

# Introduction About Unsupervised Machine Learning.

Algorithms belonging to the family of **Unsupervised Learning** have **no variable to predict** tied to the data. Instead of having an output, the data only has an input which would be multiple variables that describe the data. This is where clustering comes in.

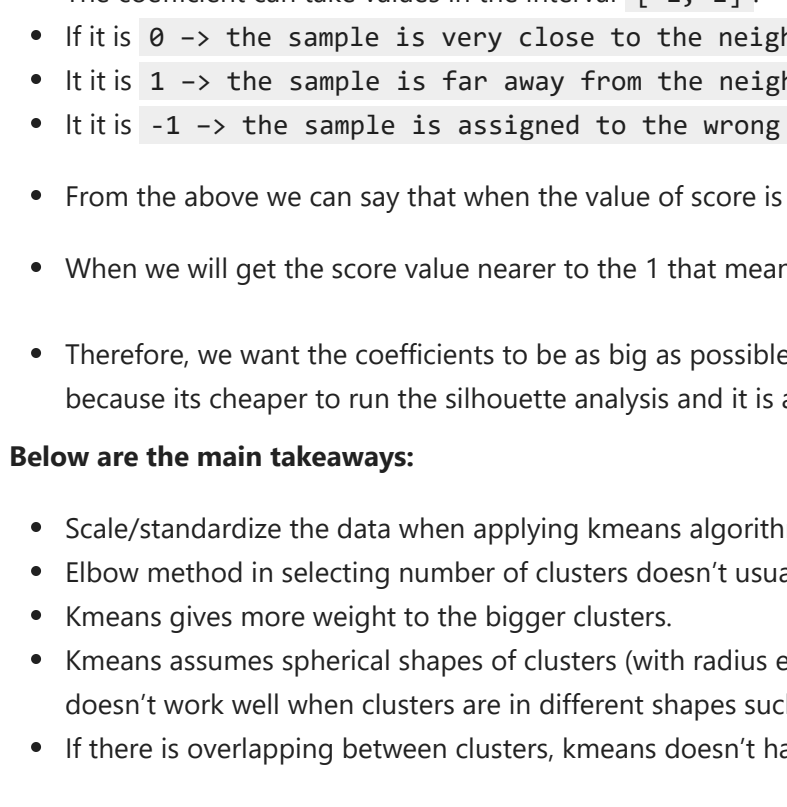
- Clustering is one of the most common **exploratory data analysis technique** used to get an intuition about the structure of the data. It can be defined as the task of **identifying subgroups** in the data such that data points in the **same subgroup (cluster)** are **very similar while data points in different clusters are very different**. In other words, we try to find homogeneous subgroups within the data such that data points in each cluster are as similar as possible according to a **similarity measure such as euclidean-based distance or correlation-based distance**. The decision of which similarity measure to use is application-specific.
- Clustering analysis can be done on the basis of features where we try to find subgroups of samples based on features or on the basis of samples where we try to find subgroups of features based on samples. We'll cover here different clustering algorithms. Clustering is used in market segmentation; where we try to find customers that are similar to each other whether in terms of behaviors or attributes, image segmentation/compression; where we try to group similar regions together, document clustering based on topics, etc.
- Unlike supervised learning, clustering is considered an unsupervised learning method since we don't have the ground truth to compare the output of the clustering algorithm to the true labels to evaluate its performance. We only want to try to investigate the structure of the data by grouping the data points into distinct subgroups.
- Clustering is the task of grouping together a set of objects in a way that objects in the same cluster are more similar to each other than to objects in other clusters. Similarity is a metric that reflects the strength of relationship between two data objects. Clustering is mainly used for exploratory data mining. It has manifold usage in many fields such as machine learning, pattern recognition, image analysis, information retrieval, bio-informatics, data compression, and computer graphics.

## 1. K-Means Clustering Algorithm

- Kmeans algorithm is an iterative algorithm that tries to partition the dataset into Kpre-defined distinct non-overlapping subgroups (clusters) where each data point belongs to only one group. It tries to make the intra-cluster data points as similar as possible while also keeping the clusters as different (far) as possible. It assigns data points to a cluster such that the sum of the squared distance between the data points and the cluster's centroid (arithmetic mean of all the data points that belong to that cluster) is at the minimum. The less variation we have within clusters, the more homogeneous (similar) the data points are within the same cluster.
- The way kmeans algorithm works is as follows:
  - Specify number of clusters K.
  - Initialize centroids by first shuffling the dataset and then randomly selecting K data points for the centroids without replacement.
  - Keep iterating until there is no change to the centroids. I.e assignment of data points to clusters isn't changing.
  - Compute the sum of the squared distance between data points and all centroids.
  - Assign each data point to the closest cluster (centroid).
  - Compute the centroids for the clusters by taking the average of the all data points that belong to each cluster.

