

Hierarchical Clustering

INDEX

DATE:

- Hierarchical clustering is very common in biology and lend itself nicely to visualizing clusters.
- It can also help the user decide on an appropriate number of clusters.

Section Overview:-

- Theory and Intuition of Hierarchical Clustering
- Coding Example of Hierarchical Clustering
- .

Hierarchical clustering

Theory and Intuition

- Like most clustering algorithms, Hierarchical clustering simply relies on measuring which data point are most "similar" to other data points.
- "Similarity" is defined by choosing a distance metric
- So why use Hierarchical clustering?

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

M T W T F S

DATE: _____

• Divide points into potential clusters:

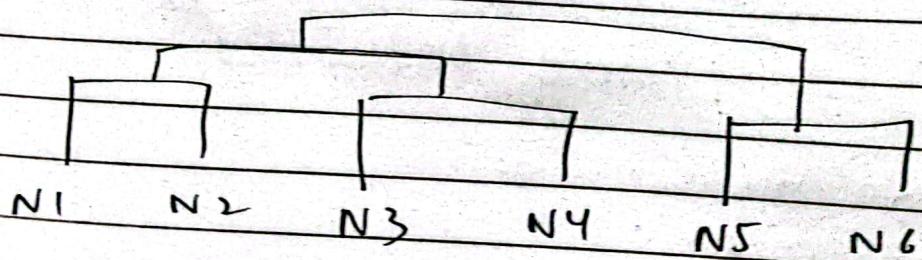
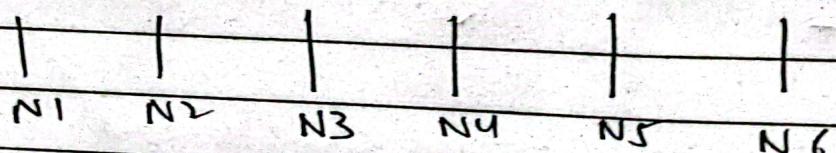
• Agglomerative Approach:

- Each point begins as its own cluster, then clusters are joined-

• Divisive Approach:

- All points begin in the same cluster, then clusters are split-

→ • Agglomerative



— Opposite of the Agglomerative approach is a Divisive approach, which start with all point belonging to the same cluster, and the begin division to separate out clusters.

Hierarchical Clustering process:

- Compare data points to find similar data point 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.

There are a few key topics we still need to understand for Hierarchical clustering:

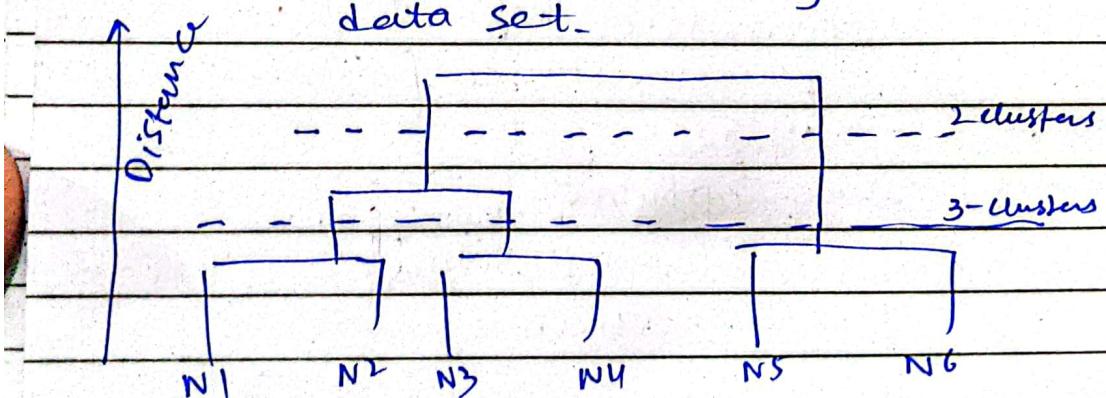
- Similarity Metric
- Dendrogram
- Linkage matrix
- Similarity metrics
 - Measure distance between two points
 - Many options
 - Euclidean Distance
 - Manhattan
 - Cosine
 - And many more
 - Default choice is Euclidean

$$d(p_1, q) = \sqrt{(q_1 - p_1)^2 + (q_2 - p_2)^2}$$

- Each dimension would be a feature
- for n data points and p features:
 - $D^2 = (x_{11} - x_{12})^2 + \dots + (x_{n_1 p_1} - x_{n_1 p_2})^2$
- Using MinMax Scaler we can scale all features to be between 0 and 1.
- This allow for maximum distance between a features to be 1.

