

## LAB 06 - Clustering Methods - Comparison ¶

### SUBMITTED BY:

Name : Kavitha.S

Reg no : 21122033

Class : 2MSDS

### LAB OVERVIEW:

Illustrate KMeans and Agglomerative Hierarchical Clustering on Iris Dataset, considering only two features - Sepal Length and Petal Width.

Use Elbow Method as a way to find optimum number of clusters

### PROBLEM DEFINITION:

To perform K-Means and Agglomerative Hierarchical Clustering on the in-built iris dataset. To find the optimum number of clusters.

### APPROACH:

Imported the necessary packages. Analyzed the dataset, drop the unnecessary columns. performed K-Means clustering on 2 features with various number of clusters. Found the optimum number of clusters using Elbow method. For Agglomerative hierarchical clustering, found the optimum number of clusters using silhouette score.

```
In [97]: 1 from sklearn import datasets
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import scipy.cluster.hierarchy as shc
5 from sklearn.cluster import AgglomerativeClustering
6 from sklearn.metrics import silhouette_score
7 import seaborn as sns
```

```
In [3]: 1 iris = datasets.load_iris()
```

```
In [11]: 1 print("Features:",iris.feature_names)
2 print("\nTarget:",iris.target_names)
```

```
Features: ['sepal length (cm)', 'sepal width (cm)', 'petal length (cm)', 'petal
width (cm)']
```

```
Target: ['setosa' 'versicolor' 'virginica']
```

```
In [56]: 1 iris_df=pd.DataFrame(iris.data)
          2 iris_df['Species']=iris.target
          3 iris_df.head()
```

Out[56]:

	0	1	2	3	Species
0	5.1	3.5	1.4	0.2	0
1	4.9	3.0	1.4	0.2	0
2	4.7	3.2	1.3	0.2	0
3	4.6	3.1	1.5	0.2	0
4	5.0	3.6	1.4	0.2	0

```
In [57]: 1 x = iris_df.iloc[:, [0, 1, 2, 3]].values
```

```
In [59]: 1 iris_df.columns=['SepalLength', 'SepalWidth', 'PetalLength', 'PetalWidth', 'Species']
          2 iris_df.head()
```

Out[59]:

	SepalLength	SepalWidth	PetalLength	PetalWidth	Species
0	5.1	3.5	1.4	0.2	0
1	4.9	3.0	1.4	0.2	0
2	4.7	3.2	1.3	0.2	0
3	4.6	3.1	1.5	0.2	0
4	5.0	3.6	1.4	0.2	0

```
In [60]: 1 iris_df.describe()
```

Out[60]:

	SepalLength	SepalWidth	PetalLength	PetalWidth	Species
count	150.000000	150.000000	150.000000	150.000000	150.000000
mean	5.843333	3.057333	3.758000	1.199333	1.000000
std	0.828066	0.435866	1.765298	0.762238	0.819232
min	4.300000	2.000000	1.000000	0.100000	0.000000
25%	5.100000	2.800000	1.600000	0.300000	0.000000
50%	5.800000	3.000000	4.350000	1.300000	1.000000
75%	6.400000	3.300000	5.100000	1.800000	2.000000
max	7.900000	4.400000	6.900000	2.500000	2.000000

```
In [61]: 1 iris_df.shape
```

Out[61]: (150, 5)

```
In [62]: 1 iris_df = iris_df.drop(['SepalWidth', 'PetalLength', 'Species'], axis=1)
```

```
In [63]: 1 iris_df.head()
```

Out[63]:

	SepalLength	PetalWidth
0	5.1	0.2
1	4.9	0.2
2	4.7	0.2
3	4.6	0.2
4	5.0	0.2

```
In [64]: 1 iris_df.SepalLength = iris_df.SepalLength.reshape(1, -1)
```

