

Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Hierarchical Clustering

Led by : Shreyas Shukla

Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Hierarchical clustering is very common in biology and lends itself nicely to visualizing clusters.

It can also help the user decide on an appropriate number of clusters.

Led by : Shreyas Shukla

Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Overview

- 1. Theory and Intuition*
- 2. Coding*

Led by : Shreyas Shukla

Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

iHUB DivyaSampark, IIT Roorkee
and
Ritvi Bharat Private Limited

Theory and Intuition

Led by : Shreyas Shukla

Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Like most clustering algorithms, Hierarchical Clustering simply relies on measuring which data points are most “similar” to other data points.

“Similarity” is defined by choosing a distance metric.

Led by : Shreyas Shukla

Benefits of Hierarchical Clustering

- Easy to understand and visualize.
- Helps users decide how many clusters to choose.
- Not necessary to choose cluster amount **before** running the algorithm.

Led by : Shreyas Shukla

Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

So why use Hierarchical Clustering?

- Divides points into ***potential*** clusters:

Led by : Shreyas Shukla

So why use Hierarchical Clustering?

- Divides points into ***potential*** clusters:
 - Agglomerative Approach:
 - Each point begins as its own cluster, then clusters are joined.
 - Divisive Approach:
 - All points begin in the same cluster, then clusters are split.

Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Agglomerative:



Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Agglomerative:

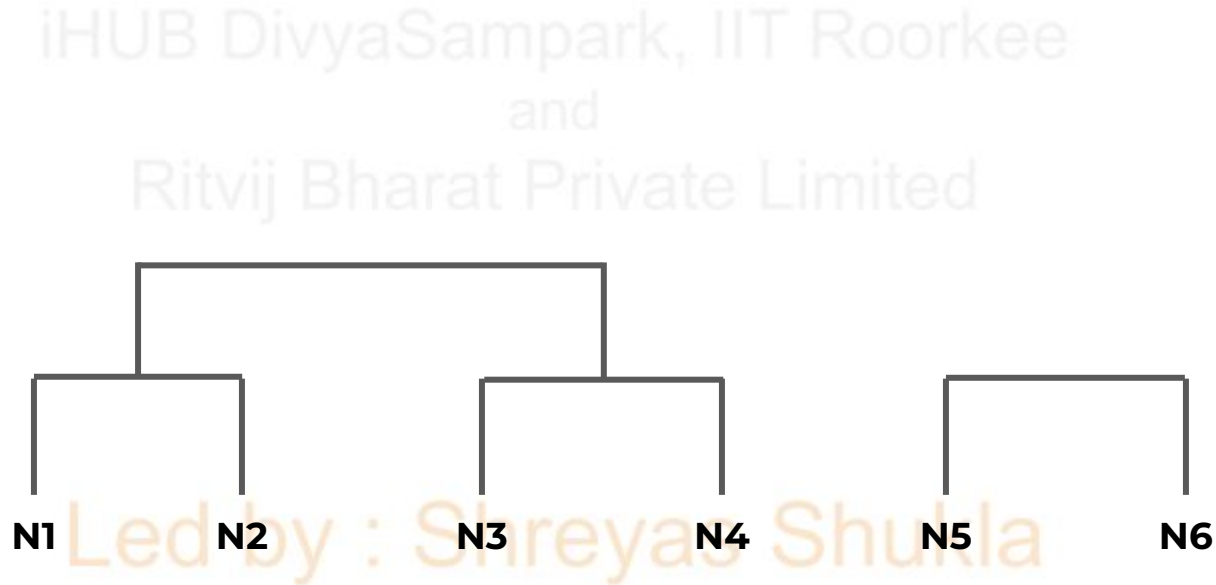
iHUB DivyaSampark, IIT Roorkee
and
Ritvij Bharat Private Limited



Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

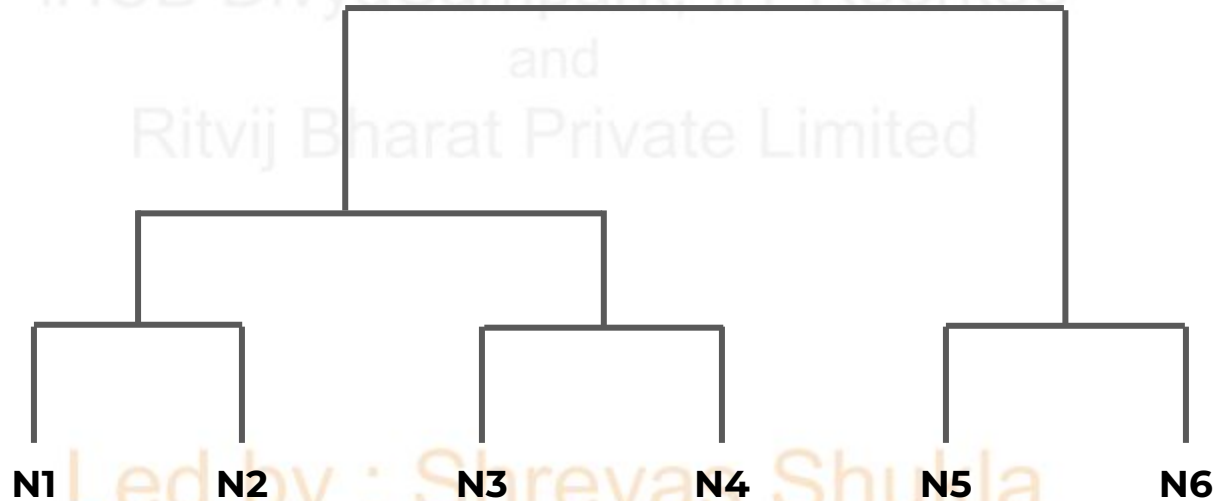
Agglomerative:



Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Agglomerative:



Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Opposite of the Agglomerative approach is a **Divisive** approach, which starts with all points belonging to the same cluster, and then begins divisions to separate out clusters.