Dendrogram:-

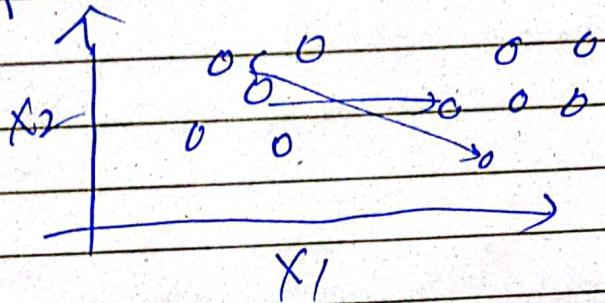
- plot displaying all potential clustering.
- Very computationally expensive to compute and display for larger data set.



Linkage:-

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

- Once two or more points are together and we want to continue agglomerative clustering to join cluster we need to decide on a linkage parameter-



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

Ward: minimize variance of clusters being merged -

Average: use average distance between two sets -

- Minimum or Maximum distances between all observation of the two sets -

003 Hierarchical clustering

Coding # part 1

M T W T F S

DATE: _____

```
df = pd.read_csv('cluster_mpg.csv')
```

df

```
df.describe()
```

```
df['origin'].value_counts()
```

```
df_w_dummies = pd.get_dummies(df.drop("name", axis=1))
```

df_w_dummies

```
from sklearn.preprocessing import MinMaxScaler
```

sqrt!

```
scaler = MinMaxScaler()
```

```
scaled_data = scaler.fit_transform(df_w_dummies)
```

```
Scaled_df = pd.DataFrame(scaled_data, columns=df_w_dummies)
```

Scaled_df

```
plt.figure(figsize=(15, 8))
```

```
sns.heatmap(Scaled_df, cmap=''))
```

```
sns.clustermap(Scaled_df)
```

, row_cluster=False

```
(Scaled_df, col_cluster=False)
```

h

(b)
cl

f

mo

clu

plt
sns

np.

#

le

n

- Hierarchical clustering

coding part two

DATE:

Cluster and Dendrogram

code

from sklearn.cluster import

Agglomerative
clustering

model = AgglomerativeClustering(n_clusters=4)

cluster_labels = model.fit_predict(scaled_df)

plt.figure(figsize

sns.scatterplot(data=df, x='mpg', y='height')

hue =

cluster
labels

palette=brgidis)

hyperparameter

affinity: "euclidean".

max distance possible?

sqrt(N-features)

len(scaled_df.columns)

np.sqrt(10) # max distance possible

between two points

(min/max scaling)

model = AgglomerativeClustering(n_clusters=
None,
distance_threshold=0)

cluster_labels = model.fit_predict
(scaled_df)

then ~~we get~~ -

```
from scipy.cluster.hierarchy import  
dendrogram  
from scipy.cluster import hierarchy
```

linkage_matrix = hierarchy.linkage(
model.children_)

plt.figure(figsize=(20, 10),)

dendro = ~~hierarchy~~ dendrogram(linkage_matrix)

truncation-mode = 'lastp', p=10
'level', p=3

np.sqrt(len(scaled_df.columns))

scaled_df['mpg'].idxmin()

car_a = scaled_df.loc[370]

car_b = scaled_df.loc[28]

MODULES

DATE: _____

$$\text{distance} = \text{np.linalg.norm}(\text{car}_a - \text{car}_b)$$

distance (2.3857)

x)

DBSCAN

M T W T F S

DATE: _____

M T W

- DBSCAN - Density-based spatial clustering of application with noise is a powerful technique which can be used for clustering and outliers detection.

- Let's review what this section will cover

- Section Overview:-

- DBSCAN vs. KMEAN clustering
- DBSCAN Hyperparameters Theory
- DBScan Hyperparameters coding
- Outlier Project Exercise
- Project solutions -

Theory and Intuition

DBSCAN = Density-based spatial clustering of application with noise

- 1972: Robert F. Ling published a closely related algorithm in "The theory and construction of k-clusters" with an expected run time of $O(n^3)$.

- This means that as n number of points grows, the run time if

DBSCAN

- DBSCAN
points are
clustered

- This
segment
have different
Consider

the algorithm grows cubically!

- 1996: Martin Ester, Hans-Peter Kriegel, Jörg Sander and Xiaowei Xu proposed the modern version of DBSCAN with a runtime of $O(n^2)$.
- 2014: DBSCAN was awarded the test of time award at the leading data mining conference, SIGKDD.

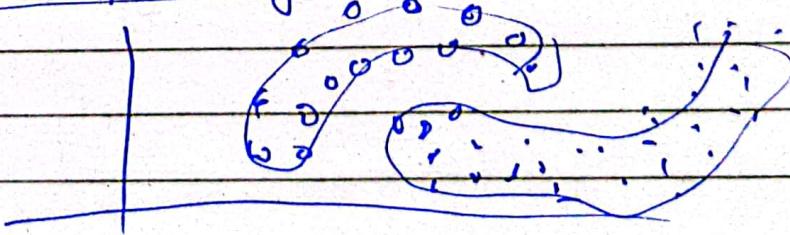
Question

- How it's work?
- Advantages and disadvantages
- How it deal outliers and noise

DBSCAN Key Idea:

- DBSCAN focuses on using density of points as its main factor for assigning cluster label.
- This creates the ability to find cluster segmentations that other algorithms have difficulty with.

