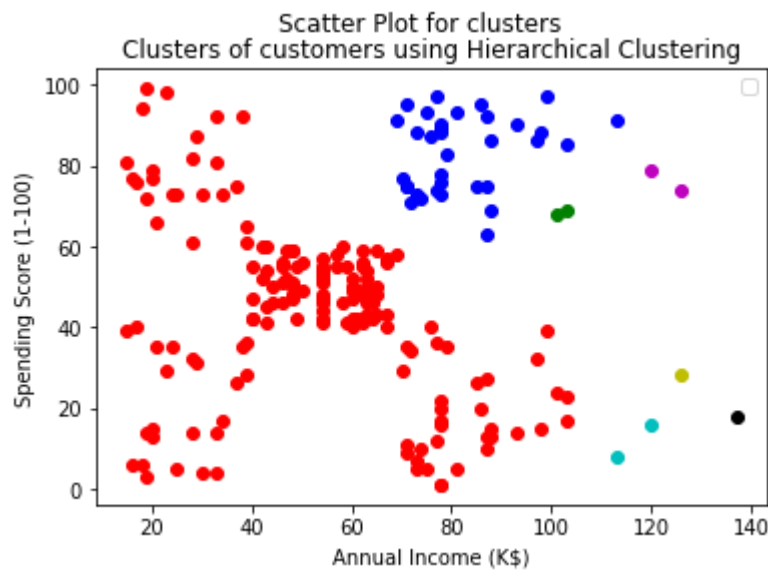
```

1 # Here we visualize all clusters starting from half of the length of points corresponding to
2 number_of_clusters = 8
3
4 # Here in we just go on decreasing the number of clusters. The points in the scatter plot
5 while number_of_clusters > 0:
6     print(number_of_clusters)
7     hierarchical_clustering(data,"single",number_of_clusters)
8     number_of_clusters -= 1
9

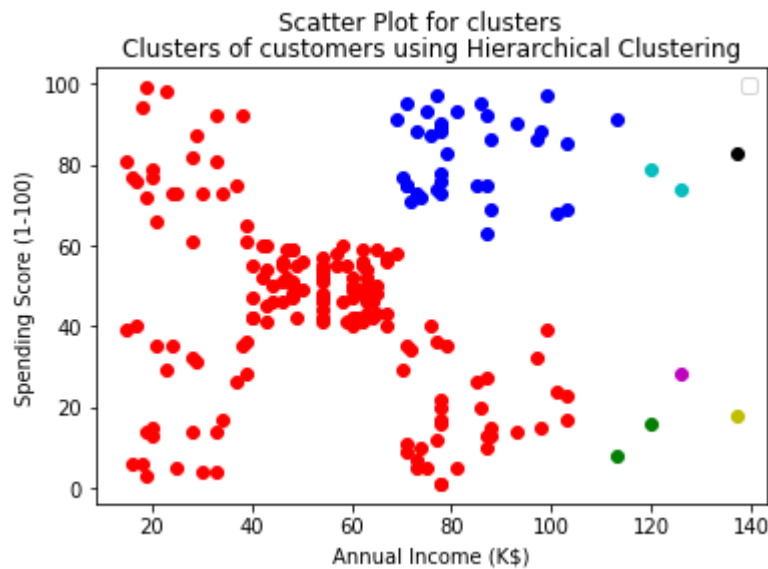
```