### 1.1 Working Principle:-

- The k-means clustering algorithm attempts to split a given anonymous data set (a set containing no information as to class identity) into a fixed number (k) of clusters.
- Initially a number of so called centroids are chosen. A centroid is a data point (imaginary or real) at the center of a cluster. In Praat each centroid is an existing data point in the given input data set, picked at random, such that all centroids are unique (that is, for all centroids ci and cj, ci ≠ cj). These centroids are used to train a kNN classifier.
- The resulting classifier is used to classify (using k = 1) the data and thereby produce an initial randomized set of clusters. Each centroid is thereafter set to the arithmetic mean of the cluster it defines. The process of classification and centroid adjustment is repeated until the values of the centroids stabilize. The final centroids will be used to produce the final classification/clustering of the input data, effectively turning the set of initially anonymous data points into a set of data points, each with a class identity.



### 1.2 Evaluation Methods

- Contrary to supervised learning where we have the ground truth to evaluate the model's performance, clustering analysis doesn't have a solid evaluation metric that we use to evaluate the outcome of different clustering algorithms. Moreover, since kmeans requires k as an input and doesn't learn it from data, there is no right answer in terms of the number of clusters that we should have in any problem. Sometimes domain knowledge and intuition may help but usually that is not the case. In the cluster-predict methodology, we can evaluate how well the models are performing based on different K clusters since clusters are used in the downstream modeling. In this post we'll cover two metrics that may give us some intuition about k:
  - Elbow method**
  - Silhouette analysis**
- Elbow Method**
  - Elbow method gives us an idea on what a good k number of clusters would be based on the sum of squared distance (SSE) between data points and their assigned clusters' centroids. We pick k at the spot where SSE starts to flatten out and forming an elbow. We'll use the geyser dataset and evaluate SSE for different values of k and see where the curve might form an elbow and flatten out.
- Silhouette Analysis**
  - Silhouette analysis can be used to determine the degree of separation between clusters. For each sample:
    - Compute the average distance from all data points in the same cluster (ai).
    - Compute the average distance from all data points in the closest cluster (bi).

Compute the coefficient:

$$\frac{b^i - a^i}{\max(a^i, b^i)}$$

- The coefficient can take values in the interval [-1, 1].
  - If it is 0 -> the sample is very close to the neighboring clusters.
  - If it is 1 -> the sample is far away from the neighboring clusters.
  - If it is -1 -> the sample is assigned to the wrong clusters.
- From the above we can say that when the value of score is below the 0 that means the data is not clustered well.
- When we want the score value nearer to the 1 that means the cluster has been done in good manner.
- Therefore, we get the coefficients to be as big as possible and close to 1 to have a good clusters. We'll use here geyser dataset again because its cheaper to run the silhouette analysis and it is actually obvious that there is most likely only two groups of data points.

Below are the main takeaways:

- Scale/standardize the data when applying kmeans algorithm.
- Elbow method in selecting number of clusters doesn't usually work because the error function is monotonically decreasing for all ks.
- Kmeans gives more weight to the bigger clusters.
- Kmeans assumes spherical shapes of clusters (with radius equal to the distance between the centroid and the furthest data point) and doesn't work well when clusters are in different shapes such as elliptical clusters.
- If there is overlapping between clusters, kmeans doesn't have an intrinsic measure for uncertainty
- For the examples belong to the overlapping region in order to determine for which cluster to assign each data point Kmeans may still cluster the data even if it can't be clustered such as data that comes from uniform distributions.

### Let's do one example:-

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import silhouette_score

In [2]: data = pd.read_csv('faithful.csv', usecols=[1,2])