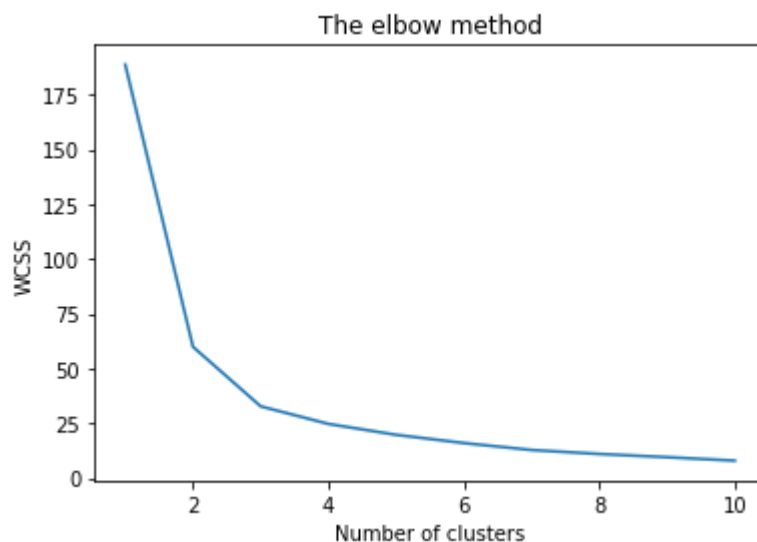
...

```
In [65]: 1 from sklearn.cluster import KMeans
2 wcss = []
3
4 for i in range(1, 11):
5     kmeans = KMeans(n_clusters = i, init = 'k-means++', max_iter = 300, n_in
6     kmeans.fit(iris_df)
7     wcss.append(kmeans.inertia_)
```

C:\Users\SRIDHAR\anaconda3\lib\site-packages\sklearn\cluster\\_kmeans.py:881: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP\_NUM\_THREADS=1.

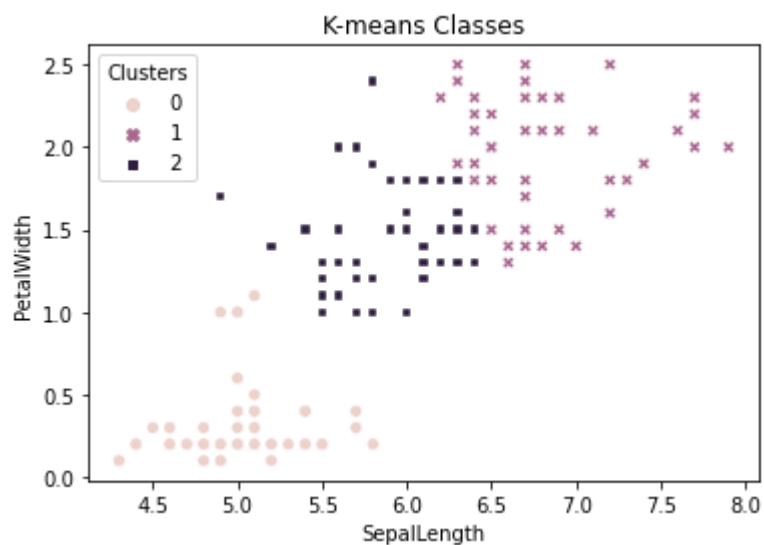
warnings.warn(

```
In [66]: 1 plt.plot(range(1, 11), wcss)
2 plt.title('The elbow method')
3 plt.xlabel('Number of clusters')
4 plt.ylabel('WCSS')
5 plt.show()
```