Led by : Shreyas Shukla

Hierarchical Clustering Process

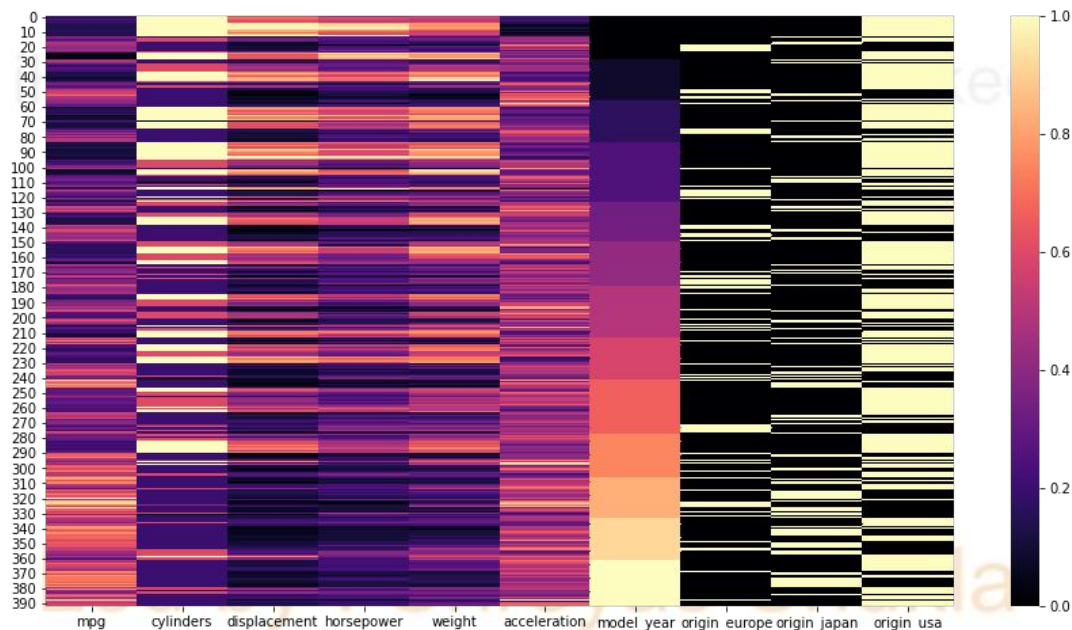
- Compare data points to find most similar data points to each other.
- Merge these to create a cluster.
- Compare clusters to find most similar clusters and merge again.
- Repeat until all points in a single cluster.

Led by : Shreyas Shukla

Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

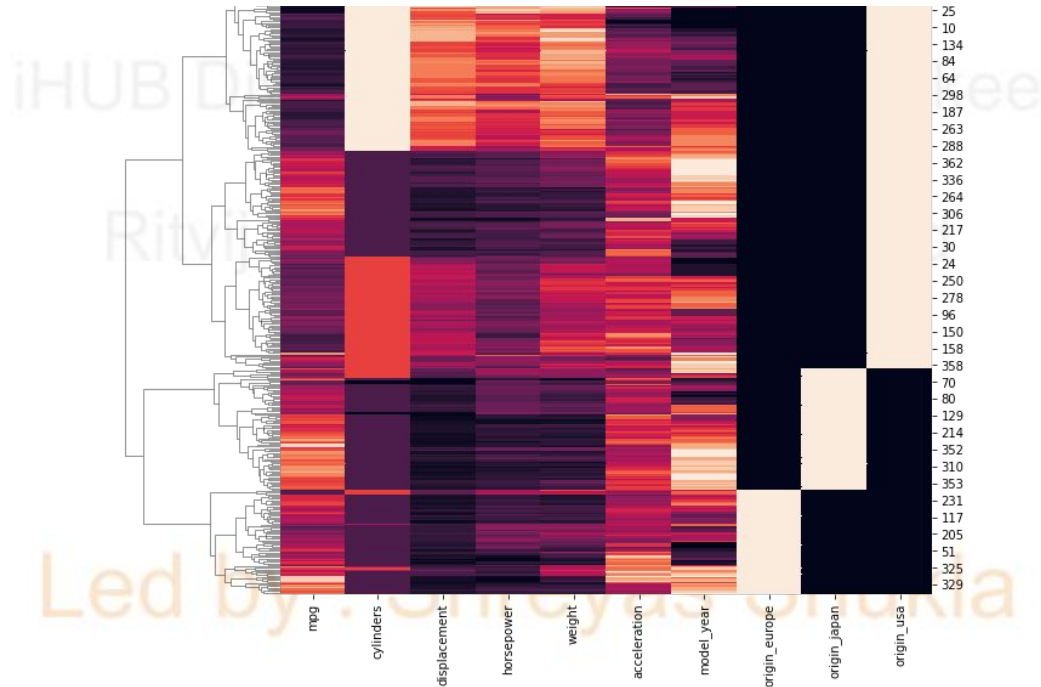
Hierarchical Clustering Process



Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Hierarchical Clustering Process



Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Topics which we still need to understand for Hierarchical Clustering:

- Similarity Metric
- Dendrogram
- Linkage Matrix

Led by : Shreyas Shukla

Similarity Metric

Measures distance between two points.

