Hierarchical Clustering Agenda In this session, we will cover the following concepts with the help of a business use case: Hierarchical Clustering Types of Hierarchical Clustering Evaluate similarities between clusters **Hierarchical Clustering** Hierarchical clustering is where we build a cluster tree (dendrogram) to represent data, where each group (node) is linked to two or more successor groups. • The groups are nested and organized as a tree, which ideally ends up as a meaningful classification scheme. Each node in the cluster tree contains a group of similar data; nodes are placed on the graph next to other similar nodes. Clusters at one level are joined with clusters in the next level up, using a degree of similarity. • The process carries on until all nodes are in the tree, which gives a visual snapshot of the data contained in the whole set. • The total number of clusters is not predetermined before you start the tree creation.

There is also evidence that divisive algorithms produce more accurate hierarchies than bottom-up algorithms. In some circumstances according to the [Stanford University Review] (https://nlp.stanford.edu/IR-book/html/htmledition/references-and-further-reading-

In both techniques, we discussed finding the similarity or dissmilarity between the two clusters. The question is, "How to measure them?"

- 2. Divisive or Top-Down Clustering

1. Agglomerative or Bottom-Up Clustering

Divisive (Top-Down) Clustering

1. Start with all examples in the same cluster

Agglomerative (Bottom-Up) Clustering

There are two major ways in which hierarchical clustering can be carried out:

Implementing Hierarchical Clustering

1. Start with each example in its own singleton cluster 2. At each time-step, greedily merge two most similar clusters

3. Stop when there is a single cluster of all examples, else go to two

3. Stop when each example is in its own singleton cluster, else go to two Which One to Use?

2. At each time-step, remove the "outsiders" from the least cohesive cluster

- Agglomerative Clustering is simpler than the divisive clustering to implement because for the latter, we need a second, flat clustering algorithm as a subroutine.
 - However, top-down routine has the advantage of being more efficient if we do not generate a complete hierarchy all the way down to individual document leaves. ullet For a fixed number of top levels, using an efficient flat algorithm like K-means. Top-down algorithms are linear in the number of documents and clusters.

Deciding (Dis)similarity between Clusters

how the distance between each cluster is measured.

Draw a horizontal line at both extremities

Before any clustering is performed, it is required that you determine the proximity matrix containing the distance between each point using a distance function. Then, the matrix is updated to display the distance between each cluster. The following three methods differ in

Dendrograms

0.65

0.6

0.55

0.5

0.45

0.4

0.25

Linkage Criteria

Single Linkage

points.

17.html#sec:hclstfurther).

For example, in the below case, best choice for no. of clusters will be 4.

• Determine the largest vertical distance that doesn't intersect any of the other clusters

We can use a dendrogram to visualize the history of groupings and figure out the optimal number of clusters.

The optimal number of clusters is equal to the number of vertical lines going through the horizontal line

- 0.35 0.3
- 0.15 9 23 17 6 11 3 15 7 14 19 16 24 22 1 13 12 5 21 4 10 20 2 18 8 25

In single linkage hierarchical clustering, the distance between two clusters is defined as the shortest distance between two points in each cluster. For example, the distance between clusters "r" and "s" to the left is equal to the length of the arrow between their two closest

In complete linkage hierarchical clustering, the distance between two clusters is defined as the longest distance between two points in each cluster. For example, the distance between clusters "r" and "s" to the left is equal to the length of the arrow between their two furthest

In average linkage hierarchical clustering, the distance between two clusters is defined as the average distance between each point in one cluster to every point in the other cluster. For example, the distance between clusters "r" and "s" to the left is equal to the average length of

$L(r,s) = \min(D(x_{ri}, x_{si}))$ **Complete Linkage**

points.

Average Linkage

 $L(r,s) = \max(D(x_{ri}, x_{si}))$

each arrow between connecting the points of one cluster to the other.

Ward Linkage

In ward linkage hierarchical clustering, the distance between clusters is the sum of squared differences within all clusters.

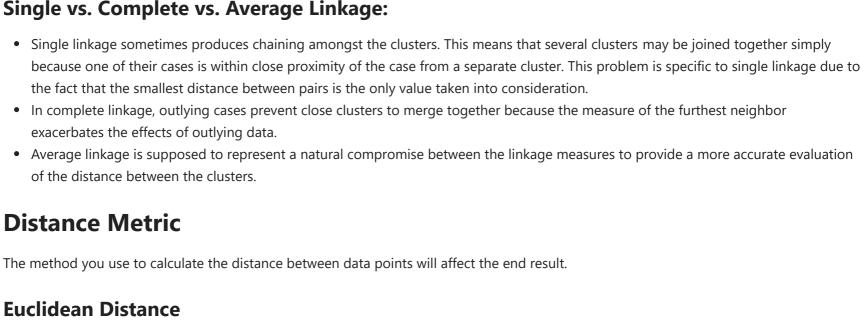
Imagine you were in the downtown center of a big city and you wanted to get from point A to point B. You wouldn't be able to cut across