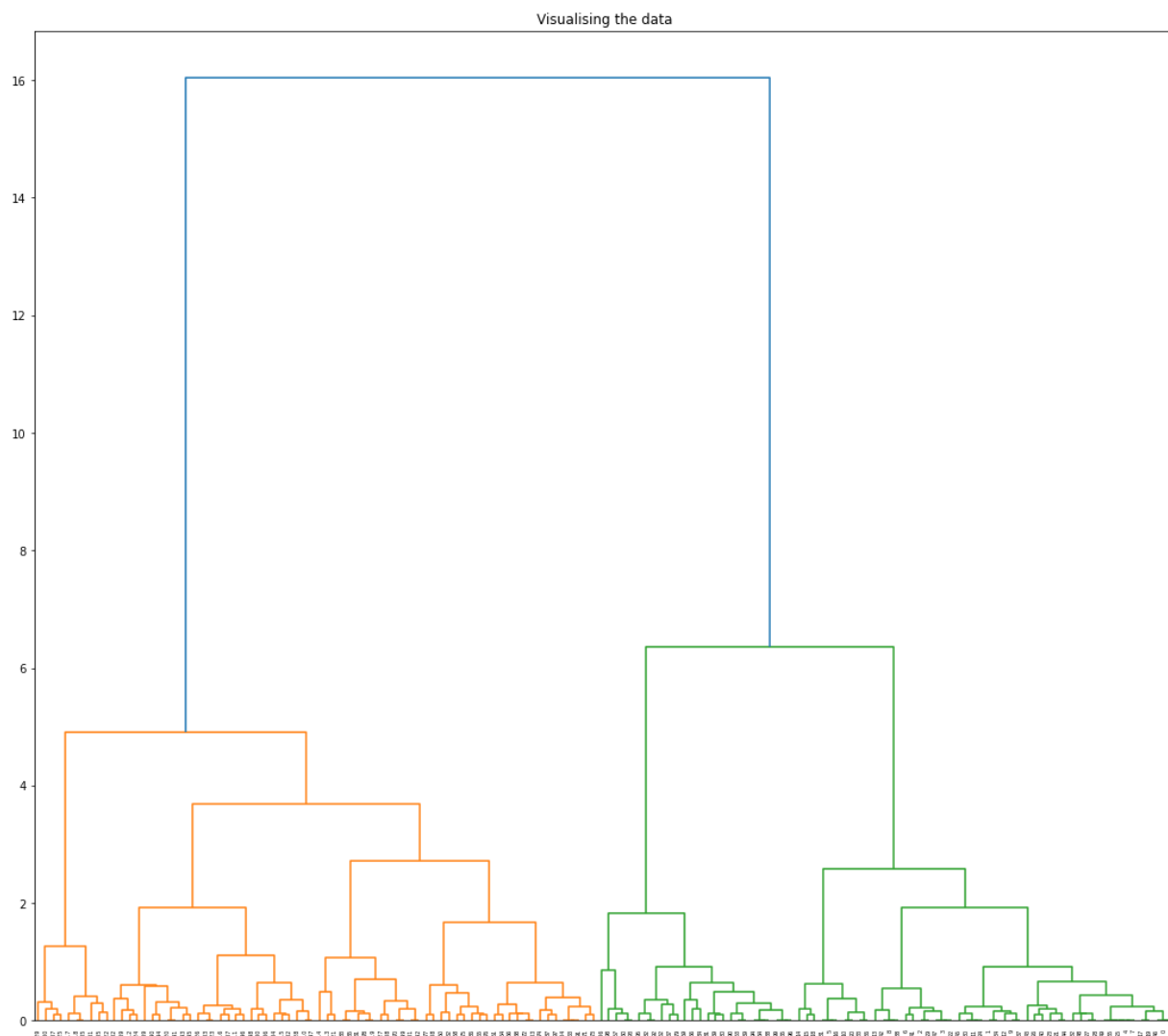


```
In [69]: 1 kmeans = KMeans(n_clusters = 3, init = 'k-means++', max_iter = 300, n_init =
2 identified_clusters = kmeans.fit_predict(iris_df)
```

```
In [102]: 1 plt.title('K-means Classes')
2 sns.scatterplot(x=data_with_clusters['SepalLength'], y=data_with_clusters['P
3 plt.show()
```



```
In [74]: 1 plt.figure(figsize =(18, 16))  
2 plt.title('Visualising the data')  
3 Dendrogram = shc.dendrogram((shc.linkage(iris_df, method ='ward')))
```



```
In [ ]: 1 df_h = df.copy(deep=True)  
2 df_h['labels_h'] = cluster2
```

```
In [80]: 1 ac2 = AgglomerativeClustering(n_clusters = 2)
        2 ac2.fit_predict(iris_df)
```

```
Out[80]: array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                0, 0, 0, 0, 0, 0, 1, 1, 1, 0, 1, 0, 1, 0, 1, 0, 1, 0, 1, 1,
                0, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 1, 0, 1, 1, 1,
                0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 1, 1, 1, 1, 1, 0, 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], dtype=int64)
```

```
In [81]: 1 ac3 = AgglomerativeClustering(n_clusters = 3)
        2 ac3.fit_predict(iris_df)
```

```
Out[81]: 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, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
                1, 1, 1, 1, 1, 0, 0, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0,
                2, 2, 0, 2, 0, 0, 0, 0, 0, 0, 0, 0, 0, 2, 2, 2, 2, 0, 2, 0, 0, 0,
                2, 2, 2, 0, 2, 2, 2, 2, 2, 0, 2, 2, 0, 0, 0, 0, 0, 0, 2, 0, 0, 0,
                0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0], dtype=int64)
```

```
In [82]: 1 ac4 = AgglomerativeClustering(n_clusters = 4)
        2 ac4.fit_predict(iris_df)
```

```
Out[82]: 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, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
                1, 1, 1, 1, 1, 0, 0, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0,
                2, 2, 0, 2, 0, 0, 0, 0, 0, 0, 0, 0, 0, 2, 2, 2, 2, 0, 2, 0, 0, 0,
                2, 2, 2, 0, 2, 2, 2, 2, 2, 0, 2, 2, 0, 0, 0, 0, 0, 3, 2, 3, 0, 0,
                0, 0, 0, 0, 0, 0, 0, 3, 3, 0, 0, 0, 3, 0, 0, 3, 0, 0, 0, 3, 3, 3,
                0, 0, 0, 3, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0], dtype=int64)
```