Out [3]: data.head()
```

	eruptions	waiting
0	3.600	79
1	1.800	54
2	3.333	74
3	2.283	62
4	4.533	85

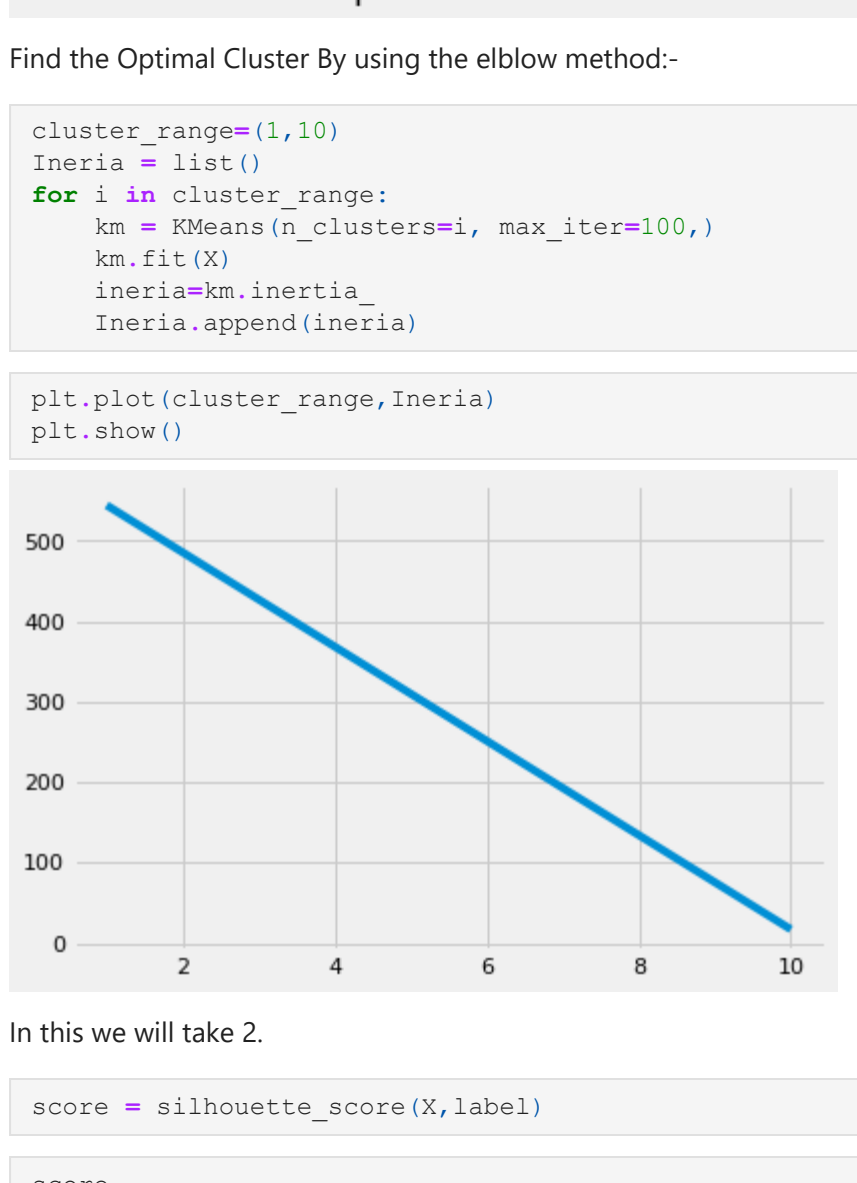
```
In [11]: data.tail()
```

```
Out [11]: eruptions  waiting
267    4.117      81
268    2.150      46
269    4.417      90
270    1.867      46
271    4.467      74
```

```
In [12]: data.shape
Out [12]: (272, 2)
```

```
In [13]: sns.set_context('notebook')
plt.style.use('fivethirtyeight')
```

```
In [14]: # Plot the data
plt.figure(figsize=(6, 6))
plt.scatter(data.iloc[:, 0], data.iloc[:, 1])
plt.xlabel('Eruption time in mins')
plt.ylabel('Waiting time to next eruption')
plt.title('Visualization of raw data');
```



```
In [25]: # Standardize the data
X = StandardScaler().fit_transform(data)