Many types:

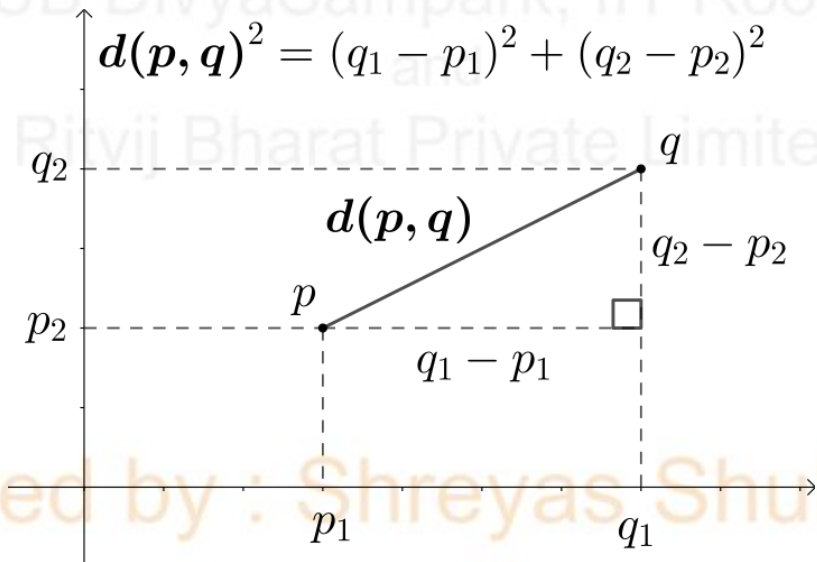
- Euclidean Distance
- Manhattan
- Cosine
- and many more...

Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Similarity Metric

Default choice is Euclidean



Similarity Metric

- Each dimension would be a feature
- For **n** data points and **p** features:

- $$D^2 = (x_{11} - x_{12})^2 + \dots + (x_{n-1p-1} - x_{np})^2$$

Led by : Shreyas Shukla

Similarity Metric

- Each dimension would be a feature
- For **n** data points and **p** features:
 - $D^2 = (x_{11} - x_{12})^2 + \dots + (x_{n-1p-1} - x_{np})^2$
- Using MinMaxScaler we can scale all features to be between 0 and 1.
- This allows for maximum distance between a feature to be 1.

Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Dendrogram:

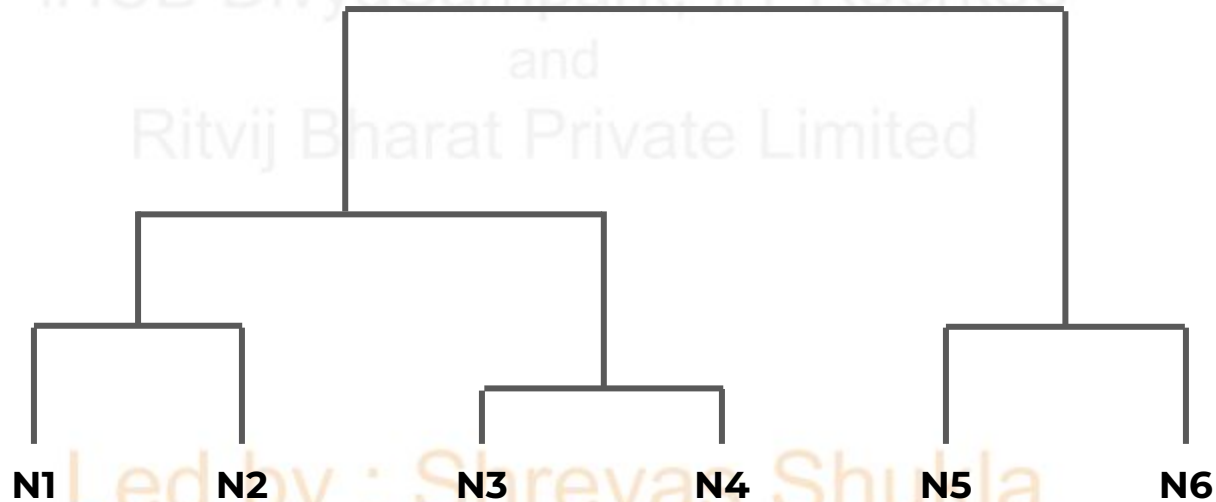
- Plot displaying all potential clusters.
- Very computationally expensive to compute and display for larger data sets.
- Very useful for deciding on number of clusters.