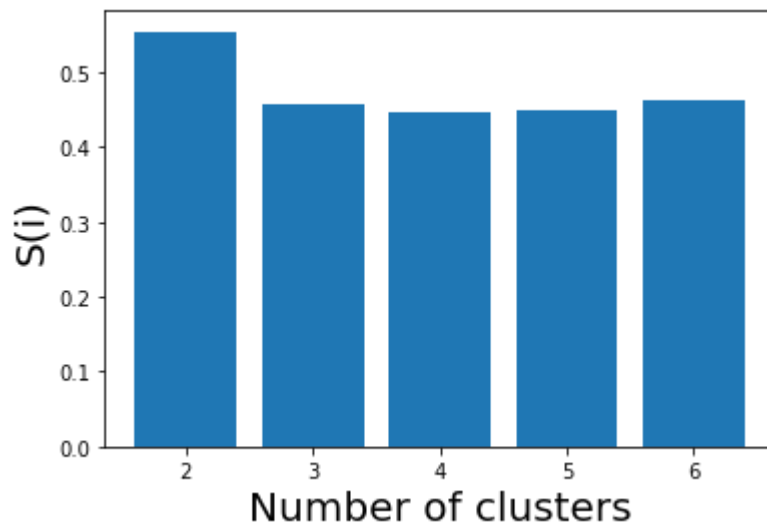
```
In [83]: 1 ac5 = AgglomerativeClustering(n_clusters = 5)
        2 ac5.fit_predict(iris_df)
```

```
Out[83]: 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, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
                1, 1, 1, 1, 1, 0, 0, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0,
                2, 2, 0, 2, 0, 0, 0, 0, 0, 0, 0, 0, 0, 2, 2, 2, 2, 0, 2, 0, 0, 0,
                2, 2, 2, 0, 2, 2, 2, 2, 2, 0, 2, 2, 4, 0, 4, 4, 4, 3, 2, 3, 0, 4,
                4, 4, 4, 0, 0, 4, 4, 3, 3, 0, 4, 0, 3, 4, 4, 3, 4, 0, 4, 3, 3, 3,
                4, 0, 0, 3, 4, 4, 0, 4, 4, 4, 0, 4, 4, 4, 4, 4, 4, 4, 0], dtype=int64)
```

```
In [84]: 1 ac6 = AgglomerativeClustering(n_clusters = 6)
        2 ac6.fit_predict(iris_df)
```

```
Out[84]: array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                0, 0, 0, 0, 0, 4, 4, 4, 2, 4, 2, 4, 2, 4, 2, 2, 5, 2, 4, 2, 4,
                2, 2, 4, 2, 5, 4, 4, 4, 4, 4, 4, 4, 5, 2, 2, 2, 2, 5, 2, 5, 4, 4,
                2, 2, 2, 4, 2, 2, 2, 2, 2, 4, 2, 2, 1, 5, 1, 1, 1, 3, 2, 3, 4, 1,
                1, 1, 1, 5, 5, 1, 1, 3, 3, 5, 1, 5, 3, 1, 1, 3, 1, 5, 1, 3, 3, 3,
                1, 4, 4, 3, 1, 1, 5, 1, 1, 1, 5, 1, 1, 1, 1, 1, 1, 5], dtype=int64)
```

```
In [89]: 1 k = [2, 3, 4, 5, 6]
2 silhouette_scores = []
3 silhouette_scores.append(
4     silhouette_score(iris_df, ac2.fit_predict(iris_df)))
5 silhouette_scores.append(
6     silhouette_score(iris_df, ac3.fit_predict(iris_df)))
7 silhouette_scores.append(
8     silhouette_score(iris_df, ac4.fit_predict(iris_df)))
9 silhouette_scores.append(
10    silhouette_score(iris_df, ac5.fit_predict(iris_df)))
11 silhouette_scores.append(
12    silhouette_score(iris_df, ac6.fit_predict(iris_df)))
13
14 plt.bar(k, silhouette_scores)
15 plt.xlabel('Number of clusters', fontsize = 20)
16 plt.ylabel('S(i)', fontsize = 20)
17 plt.show()
```