Consider the following :-



- DBSCAN iterates through points and uses two key value hyperparameters (epsilon and minimum number of points) to assign clusters labels.
 - Unlike Kmean, it focuses on density as the main factor for clustering assignment of points.

- DBSCAN key hyperparameter:

- Epsilon
 - Distance extended from a point
- Minimum number of points
 - minimum number of points in an epsilon distance -

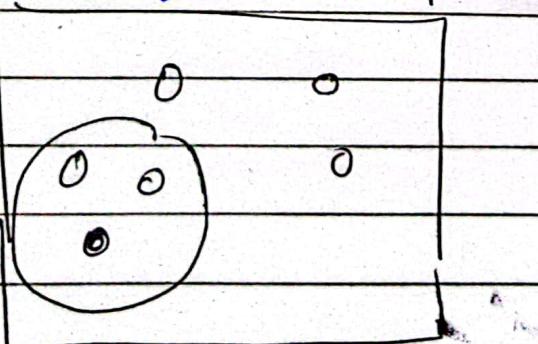
- DBSCAN point Types:

- Core
- Border
- Outlier

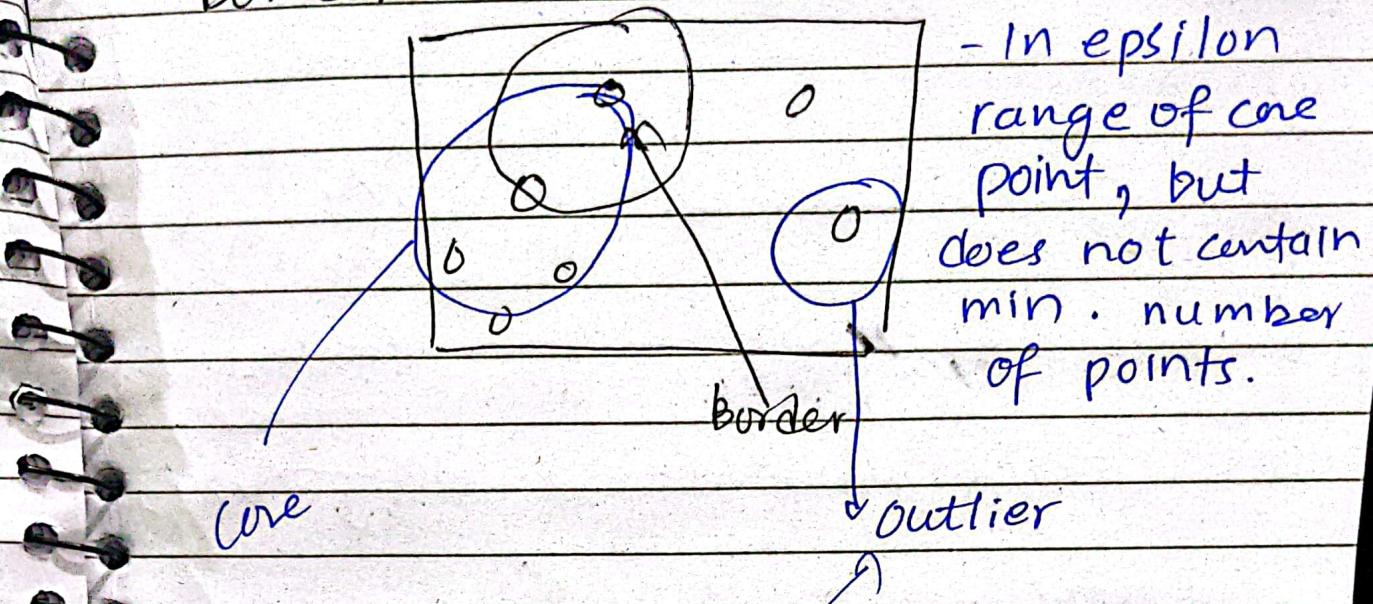
$$\epsilon = 1 \text{ and Min points} = 2$$

Core :

points with min points in epsilon range -
(include itself)



Border: $\epsilon = 1$ and min points = 3



Outlier:-

Can not be "reached" by points in a cluster assignment.