Led by : Shreyas Shukla

Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

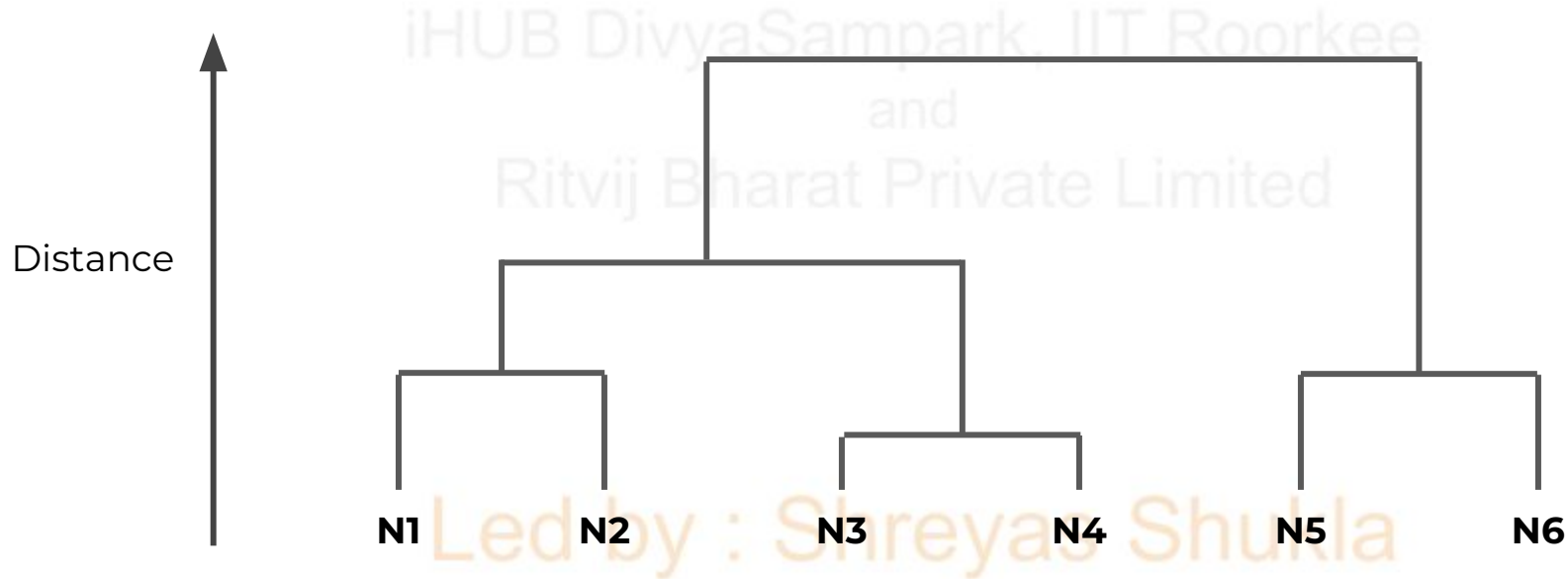
Dendrogram:



Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

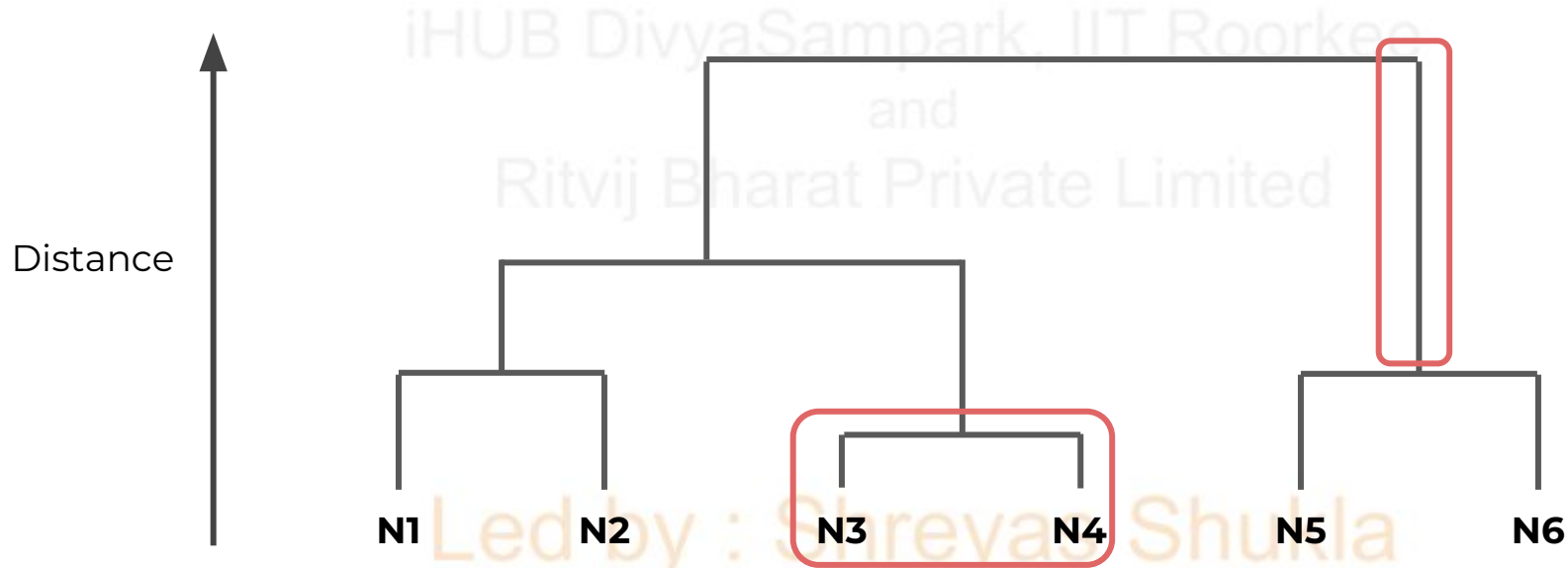
Dendrogram:



Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Dendrogram:

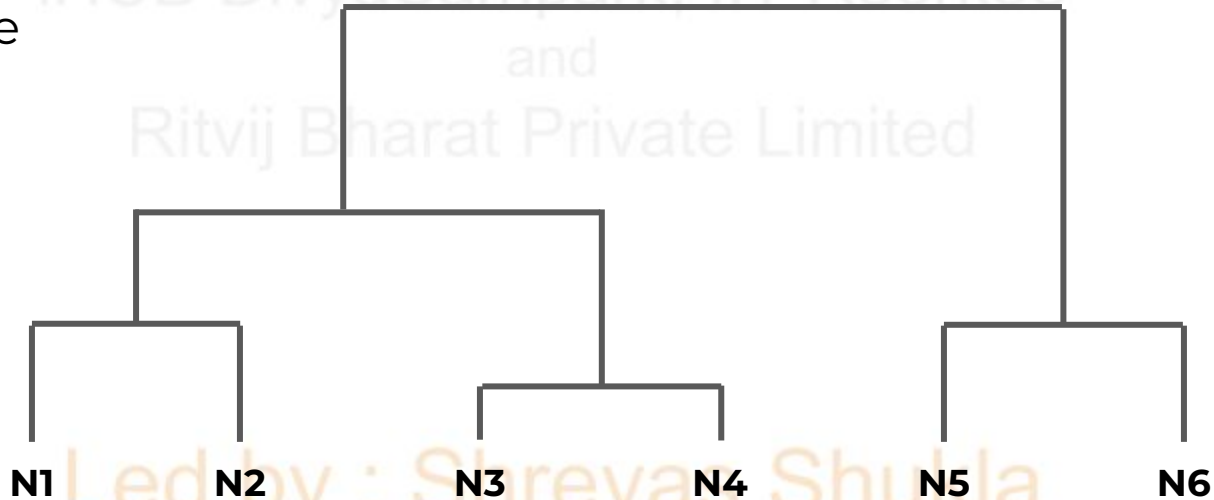


Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Dendrogram:

“Slice” to decide
cluster count

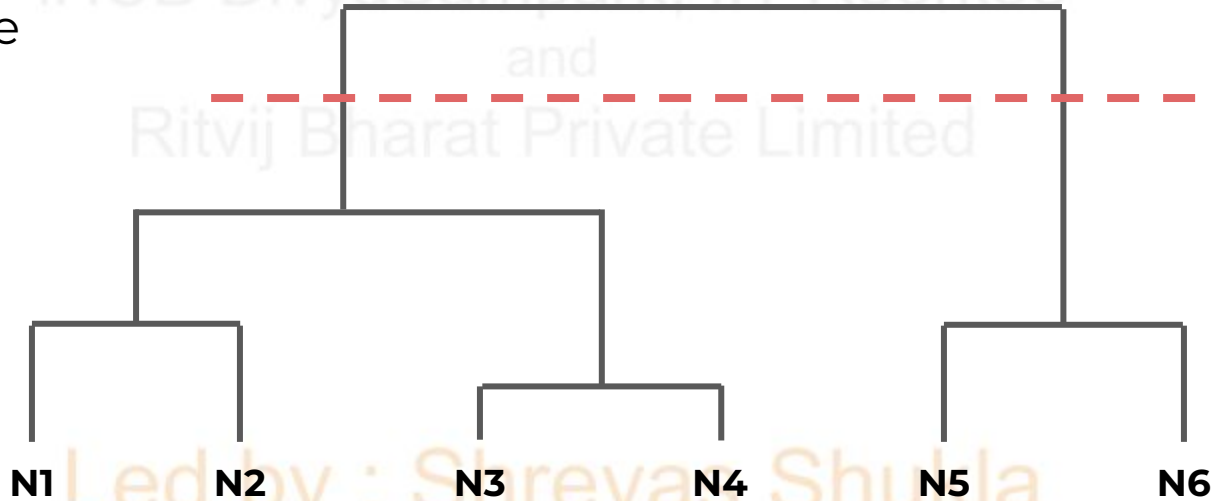


Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Dendrogram:

“Slice” to decide
cluster count

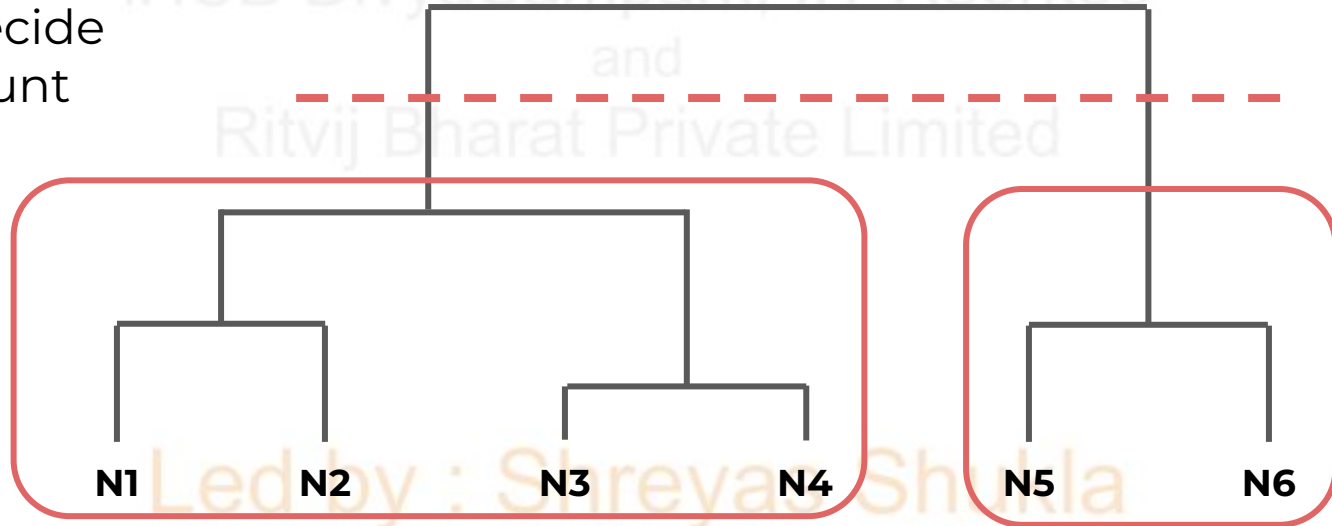


Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Dendrogram:

“Slice” to decide
cluster count

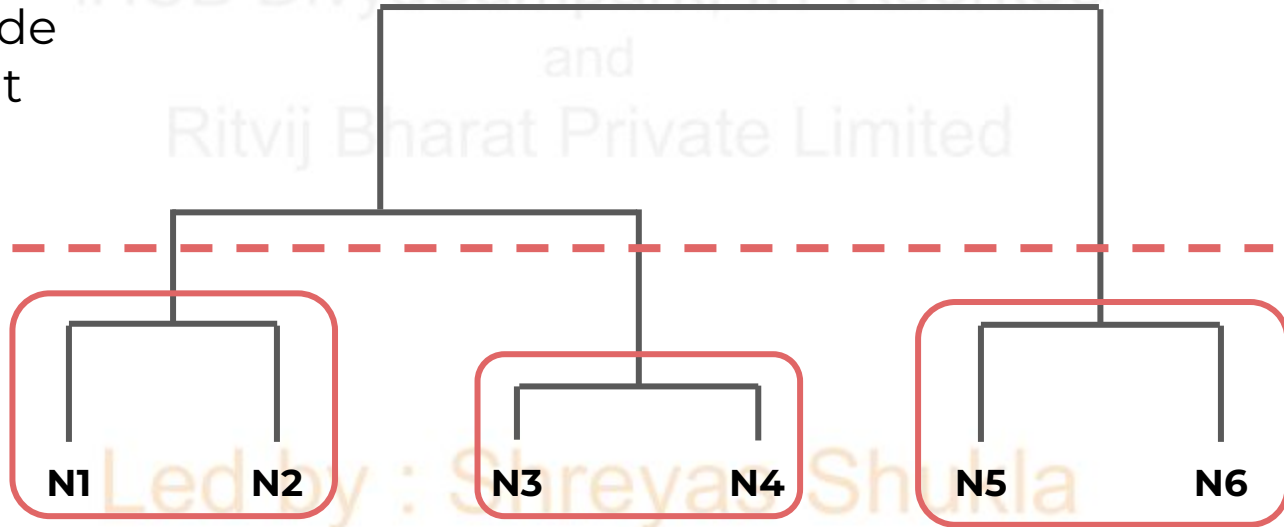


Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Dendrogram:

“Slice” to decide
cluster count



Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Linkage

- How do we measure distance from a point to an entire cluster?
- How do we measure distance from a cluster to another cluster?

