# Hierarchical Clustering

## Example 1

In our first example we will cluster the X numpy array of data points that we created in the previous section.

The process of clustering is similar to any other unsupervised machine learning algorithm. We start by importing the required libraries:

```python
In [14]:   import matplotlib.pyplot as plt
           import pandas as pd
           %matplotlib inline
           import numpy as np
```

```python
In [15]:   X = np.array([[5,3],
               [10,15],
               [15,12],
               [24,10],
               [30,30],
               [85,70],
               [71,80],
               [60,78],
               [70,55],
               [80,91],])
```

The next step is to import the class for clustering and call its fit_predict method to predict the clusters that each data point belongs to.

Take a look at the following script:

```python
In [16]:   from sklearn.cluster import AgglomerativeClustering

           cluster = AgglomerativeClustering(n_clusters=2, affinity='euclidean', linkage='
           cluster.fit_predict(X)
```

```
Out[16]:   array([1, 1, 1, 1, 1, 0, 0, 0, 0, 0])
```

In the code above we import the AgglomerativeClustering class from the "sklearn.cluster" library. The number of parameters is set to 2 using the n_clusters parameter while the affinity is set to "euclidean" (distance between the datapoints). Finally linkage parameter is set to "ward", which minimizes the variant between the clusters.

Next we call the fit_predict method from the AgglomerativeClustering class variable cluster. This method returns the names of the clusters that each data point belongs to. Execute the following script to see how the data points have been clustered.

```python
In [17]:   print(cluster.labels_)
```
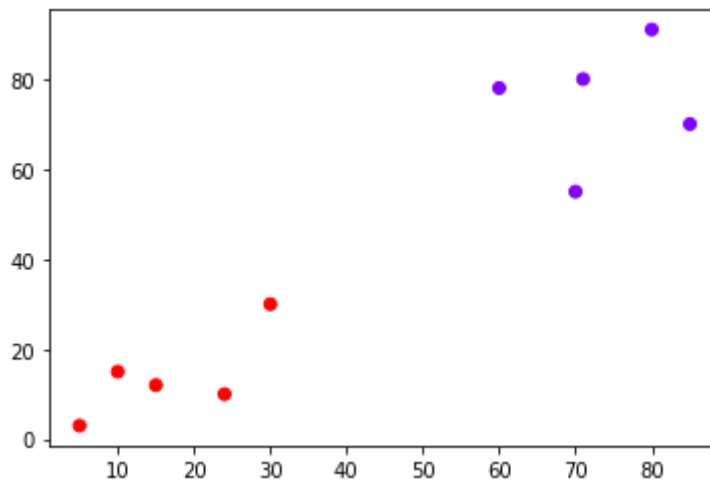
```
[1 1 1 1 1 0 0 0 0 0]
```

The output is a one-dimensional array of 10 elements corresponding to the clusters assigned to our 10 data points.

As expected the first five points have been clustered together while the last five points have been clustered together. It is important to mention here that these ones and zeros are merely labels assigned to the clusters and have no mathematical implications.

Finally, let's plot our clusters. To do so, execute the following code:

```python
In [18]: plt.scatter(X[:,0],X[:,1], c=cluster.labels_, cmap='rainbow')
```

```
Out[18]: <matplotlib.collections.PathCollection at 0x7fd7d1f08970>
```



# Example 2

In the last section we performed hierarchical clustering on dummy data. In this example, we will perform hierarchical clustering on real-world data and see how it can be used to solve an actual problem.

The problem that we are going to solve in this section is to segment customers into different groups based on their shopping trends.

```python
In [19]: import matplotlib.pyplot as plt
         import pandas as pd
         %matplotlib inline
         import numpy as np
```

```python
In [20]: customer_data = pd.read_csv('shopping_data.csv')
```

```python
In [21]: customer_data.head()
```

Out[21]:

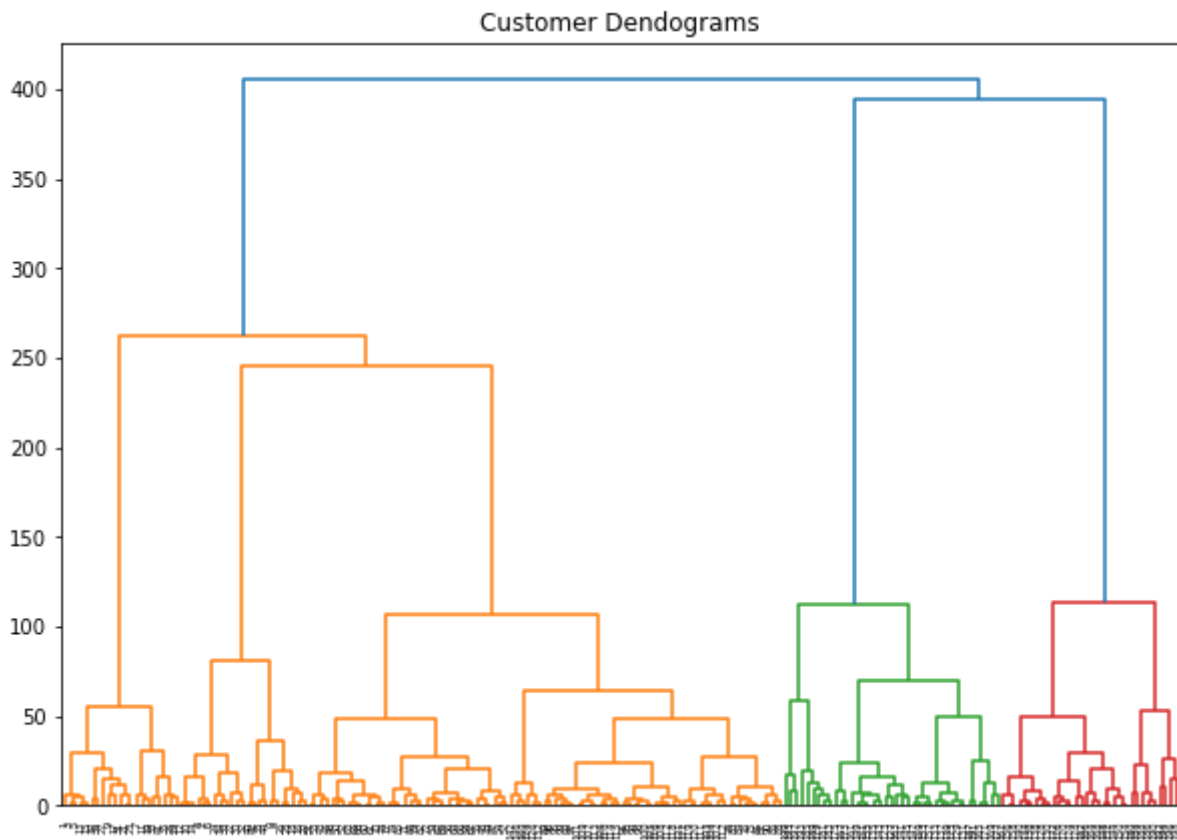| | CustomerID | Genre | Age | Annual Income (k$) | Spending Score (1-100) |
|---|---|---|---|---|---|
| **0** | 1 | Male | 19 | 15 | 39 |
| **1** | 2 | Male | 21 | 15 | 81 |
| **2** | 3 | Female | 20 | 16 | 6 |
| **3** | 4 | Female | 23 | 16 | 77 |
| **4** | 5 | Female | 31 | 17 | 40 |

Our dataset has five columns: CustomerID, Genre, Age, Annual Income, and Spending Score. To view the results in two-dimensional feature space, we will retain only two of these five columns. We can remove CustomerID column, Genre, and Age column. We will retain the Annual Income (in thousands of dollars) and Spending Score (1-100) columns. The Spending Score column signifies how often a person spends money in a mall on a scale of 1 to 100 with 100 being the highest spender. Execute the following script to filter the first three columns from our dataset:

In [22]:
```python
data = customer_data.iloc[:, 3:5].values
```

Next, we need to know the clusters that we want our data to be split to. We will again use the scipy library to create the dendrograms for our dataset. Execute the following script to do so:

In [23]:
```python
import scipy.cluster.hierarchy as shc

plt.figure(figsize=(10, 7))
plt.title("Customer Dendograms")
dend = shc.dendrogram(shc.linkage(data, method='ward'))
```

## Customer Dendograms



In the script above we import the hierarchy class of the scipy.cluster library as shc. The hierarchy class has a dendrogram method which takes the value returned by the linkage method of the same class. The linkage method takes the dataset and the method to minimize distances as parameters. We use 'ward' as the method since it minimizes then variants of distances between the clusters.

The output of the script above looks like the graph above.