# Run local implementation of kmeans
km = KMeans(n_clusters=2, max_iter=100,
            random_state=None, tol=0.0001, verbose=0)
```

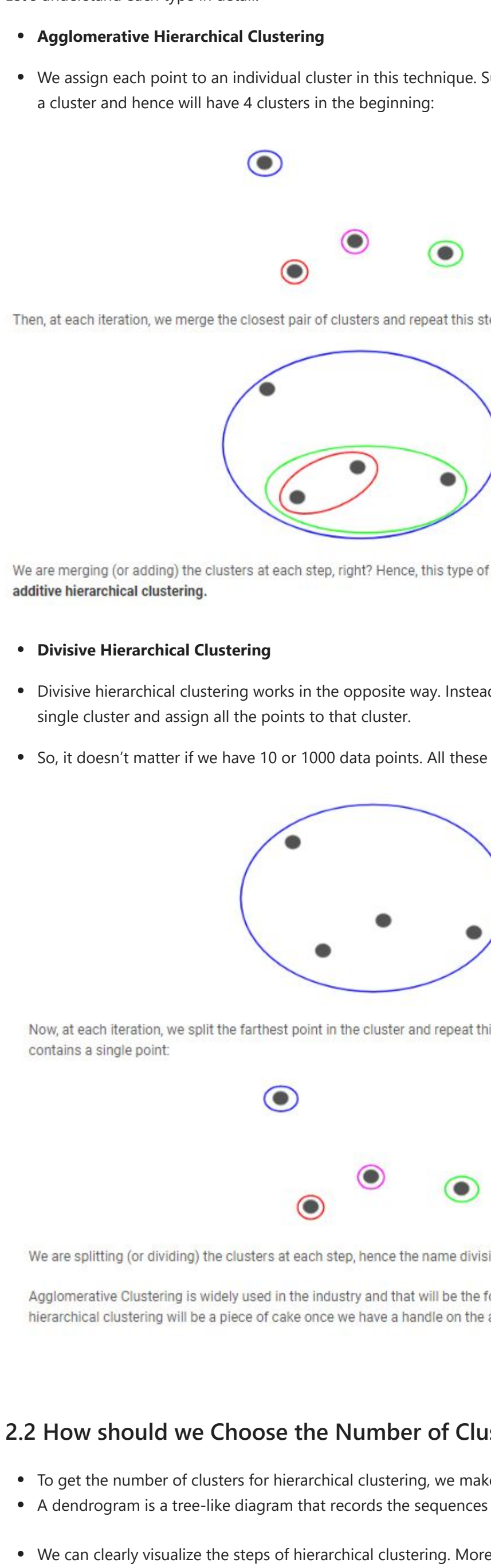
```
Out [25]: KMeans(algorithm='auto', copy_x=True, init='k-means++', max_iter=100,
               n_clusters=2, n_init=10, n_jobs=None, precompute_distances='auto',
               random_state=None, tol=0.0001, verbose=0)
```

```
In [27]: label=km.labels_

In [61]: Centroids=km.cluster_centers_
Centroids=list(Centroids)
```

```
Out [61]: [[array([-1.26008539, -1.20156744]), array([0.70970327, 0.67674488])]]
```

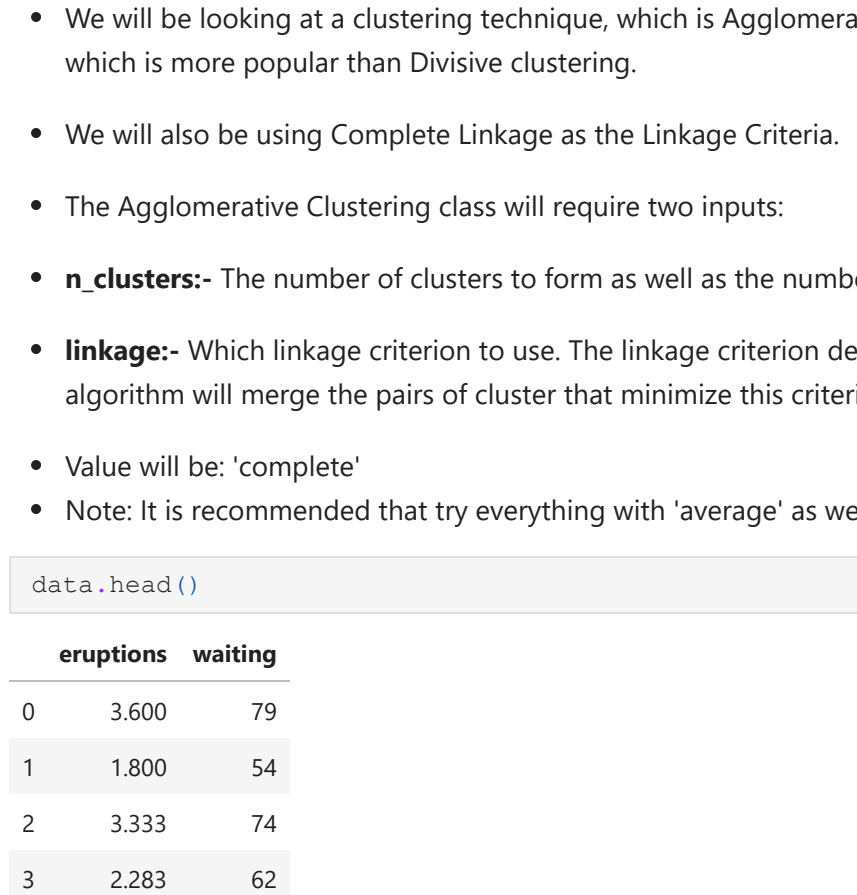
```
In [68]: # Plot the clustered data
fig, ax = plt.subplots(figsize=(6, 6))
plt.scatter(X[label == 0, 0], X[label == 0, 1],
            c='green', label='cluster 1')
plt.scatter(X[label == 1, 0], X[label == 1, 1],
            c='blue', label='cluster 2')
plt.scatter(Centroids[0][0], Centroids[0][1], marker='*', s=300,
            c='r', label='centroid1')
plt.scatter(Centroids[1][0], Centroids[1][1], marker='*', s=300,
            c='r', label='centroid2')
```



Find the Optimal Cluster By using the elbow method:-