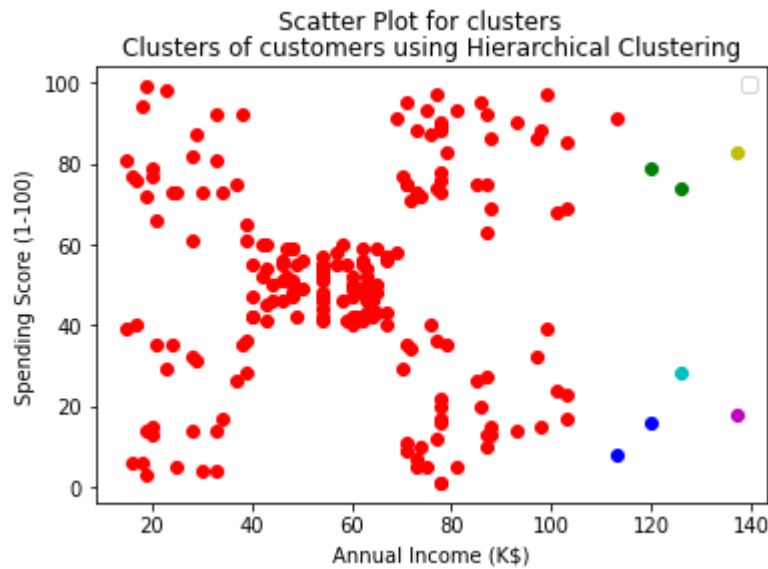
8  
No handles with labels found to put in legend.



7  
No handles with labels found to put in legend.

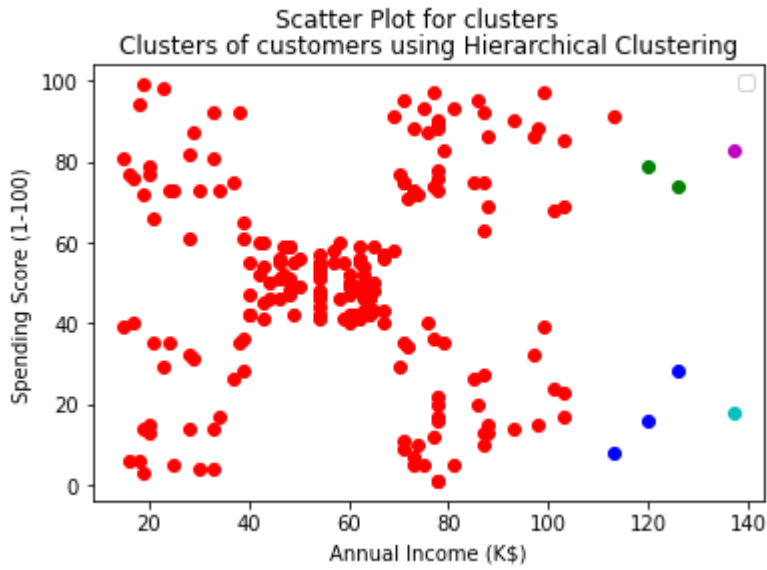


6  
No handles with labels found to put in legend.



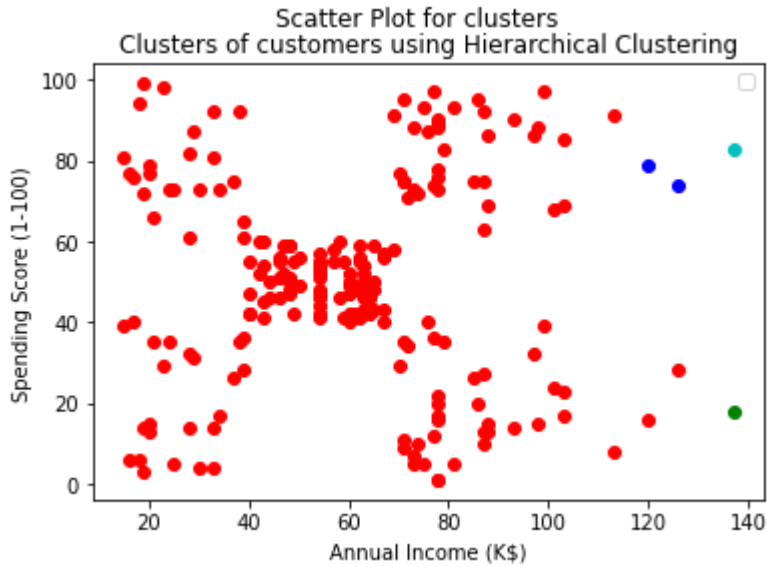
5

No handles with labels found to put in legend.



