

Male:1 & Female:0

To label the Gender text data to numerical values

4

Label Encoding

sklearn.preprocessing

In data mining and statistics, hierarchical clustering is a method of cluster analysis which seeks to build a hierarchy of clusters

Definition

B	16				
C	47	37			
D	72	57	40		
E	77	65	30	31	
F	79	66	35	23	10
	A	B	C	D	E

Compute a distance matrix

EUCLIDEAN DISTANCE

The distance between two clusters has been computed based on the length of the straight line drawn from one cluster to another

Need to determine from where distance is computed

Two most similar parts of a cluster

single-linkage

Two least similar bits of a cluster

complete-linkage

Center of the clusters

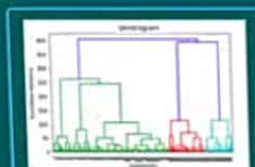
mean or average-linkage

Method

Linkage Criteria

DEFAULT

ward



Hierarchical Clustering

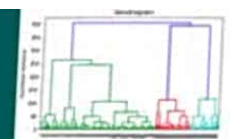
5

Algorithm



5

Algorithm



Hierarchical Clustering

Method

Linkage Criteria

complete-linkage

Center of the clusters

mean or average-linkage

DEFAULT

This method works out which observations to group based on reducing the sum of squared distances of each observation from the average observation in a cluster

ward

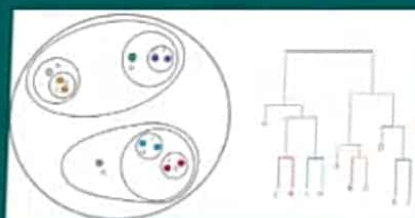
Types

Sequentially merging similar clusters

Agglomerative Hierarchical clustering

Initially grouping all the observations into one cluster, and then successively splitting these clusters

Divisive Hierarchical clustering

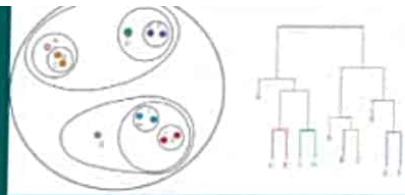


Steps

1 Identify the two clusters that are closest

2 Merge the two most similar clusters





Steps

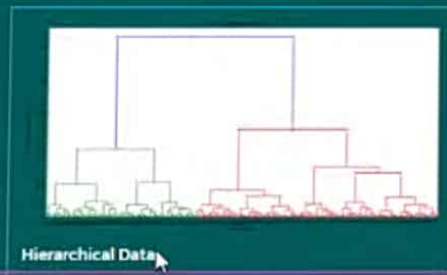
- 1 Identify the two clusters that are closest
- 2 Merge the two most similar clusters

6

Data Visualization - O/P of Hierarchical Clustering



Dendrogram



Hierarchical Data

7

Fitting Model to Clustering algorithm



8

Visualizing Clustered Result

Microsoft Excel interface showing a dataset with columns: CustomerID, Gender, Age, Annual Inc, and Spending Score. The ribbon includes options like Cut, Copy, Paste, Format Painter, Font, Alignment, Number, Styles, Conditional Formatting, and Insert. A small video feed of a person is visible in the bottom right corner.

CustomerID	Gender	Age	Annual Inc	Spending Score
1	Male	19	15	39
2	Male	21	15	81
3	Female	20	16	6
4	Female	23	16	77
5	Female	31	17	40
6	Female	22	17	76
7	Female	35	18	6
8	Female	23	18	94
9	Male	64	19	3
10	Female	30	19	72
11	Male	67	19	14
12	Female	35	19	99
13	Female	58	20	15
14	Female	24	20	77
15	Male	37	20	13
16	Male	22	20	79
17	Female	35	21	35
18	Male	20	21	66
19	Male	52	23	29
20	Female	35	23	98
21	Male	35	24	35
22	Male	25	24	73
23	Female	46	25	5
24	Male	31	25	73
25	Female	54	28	14
26	Male	29	28	82
27	Female	45	28	32
28	Male	35	28	61
29	Female	40	29	31

```
[1] import pandas as pd
import matplotlib.pyplot as plt
```

Load Dataset from Local Directory

```
from google.colab import files
uploaded = files.upload()
```

Choose Files dataset.csv

- dataset.csv(application/vnd.ms-excel) - 3973 bytes, last modified: 6/25/2021 - 100% done

Saving dataset.csv to dataset.csv





+ Code + Text

✓ RAM
Disk

Editing

✓ [3] dataset = pd.read_csv('dataset.csv')

Summarize Dataset



```
print(dataset.shape)
print(dataset.describe())
print(dataset.head(5))
```

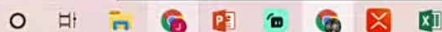


Label Encoding

✓ 0s completed at 19:24

Waiting for colab.research.google.com...