```
In [86]: cluster_range=(1,10)
Inertia = list()
for i in cluster_range:
    km = KMeans(n_clusters=i, max_iter=100,)
    km.fit(X)
    inertia=km.inertia_
    Inertia.append(inertia)
```

```
In [88]: plt.plot(cluster_range, Inertia)
plt.show()
```



In this we will take 2.

```
In [89]: score = silhouette_score(X, label)
```

```
In [90]: score
Out [90]: 0.7451774401207985
```

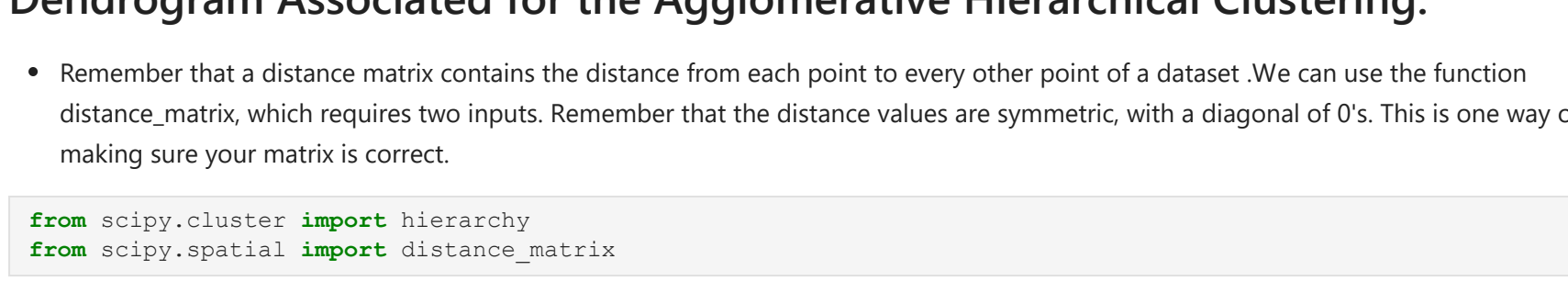
This is indicating that the clustering has been done properly.

## 2. Heirarchical Clustering.

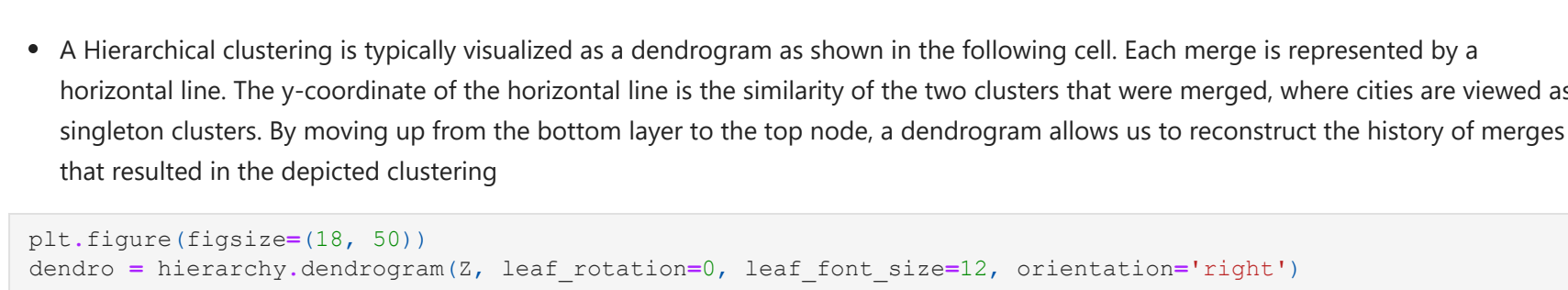
- In the KMeans clustering, we provide first number of cluster that we required in order to get the Cluster. But there are certain challenges with K-means. It always tries to make clusters of the same size. Also, we have to decide the number of clusters at the beginning of the algorithm. Ideally, we would not know how many clusters should we have, in the beginning of the algorithm and hence it's a challenge with K-means.
- This is a gap hierarchical clustering bridges with aplobm. It takes away the problem of having to pre-define the number of clusters. Sounds like a dream! So, let's see what hierarchical clustering is and how it improves on K-means.

### 2.1 What is Heirarchical Clustering:-

Let's say we have the below points and we want to cluster them into groups:



We can assign each of these three points to a separate cluster:



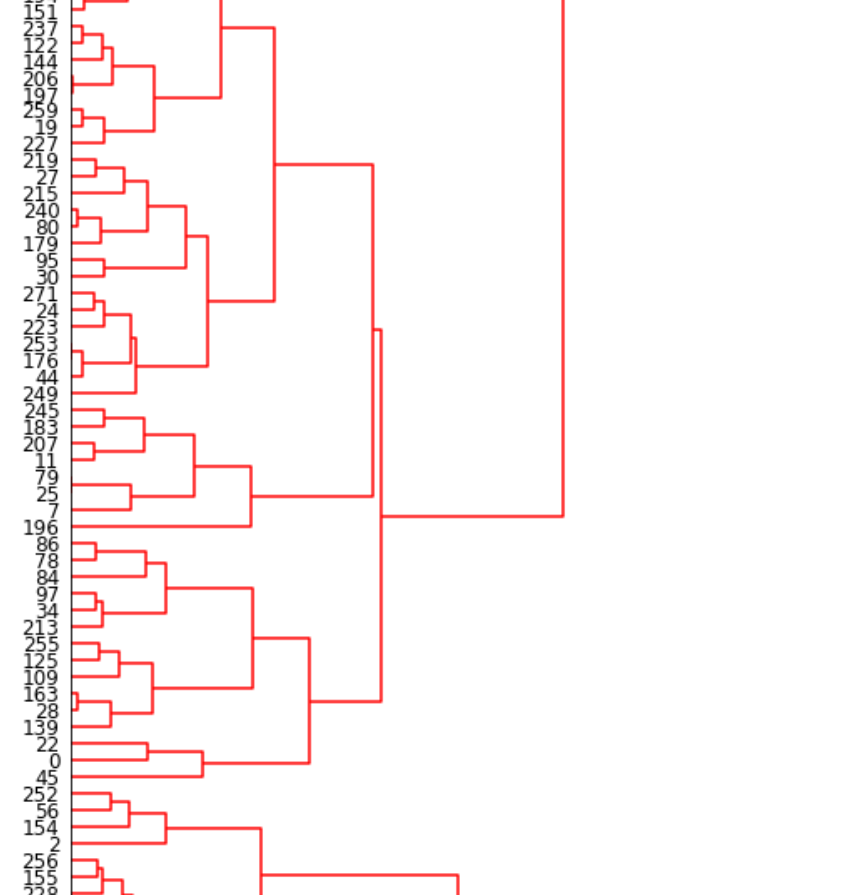
Now, based on the similarity of these clusters, we can combine the most similar clusters together and repeat this process until only a single cluster is left:



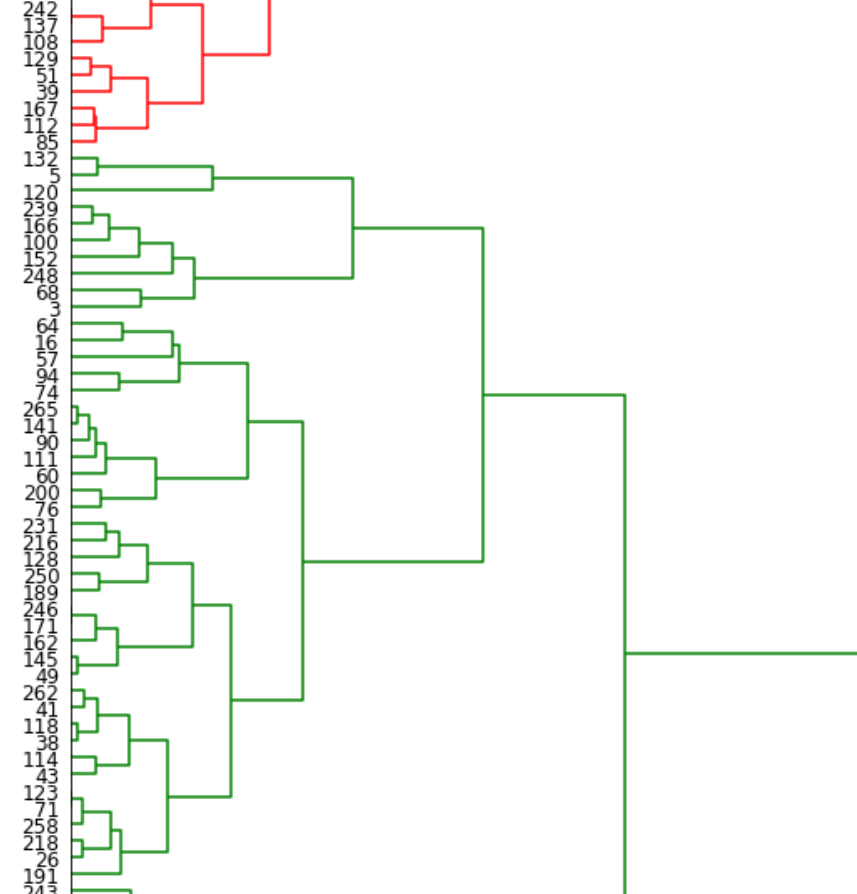
- We are essentially building a hierarchy of clusters. That's why this algorithm is called hierarchical clustering. For now, let's look at the different types of hierarchical clustering.