## INFERENCE:

Thus, with the help of the silhouette scores, it is concluded that the optimal number of clusters for the given data and clustering technique is 2.

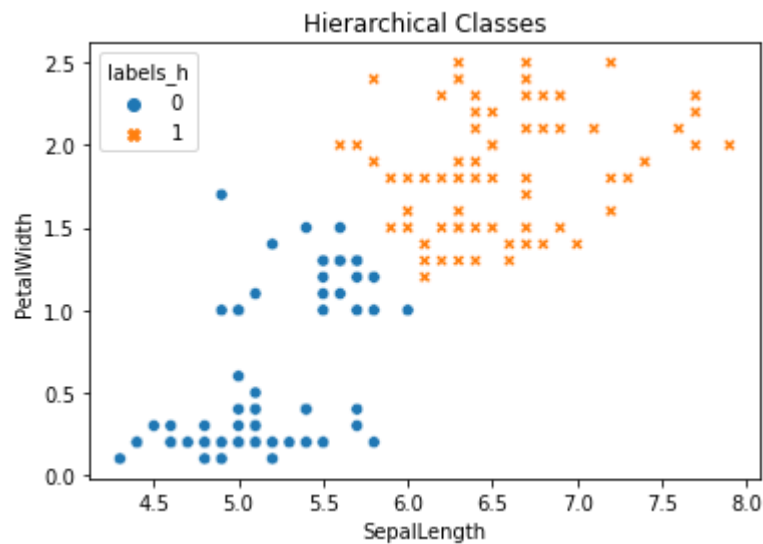
```
In [93]: 1 hie_clus = AgglomerativeClustering(n_clusters=2, affinity='euclidean', linka
2 cluster2 = hie_clus.fit_predict(iris_df)
3 df_h = iris_df.copy(deep=True)
4 df_h['labels_h'] = cluster2
```

In [95]: 1 df\_h.head()

Out[95]:

	SepalLength	PetalWidth	labels_h
0	5.1	0.2	0
1	4.9	0.2	0
2	4.7	0.2	0
3	4.6	0.2	0
4	5.0	0.2	0

In [103]: 1 plt.title('Hierarchical Classes')  
2 sns.scatterplot(x=df\_h['SepalLength'], y=df\_h['PetalWidth'], hue='labels\_h',  
3 plt.show())



## CONCLUSION:

For K-means, the optimum number of clusters was 3 and for Hierarchical, it was 2.

## REFERENCES:

<https://www.geeksforgeeks.org/elbow-method-for-optimal-value-of-k-in-kmeans/>  
[\(https://www.geeksforgeeks.org/elbow-method-for-optimal-value-of-k-in-kmeans/\)](https://www.geeksforgeeks.org/elbow-method-for-optimal-value-of-k-in-kmeans/)  
<https://www.kaggle.com/code/khotijahs1/hierarchical-agglomerative-clustering>  
 [\(https://www.kaggle.com/code/khotijahs1/hierarchical-agglomerative-clustering\)](https://www.kaggle.com/code/khotijahs1/hierarchical-agglomerative-clustering)  
<https://www.kaggle.com/code/khotijahs1/k-means-clustering-of-iris-dataset/notebook>  
 [\(https://www.kaggle.com/code/khotijahs1/k-means-clustering-of-iris-dataset/notebook\)](https://www.kaggle.com/code/khotijahs1/k-means-clustering-of-iris-dataset/notebook)