Led by : Shreyas Shukla

Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

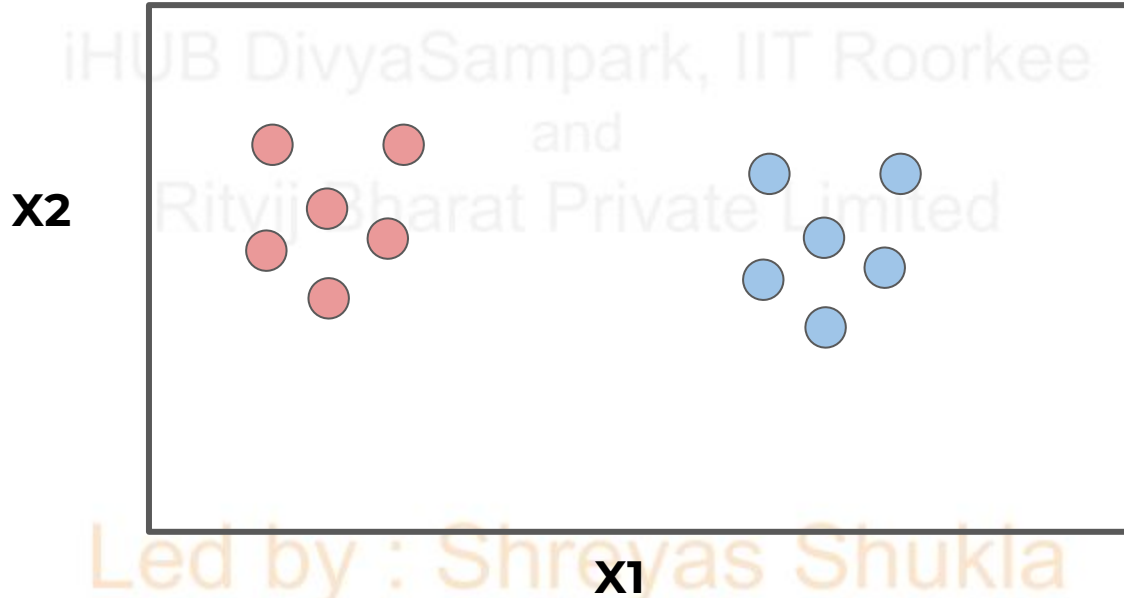
Linkage

Once two or more points are together and we want to continue agglomerative clustering to join clusters, we need to decide on a **linkage** parameter.

Led by : Shreyas Shukla

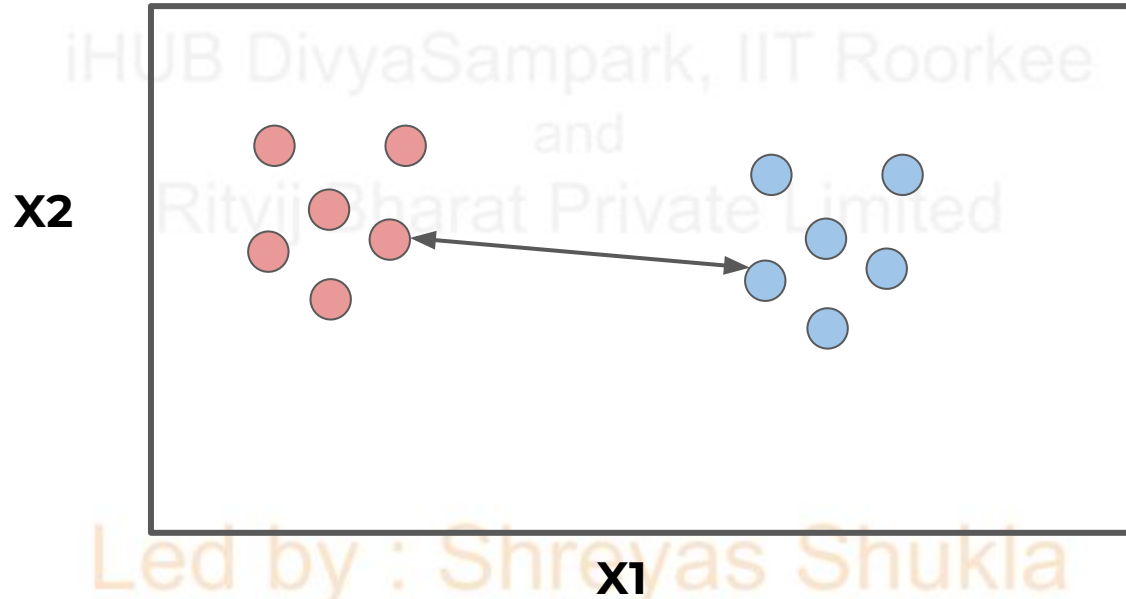
Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)



Mastering Machine Learning with Python

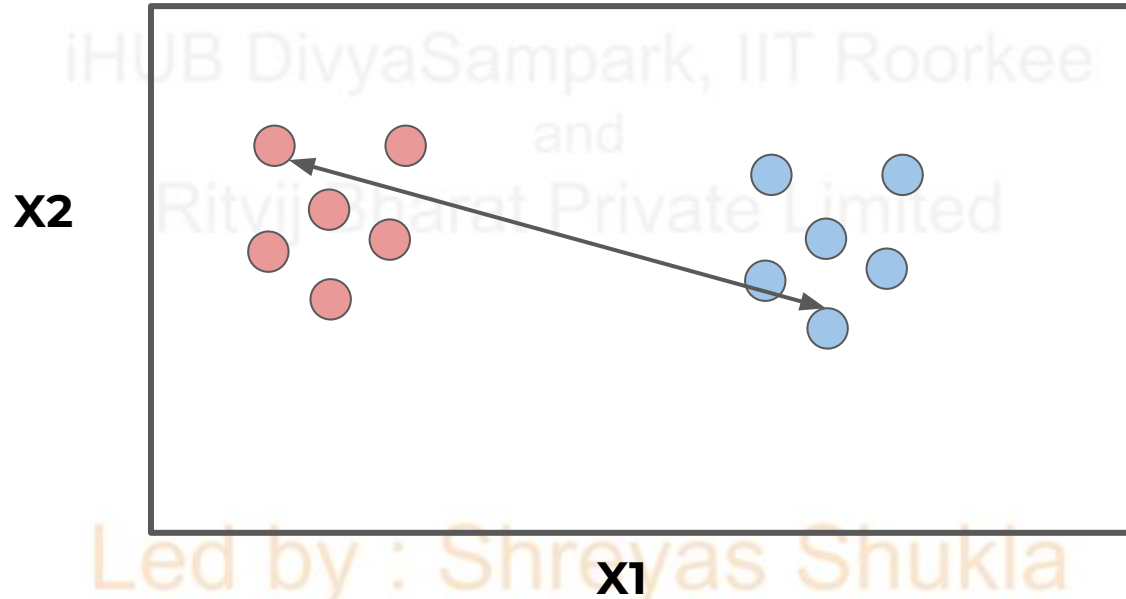
(27th Aug 2024 - 18th Oct 2024)