One of the biggest drawbacks of hierarchical clustering is it is extremely calculation heavy. Therefore, they are not scalable. That also means

Scipy has a really convenient api for carrying out hierarchical clustering. Let's see how it works. We will start with necessary imports:

In the above code, we are using the sklearn library, which contains a lot of tools for machine learning and statistical modeling, including

Load_iris is a function from sklearn that loads and returns the iris dataset. The libraries used above already contains it. Just by loading the



The shortest distance between two points.

the Euclidean distance between x and y is $\sqrt{(a-c)^2+(b-d)^2}$

the Manhattan distance between x and y is |a-c|+|b-d|

that they are not very useful for larger datasets.

import pandas as pd, numpy as np from sklearn.cluster import KMeans from sklearn.datasets import load iris

classification, regression, clustering, and dimensionality reduction.

library, a data frame named "iris" will be made available and can be used straight away.

#Import dendrogram and linkage module from scipy library from scipy.cluster.hierarchy import dendrogram, linkage

buildings, rather you'd have to make your way by walking along the various streets.

For example, if x=(a,b) and y=(c,d),

Manhattan Distance

For example, if x=(a,b) and y=(c,d),

Shortcomings

Import Libraries

In [1]:

In [2]:

In [5]:

In [6]:

 $L(r,s) = \frac{1}{n_r n_s} \sum_{i=1}^{n_r} \sum_{i=1}^{n_s} D(x_{ri}, x_{sj})$

In [3]: X = iris.data y = iris.target In [4]:

#Generate the linkage matrix Z = linkage(X, 'average')

#Calculate full dendrogram import matplotlib.pyplot as plt

plt.figure(figsize=(25, 10))

plt.xlabel('sample index') plt.ylabel('distance')

plt.title('Hierarchical Clustering Dendrogram')

leaf rotation=90., #Rotates the x axis labels

leaf font size=8., #Font size for the x axis labels

Now we are going to see the concept of Agglomerative (Bottom-Up) Clustering with a business use case.

Before reading the data from a .csv file, you need to download "housing_data.csv" dataset from the course resource and upload it into the

lab. We must use the Up arrow icon, which is shown in the left side under View icon. Click on the Up arrow icon and upload the file

Spending Score (1-100)

39

81

77

40

space, we will retain only two of these five columns. We can remove CustomerID column, Genre, and Age column. We will retain the Annual

Customer Dendograms

If we draw a horizontal line that passes through longest distance without a horizontal line, we get 5 clusters.

cluster = AgglomerativeClustering(n clusters=5, affinity='euclidean', linkage='ward')

array([4, 3, 4, 3,

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

<matplotlib.collections.PathCollection at 0x1eb265db190>

plt.scatter(data[:,0], data[:,1], c=cluster.labels , cmap='rainbow')

4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 1, 1, 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,

We create an instance of Agglomerative Clustering using the euclidean distance as the measure of distance between points and ward

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.

100

Note: In this lesson, we saw the use of the unservised learning methods, but in the next lesson we will be working on "Time-Series

When the shopping data is grouped using the agglomerative clustering technique, we can observe that there are five groups for

120

140

15

15

17

Segment customers into different groups based on their shopping trends.

iris = load iris()

dendrogram(

plt.show()

3.5

3.0

1.5

Problem Statement: An ecommerce company has prepared a rough dataset containing shopping details of their customers, which includes CustomerID, Genre, Age, Annual Income (k\$), Spending Score (1-100). The company is unable to target a specific set of customers with a particular set of SKUs. **Objective:**

Dataset

In [7]:

In [8]:

import matplotlib.pyplot as plt import pandas as pd %matplotlib inline import numpy as np

wherever it is downloaded into your system.

customer data = pd.read csv('shopping data.csv') In the above code, pd.read_csv function is used to read the "shopping_data.csv" file, and customer_data is a variable that will store the data read by the .csv file. In [9]: customer data.shape

CustomerID

Male

Male

4 Female

5 Female

19

21

23

31

(200, 5)

Out[9]: Here, shape is an attribute that returns a tuple representing the dimensionality of the customer_data. It is used to define the number of rows and columns in customer_data. In [10]: customer data.head()

Out[10]:

Our dataset has five columns: CustomerID, Genre, Age, Annual Income, and Spending Score. To view the results in two-dimensional feature

0

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 [11]:

In [12]: plt.title("Customer Dendograms") dend = shc.dendrogram(shc.linkage(data, method='ward')) 350

300

200

150

100

In [13]:

Out[13]:

In [14]:

Out[14]:

100

80

60

40

20

0

20

Conclusion

Modelling".

40

consumers whose labels range from 0 to 4.

Powered by simplilearn

60

linkage to calculate the proximity of clusters.

0, 2], dtype=int64)

plt.figure(figsize=(10, 7))

cluster.fit predict(data)

from sklearn.cluster import AgglomerativeClustering

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: import scipy.cluster.hierarchy as shc plt.figure(figsize=(25, 10))

Genre Age Annual Income (k\$)