- Types of Hierarchical Clustering There are mainly two types of hierarchical clustering:
  - Agglomerative hierarchical clustering**
  - Divisive Hierarchical clustering**

Let's understand each type in detail.

- Agglomerative Hierarchical Clustering**
  - We assign each point to an individual cluster in this technique. Suppose there are 4 data points. We will assign each of these points to a cluster and hence will have 4 clusters in the beginning:

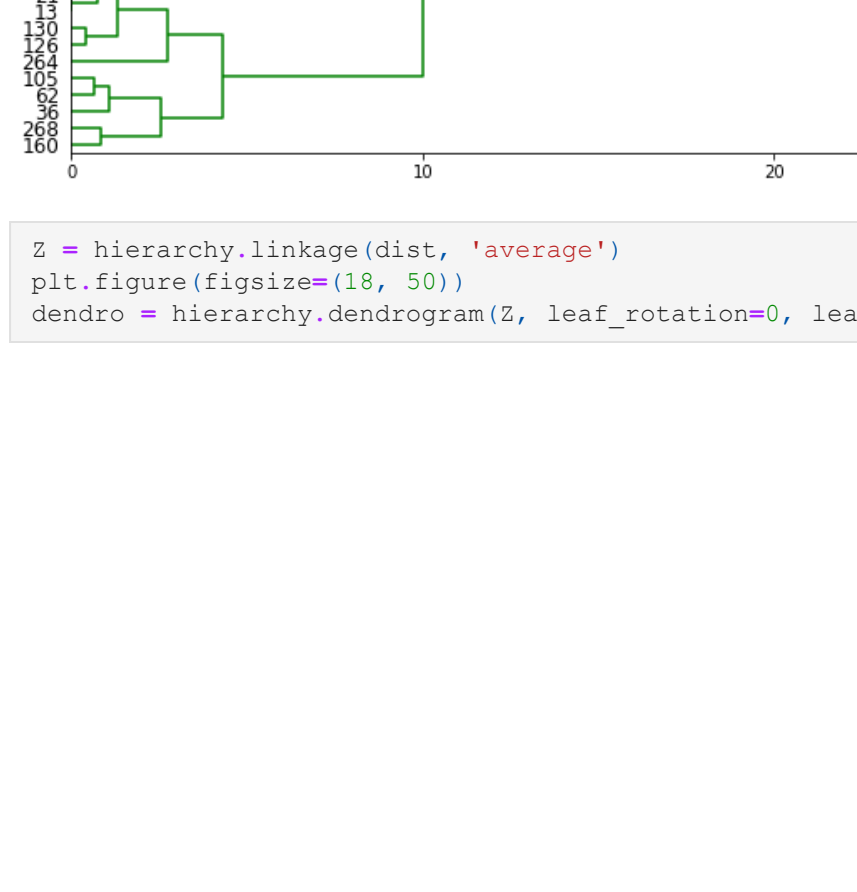


Then, at each iteration, we merge the closest pair of clusters and repeat this step until only a single cluster is left:



We are merging (or adding) the clusters at each step, right? Hence, this type of clustering is also known as **additive hierarchical clustering**.

- Divisive Hierarchical Clustering**
  - Divisive hierarchical clustering works in the opposite way. Instead of starting with n clusters (in case of n observations), we start with a single cluster and assign all the points to that cluster.
  - So, it doesn't matter if we have 10 or 1000 data points. All these points will belong to the same cluster at the beginning:-



Now, at each iteration, we split the farthest point in the cluster and repeat this process until each cluster only contains a single point:



We are splitting (or dividing) the clusters at each step, hence the name **divisive hierarchical clustering**.

Agglomerative Clustering is widely used in the industry and that will be the focus in this article. Divisive hierarchical clustering will be a piece of cake once we have a handle on the agglomerative type.