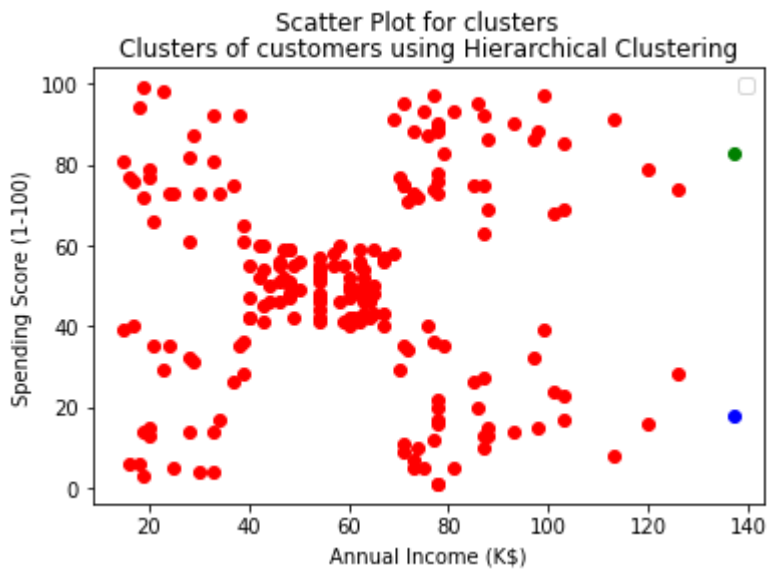
4

No handles with labels found to put in legend.



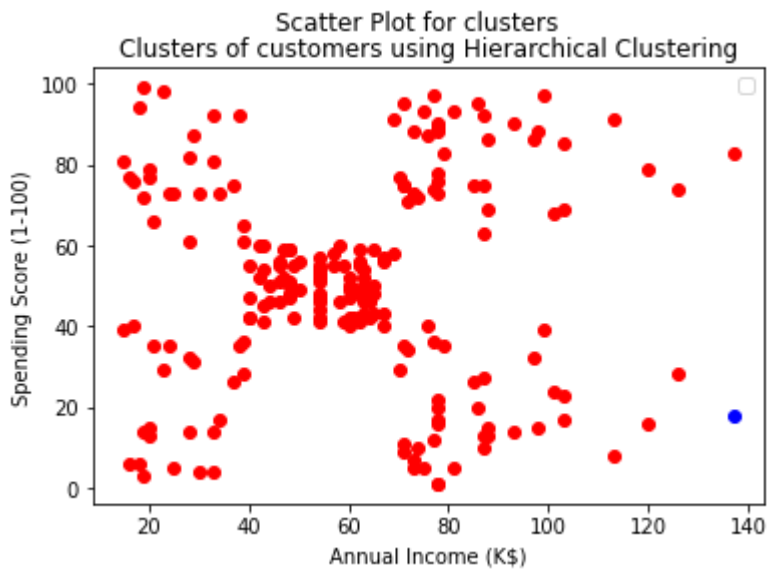
3

No handles with labels found to put in legend.



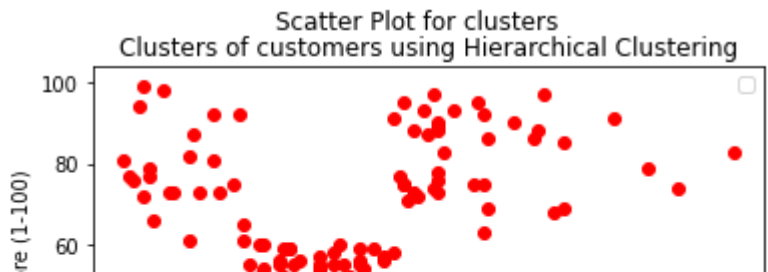
2

No handles with labels found to put in legend.