Now we know the number of clusters for our dataset, the next step is to group the data points into these five clusters. To do so we will again use the AgglomerativeClustering class of the sklearn.cluster library. Take a look at the following script:

```
In [24]:   from sklearn.cluster import AgglomerativeClustering

           cluster = AgglomerativeClustering(n_clusters=5, affinity='euclidean', linkage='
           cluster.fit_predict(data)
```

```
Out[24]:   array([4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3,
                  4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 1,
                  4, 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, 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, 1, 1, 1, 1, 1, 1, 1,
                  1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 2, 1, 2, 1, 2, 0, 2, 0, 2,
                  1, 2, 0, 2, 0, 2, 0, 2, 0, 2, 1, 2, 0, 2, 1, 2, 0, 2, 0, 2, 0, 2,
                  0, 2, 0, 2, 0, 2, 1, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2,
                  0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2,
                  0, 2])
```
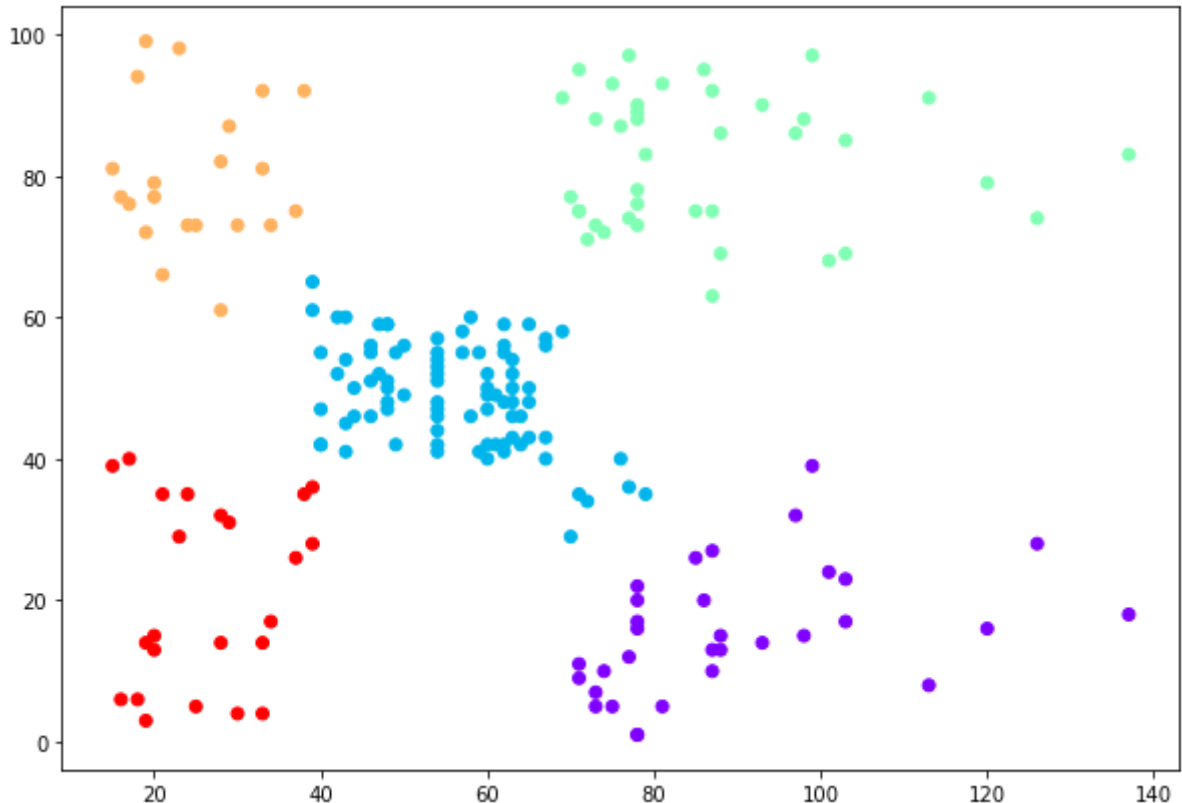
You can see the cluster labels from all of your data points. Since we had five clusters, we

have five labels in the output i.e. 0 to 4.

As a final step, let's plot the clusters to see how actually our data has been clustered:

```
In [25]:  plt.figure(figsize=(10, 7))
          plt.scatter(data[:,0], data[:,1], c=cluster.labels_, cmap='rainbow')
```

Out[25]:  `<matplotlib.collections.PathCollection at 0x7fd7d191ce80>`

You can see the data points in the form of five clusters. The data points in the bottom right belong to the customers with high salaries but low spending. These are the customers that spend their money carefully. Similarly, the customers at top right (green data points), these are the customers with high salaries and high spending. These are the type of customers that companies target. The customers in the middle (blue data points) are the ones with average income and average salaries. The highest numbers of customers belong to this category. Companies can also target these customers given the fact that they are in huge numbers, etc.

```
In [ ]:
```