Type here to search



30°C Mostly cloudy



✓ 0s	[4]	4	5	Female	20	10	0
		3	4	Female	23	16	77
		4	5	Female	31	17	40

Label Encoding

```
from sklearn import preprocessing
label_encoder = preprocessing.LabelEncoder()
dataset['Gender'] = label_encoder.fit_transform(dataset['Gender'])
dataset.head()
```

	CustomerID	Gender	Age	Annual Income (k\$)	Spending Score
0	1	1	19	15	39

✓ 1s completed at 19:25

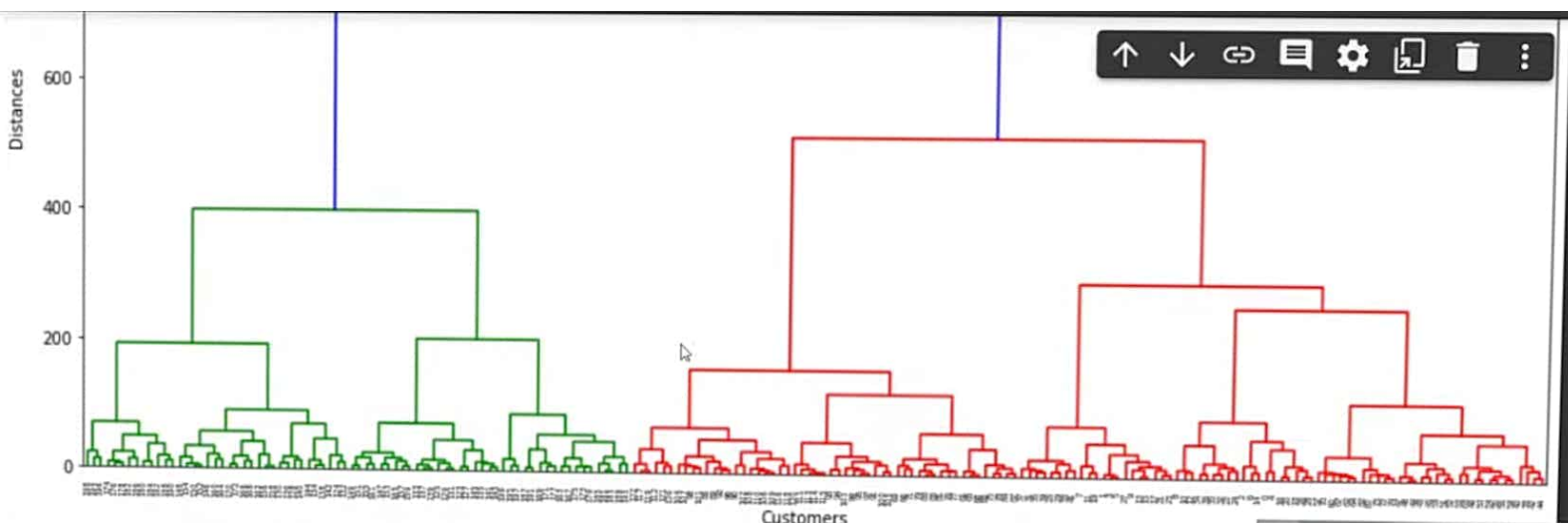


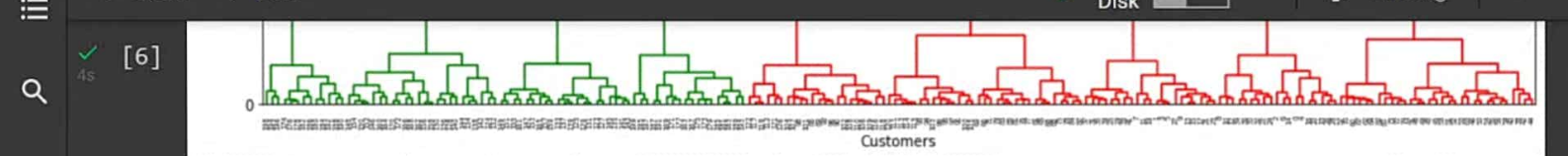
Dendrogram Data visualization

```
import scipy.cluster.hierarchy as clus
plt.figure(1, figsize = (16 ,8))
dendrogram = clus.dendrogram(clus.linkage(dataset, method = "ward"))

plt.title('Dendrogram Tree Graph')
plt.xlabel('Customers')
plt.ylabel('Distances')
plt.show()
```







▼ Fitting the Hierarchial clustering to the dataset with n=5

```
from sklearn.cluster import AgglomerativeClustering
model = AgglomerativeClustering(n_clusters=5, affinity='euclidean', linkage='average')
y_means = model.fit_predict(dataset)
y_means
```

▼ Visualizing the number of clusters n=5

✓ 4s completed at 19:28





+ Code + Text

RAM
Disk

Editing

Cluster 1: Customers with Medium Income and Medium Spending

Cluster 2: Customers with High Income and High Spending

Cluster 3: Customers with Low Income and Low Spending

Cluster 4: Customers with High Income and Low Spending

Cluster 5: Customers with Low Income and High Spending

```
[ ] X = dataset.iloc[:, [3,4]].values
plt.scatter(X[y_means==0, 0], X[y_means==0, 1], s=50, c='purple', label = 'Cluster 1')
plt.scatter(X[y_means==1, 0], X[y_means==1, 1], s=50, c='orange', label = 'Cluster 2')
plt.scatter(X[y_means==2, 0], X[y_means==2, 1], s=50, c='red', label = 'Cluster 3')
plt.scatter(X[y_means==3, 0], X[y_means==3, 1], s=50, c='green', label = 'Cluster 4')
plt.scatter(X[y_means==4, 0], X[y_means==4, 1], s=50, c='blue', label = 'Cluster 5')
```

✓ 0s completed at 19:30



```
+ Code + Text RAM Disk Editing
plt.scatter(X[y_means==4, 0], X[y_means==4, 1], s=50, c='blue', label
plt.title('Income Spent Analysis - Hierarchical Clustering')
plt.xlabel('Income')
plt.ylabel('Spent')
plt.show()
```