### 2.2 How should we Choose the Number of Clusters in Hierarchical Clustering?

- To get the number of clusters for hierarchical clustering, we make use of an awesome concept called a Dendrogram.
- A dendrogram is a tree-like diagram that records the sequences of merges or splits.
- We can clearly visualize the steps of hierarchical clustering. More the distance of the vertical lines in the dendrogram, more the distance between those clusters.
- Now, we can set a threshold distance and draw a horizontal line (Generally, we try to set the threshold in such a way that it cuts the tallest vertical line). Let's set this threshold as 12 and draw a horizontal line:



The number of clusters will be the number of vertical lines which are being intersected by the line drawn using the threshold. In the above example, since the red line intersects 2 vertical lines, we will have 2 clusters. One cluster will have a sample (1,2,4) and the other will have a sample (3,5). Pretty straightforward, right?

This is how we can decide the number of clusters using a dendrogram in Hierarchical Clustering.

### Agglomerative Clustering Method.

- We will be looking at a clustering technique, which is Agglomerative Hierarchical Clustering. Agglomerative is the bottom up approach which is more popular than Divisive clustering.
- We will also be using Complete Linkage as the Linkage Criteria.
- The Agglomerative Clustering class will require two inputs:
  - n\_clusters**: The number of clusters to form as well as the number of centroids to generate.
  - linkage**: Which linkage criterion to use. The linkage criterion determines which distance to use between sets of observation. The algorithm will merge the pairs of cluster that minimize this criterion.
  - Value** will be: 'complete'
  - Note**: It is recommended that try everything with 'average' as well