Led by : Shreyas Shukla

Mastering Machine Learning with Python

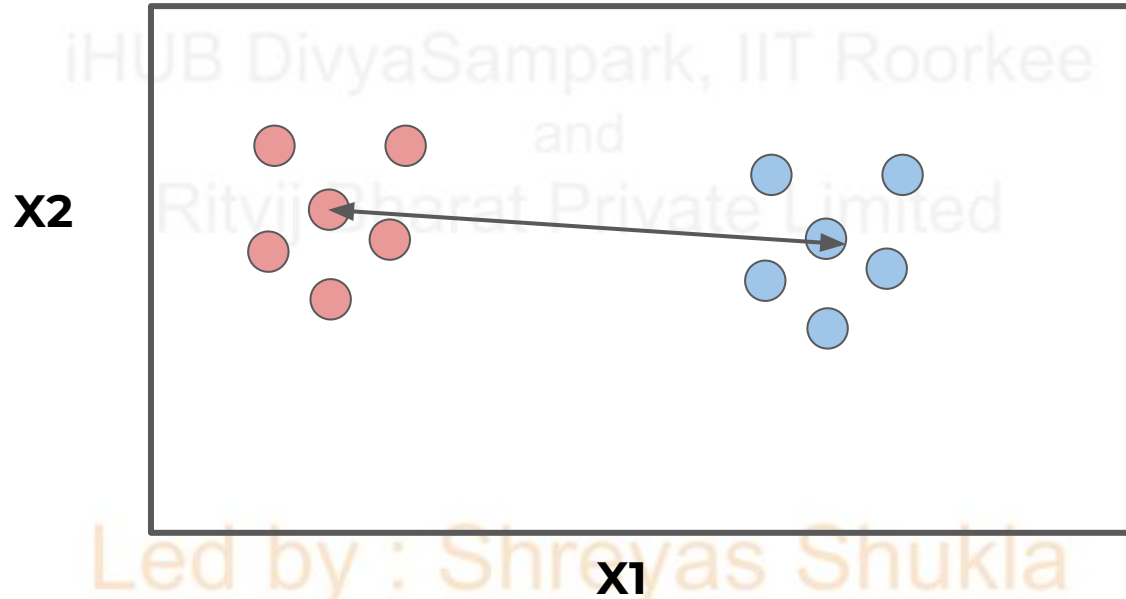
(27th Aug 2024 - 18th Oct 2024)



Led by : Shreyas Shukla

Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

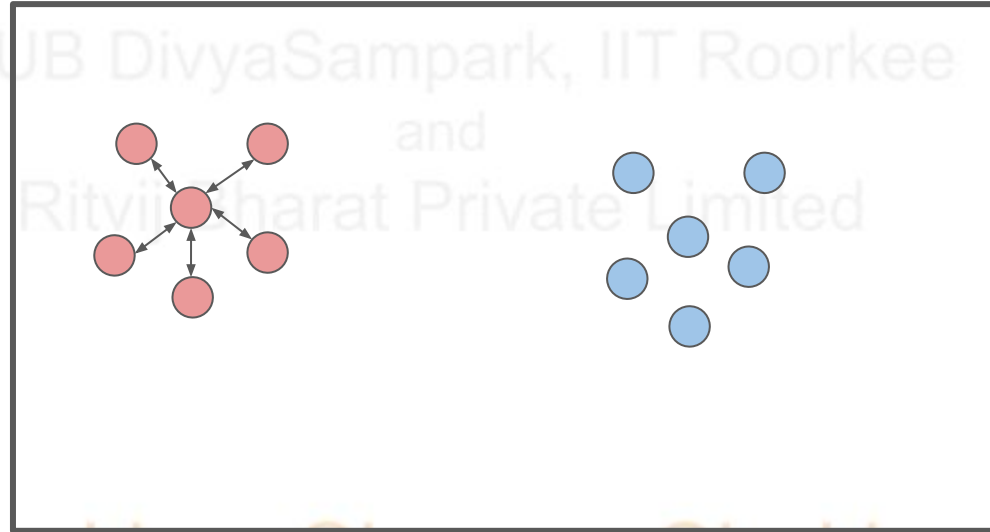


Led by : Shreyas Shukla

Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

x2



x1

Led by : Shreyas Shukla

Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Linkage

- Criterion determining which distance to use between sets of observation.
- Algorithm will merge pairs of clusters that minimizes the criterion.

Led by : Shreyas Shukla

Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Linkage:

- **Ward:** minimizes variance of clusters being merged.
- **Average:** uses average distances between two sets.
- **Minimum** or **Maximum** distances between all observations of the two sets.

led by : Shreyas Shukla

Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

iHUB DivyaSampark, IIT Roorkee

Let's code!!

Ritvij Bharat Private Limited

Led by : Shreyas Shukla