- DBSCAN procedure:

- pick a random point not yet assigned-
- determine the point type
- Once a core point has been found, add all directly reachable points to the same cluster as core-
- Repeat until all points have been assigned to a cluster or an outlier.-

Coding

M T W T F S

DATE: _____

```
blobs = pd.read_csv('.../data/cluster_blobs')
blobs.head()
```

```
sns.scatterplot(data=blobs, x='X1',
                  y='X2')
```

```
moons = pd.read_csv('.../data/cluster-
                      moons.csv')
```

```
sns.scatterplot(data=moons, x='X1', y='X2')
```

```
circles = pd.read_csv('.../cluster_circle.
                       csv')
circle.head()
```

```
def display_categorical(model, data):
    labels = model.fit_predict(data)
    sns.scatterplot(data=data, x='X1',
                    x='X2',
                    hue=labels,
                    palette='Set1')
from sklearn.cluster import KMeans
```

```
disp: model = KMeans(n_clusters=2)
display_categorical(model, moons)
```

```
from sklearn.cluster import DBSCAN
```

```
model = DBSCAN()
```

```
display_categorical(model, blobs)
```

model = DBSCAN(eps=0.15)
display = categories(model, moons)

Hyperparameter theory :-

Epsilon :

- Distance extended from a point to search for Min. Number of points.

Min. Number of points :

- Minimum number of points within Epsilon distance to be a core point

Adjusting these hyperparameter have two main outcome :

- changing number of clusters.
- changing what is an outlier point.

◦ Epsilon Intuition :-

- Increasing epsilon allows more points to be core points which also result in more border points and less outlier points

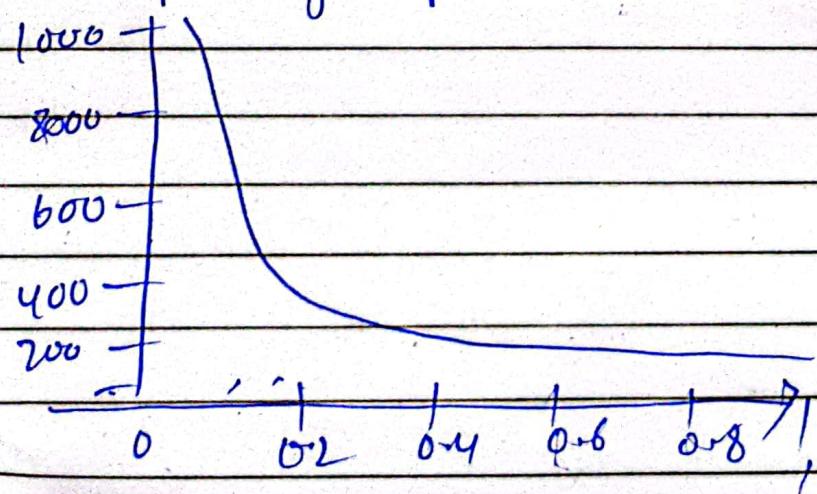
- Imagine a huge epsilon ; all points would be within the neighbourhood and classified as the same cluster !

- Decreasing epsilon causes more points not to be in range of each other, creating more unique clusters-
- Imagine a tiny epsilon , the range would ~~be~~ not extend far out enough to come into contact with any other pointc

Methods for finding an epsilon value:

- o Extremely dependent on the particular dataset and domain spaces-
- o Requires user to have some expectation or intuition about number of clusters and relative percentage and relative percentage of outliers-

plot 'elbow/knee" diagram
comparing epsilon value



- Minimum Number of Samples/points
- Number of 1
- Increasing to a larger number of samples needed to be considered a one point, causes more points to be considered unique outliers.
- Imagine if min. number of sample was close to total number of points available, then very likely all points would become outliers.

Choosing Minimum number of Sample

- Test multiple potential values and chart against number of outliers labeled-

useful to increase to create potential new small clusters, instead of complete outliers.

Coding

M T W T F S

DATE: _____

```
two_blobs = pd.read_csv("two-blobs")  
two_blobs_outlier = pd.read_csv("i")  
sns.scatterplot(data=two_blobs,  
                 x='X1',  
                 y='X2')
```

SOS: two_blobs_outlier

```
def display_categories(model, data):  
    labels =
```

```
    sns.scatterplot(data=
```

```
from model= DDSCAN()
```

```
display_categorical(dbSCAN, two_blobs_cu)
```

```
dbSCAN = DBSCAN(eps=10.0)  
only outliers:
```

```
dbSCAN = DDSCAN(eps=10)  
one whole category
```

```
np.sum  
(dbSCAN.labels == -1)
```

Percent of Outliers

```
100 * np.sum(dbSCAN.labels == -1) / len
```

out

dbSCAN.labels

MONDAYS

DATE: _____

outliers_percent = []

number_of_outliers = []

for eps in np.linspace(0.001, 1, 100):

//

// check
// jupyter
notebook

s-a)

-1)

1. ans. code