

Clustering \rightarrow Grouping similar items together.

why clustering \rightarrow To group similar items together & take actions as a group.

in Industry

Action

1) Retail

- Group similar customers.
- Initiate different campaigns for different groups of customers.

2) Banking

- Group similar transactions together.
- Easier to identify fraudulent transactions.

3) Healthcare

- Group similar patients together.
- Prescribe medicine based on the group a patient is part of.

K-means \rightarrow Group K groups, to allocate each data point to one of the K groups.

Algorithm \rightarrow

1) Choose a K value.

2) Assign K points randomly as cluster centers.

3) Assign each data point to one of the groups.

4) Update the centroid for groups by calculating the average point of all the points in the group.

5) Repeat step (3-4) until clusters stop changing.

Similarity \rightarrow Euclidean distance b/w two points.

Distance \propto $\sqrt{1 - \text{similarity}}$.

Centroid \rightarrow Center of a cluster. (similar to center of gravity).

- May or may not be an actual data point.

- Representative of a group.

Goal of K means \rightarrow 1) Minimize intra cluster distances \rightarrow

Data points inside a cluster are as close to the centroid

as possible.

2) Minimize inter-cluster distance & make clusters as far apart from each other as possible.

\rightarrow Metrics in K-means ->

1) WSS -> Within cluster sum of square.

- measure how close the points are to the centroid.
- not a good metric -> A large cluster can have large WSS value even if it is tightly packed.

WSS & compactness (closeness of points).

2) BCSS -> Between cluster sum of square.

- how far apart the centers of each cluster are from a same center point.

3) TSS -> Total sum of square.

- Measures the total variability in the dataset.

$$TSS = WSS + BSS.$$

\rightarrow K-means works both on small & large datasets.

1) Calculate distance b/w points & centroids.

& Assign each point to the nearest centroid.

2) Recompute the centroid for each group of points.

def K-mean(points, K, ^{Initial}Centroids, iterations)

(List of points

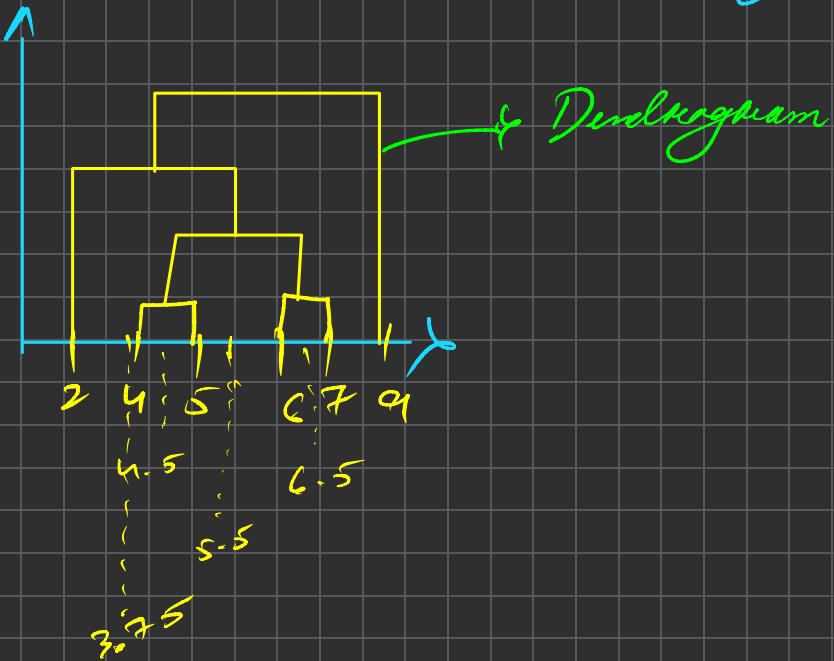
tuple of (x,y)

) list of points of size K

Hierarchical Clustering

- works poorly on larger datasets.
- K-means works both on small & large datasets.

Dendrogram \rightarrow The graph created by hierarchical clustering.



Hierarchy \rightarrow Form a tree like structure.

at agglomeration \rightarrow (Bottom up) \rightarrow Merge smaller clusters into bigger ones.

- Common
 - 1) Each data point starts as a single cluster.
 - 2) At each step two most similar clusters are combined into a single cluster.
 - 3) Repeat until there's only a single cluster remaining.

by D divisive \rightarrow (Top Down) \rightarrow Split a single cluster into two & repeat until no more splits are possible.

when to use hierarchical clustering \rightarrow

- 1) When the dataset is small.
- 2) When detailed cluster hierarchy is required.
 - To visualize & make decisions.

K-Mean Clustering using sklearn library →

from sklearn.cluster import KMeans
 from sklearn.preprocessing import StandardScaler
 from sklearn.metrics import silhouette_score

$$x = df[[\text{feature}x, \text{feature}y]]$$

scaler = StandardScaler()

$$xs = \text{scaler}.fit_transform(x)$$

WRS = []
 silhouettes = []

krange = range(2, 11)

for k in krange:

Kmean = KMeans(n_clusters=k, random_state=42, n_init=10)

Kmean.fit(xs)

WRS.append(Kmean.inertia_)

silhouettes.append(silhouette_score(xs, Kmean.labels_))

lastK = Krange[silhouettes.index(max(silhouettes))]

→ This was to find the best k.

→ Below code is for using the k-mean model.

- Finding clusters.

Kmean.cluster_centers_ → Position of every cluster center.

list of

scaler.inverse_transform(Kmean.cluster_centers_)

- Position of cluster centers scaled to the dataset.

clusters = Kmeans.fit_predict(xs)

↓ Create k-means with the last cluster.

Kmeans = KMeans(n_clusters=lastK, random_state=42, n_init=10)

Clustering on text data

from naderSentiment import SentimentIntensityAnalyzer
from sklearn.feature_extraction.text import CountVectorizer

tokens = df['tokens'].tolist()

analyser = SentimentIntensityAnalyzer()

for tweet in tokens:

pos = analyser.polarity_scores(tweet)

sentiments.append({ 'tweet': tweet, 'compound': pos['compound'] })

df = pd.DataFrame(sentiments)

sentiments = []

for tweet in df['compound']:

if tweet > 0 :

sentiment += 'Pos'

elif tweet < 0 :

sentiment += 'Neg'

else

sentiment += 'neutral'