1

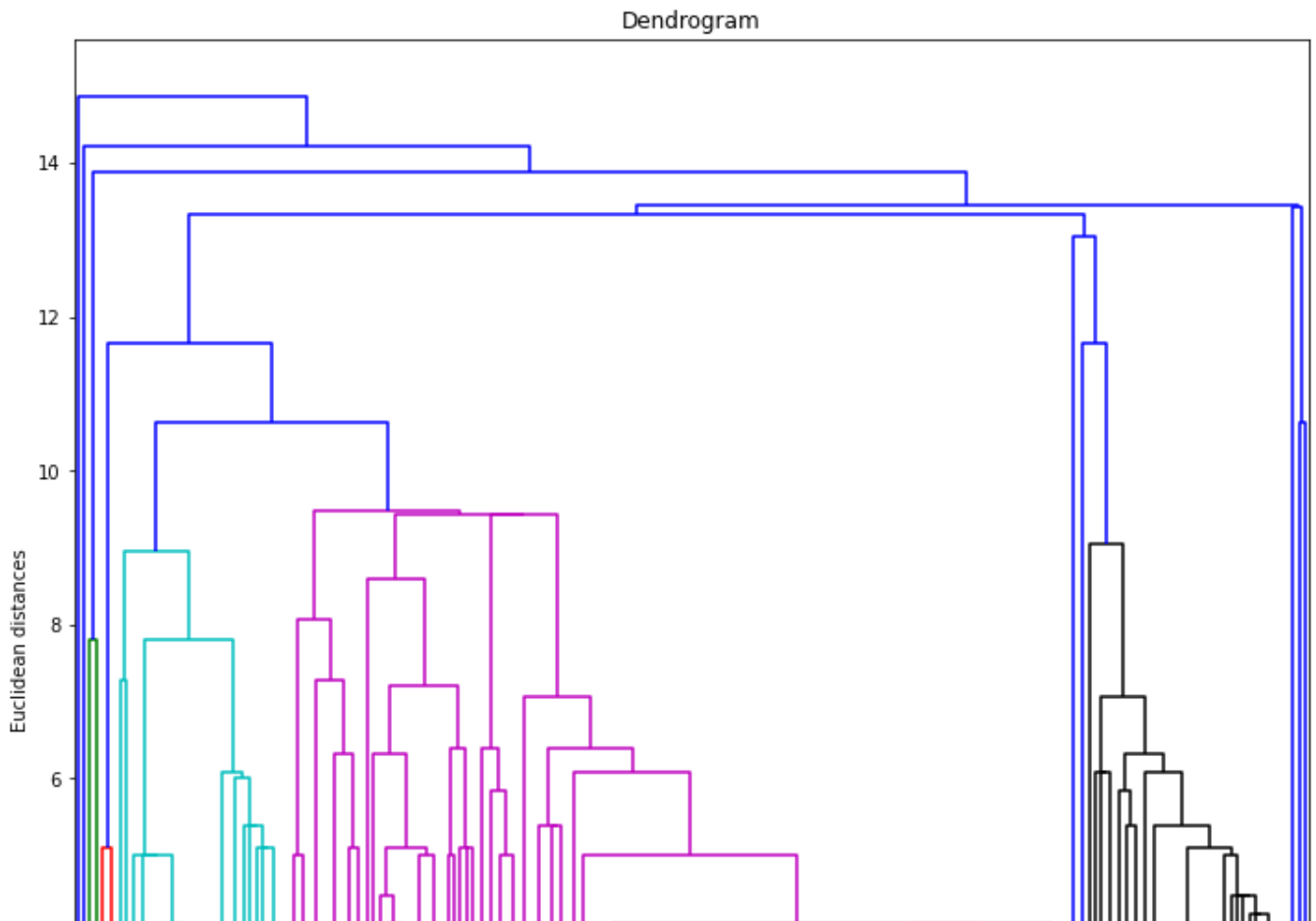
No handles with labels found to put in legend.



## ▼ Dendrogram for the the dataset



```
1 plt.figure(figsize=(12,12))
2 dendrogram = shc.dendrogram(shc.linkage(data,method = "single"))
3 plt.title('Dendrogram')
4 plt.xlabel('Customers')
5 plt.ylabel('Euclidean distances')
6 plt.show()
```



## Conclusions

(Results are on basis of single linkage for 8 clusters)

- **Red** color shows unusual behaviour of customers that have low income and high spending and high income and low spending so they are sensible and careless customers.
- **Blue and Green** color has a stable customers and they are **targeted customers** as they have good income and their spending is high.
- **Cyan and yellow** coloured customers are **careful customers** as they have high income but less spending.