```
In [5]: data.head()
```

```
Out [5]: eruptions  waiting
0    3.600      79
1    1.800      54
2    3.333      74
3    2.283      62
4    4.533      85
```

```
In [10]: scaler = StandardScaler()
X = scaler.fit_transform(data)
```

```
In [6]: from sklearn.cluster import AgglomerativeClustering
```

```
In [7]: aggllo = AgglomerativeClustering(n_clusters=2, linkage='complete')
```

```
In [11]: aggllo.fit(X)
```

```
Out [11]: AgglomerativeClustering(affinity='euclidean', compute_full_tree='auto',
                                connectivity=None, distance_threshold=None,
                                linkage='complete', memory=None, n_clusters=2)
```

```
In [12]: label=aggllo.labels_

In [21]: sns.scatterplot('eruptions', 'waiting', hue=label, data=data)
plt.title('Agglomerative with 2 Clusters')
plt.grid()
plt.show()
```



### Dendrogram Associated for the Agglomerative Hierarchical Clustering.

- Remember that a distance matrix contains the distance from each point to every other point of a dataset. We can use the function distance\_matrix, which requires two inputs. Remember that the distance values are symmetric, with a diagonal of 0's. This is one way of making sure your matrix is correct.

```
In [22]: from scipy.cluster import hierarchy
from scipy.spatial import distance_matrix
```

```
In [25]: dist = distance_matrix(X, X)
```

- Using the linkage class from hierarchy, pass in the parameters:
  - The distance matrix
  - 'complete' for complete linkage

```
In [26]: Z = hierarchy.linkage(dist, 'complete')
```

- A Hierarchical clustering is typically visualized as a dendrogram as shown in the following cell. Each merge is represented by a horizontal line. The y-coordinate of the horizontal line is the similarity of the two clusters that were merged, where cities are viewed as singleton clusters. By moving up from the bottom layer to the top node, a dendrogram allows us to reconstruct the history of merges that resulted in the depicted clustering

```
In [28]: plt.figure(figsize=(18, 50))
dendro = hierarchy.dendrogram(Z, leaf_rotation=0, leaf_font_size=12, orientation='right')
```



```
In [33]: Z = hierarchy.linkage(dist, 'average')
plt.figure(figsize=(18, 50))
dendro = hierarchy.dendrogram(Z, leaf_rotation=0, leaf_font_size=12, orientation='right')
```





