

3. put an Edge between all core points that are within ϵ
4. Make Each group of Connected Core points to Separate Cluster.
5. Assign Each Border point to one of Cluster of its Associated Core points.

* Strength & Weakness:- DBSCAN uses a density-based defn of a cluster

* It is relatively resistant to Noise & can handle clusters of arbitrary Shapes & Sizes

* DBSCAN has trouble when the clusters have widely varying densities.

* It also has trouble with high-dimensional data because density is more difficult to define for such data.

* DBSCAN can be expensive when the computation of Nearest Neighbours requires computing all pairwise proximities.

* Cluster Evaluation :- (Cluster Validation)

The Evaluation of the resulting Classification Model is an integral part of the process of developing a Classification Model & there are well-separated Measures & procedures Ex "Accuracy & Cross-validation" etc.

→ "Cluster Evaluation" should be a part of any cluster Analysis, A key Motivation is that almost Every clustering Algo will find clusters in a data set, Even if that data set has no Natural Cluster Structure.

→ The Cluster Evaluation/Validity are traditionally classified into the following three types

- * Unsupervised Cluster Evaluation

- * Supervised Cluster Evaluation

- * Relative Cluster Evaluation

- * Unsupervised Cluster Evaluation :- Measures the goodness of a clustering structure without considering External Info. Ex: Sum of Squared Error (SSE)

→ Unsupervised Measures of Cluster Validity are often further divided into two classes

- * Measure of Cluster Cohesion (Compactness)

- * Measure of Cluster Separation (Isolation)

- * Measure of Cluster Cohesion & Separation → "Cluster Cohesion" determines how closely related the objects in a cluster.

"Cluster Separation" determines how distinct or well separated a cluster is from other clusters.

→ Cohesion & Separation for a cluster can be expressed using the following equations

$$\text{Cohesion } (C_x) = \sum_{\substack{x \in C_x \\ y \in C_x}} \text{proximity}(x, y)$$

$$\text{Separation } (C_x, C_y) = \sum_{\substack{x \in C_x \\ y \in C_y}} \text{proximity}(x, y)$$

→ In general, we can Express overall Cluster Validity for a Set of k -clusters as a Weighted Sum of Validity of Individual Clusters

$$\text{Overall Validity} = \sum_{i=1}^k w_i \text{ Validity}(C_i)$$

The "Validity" function can be Cohesion, Separation or Some Combination of these Quantities.

* "Silhouette Co-efficient" → is a popular Method which combines both Cohesion & Separation.

The foll Steps Show how to Compute It :

- 1) for i^{th} object, Calculate its Average distance to all other objects, call this value " a_i "
- 2) for i^{th} object, Calculate its Average distance to all other objects & find the Min such value, call this value " b_i "
- 3) Compute Silhouette Co-efficient. $S_i = (b_i - a_i) / \max(a_i, b_i)$

The Value of the Silhouette Co-efficient can Vary between -1 & 1.

* Measuring Cluster Validity via Correlation → If we are given Similarity Matrix for a data Set & the cluster labels from a Cluster Analysis of data Set.

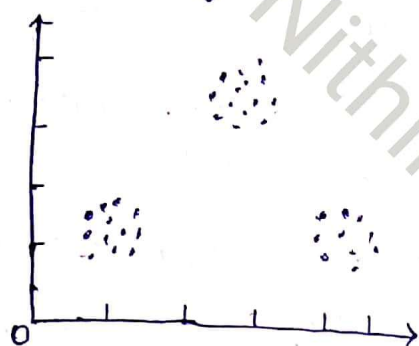
then we can Evaluate "goodness" of the clustering by looking at the Correlation between Similarity Matrix & an Ideal version of Similarity Matrix based on Cluster labels.

"1" to all points within the cluster & similarity of "0" to all points in other clusters.

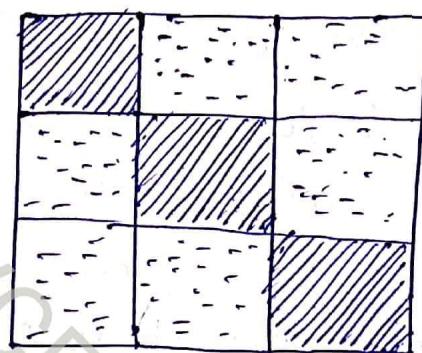
Thus if we sort the rows & columns of similarity matrix so that all objects belonging to same class together, then an ideal similarity matrix has "block diagonal" structure

→ High Correlation b/w Ideal & Actual Similarity Matrices indicates that the points that belong to same cluster are close to each other while low-correlation indicates the opposite.

→ Consider the foll example



(Well-Separated clusters)



(Similarity Matrix Sorted)

If we have well-separated clusters, then the similarity matrix should be "roughly block-diagonal".

→ This approach may seem hopelessly expensive for large data sets, since the computation of proximity matrix takes $O(n^2)$ time.

* Supervised Measures of Cluster Evaluation :- These approaches measure the extent to which two objects that are in the same class are in the same cluster & vice versa.

→ There are two types of supervised measures.

* Classification-Oriented

* Similarity-Oriented

* Classification-oriented Measures \rightarrow are a no of Measures such as "Entropy", "purity", "precision", "Recall" & "F-Measure" that are commonly used to Evaluate the performance of a Classification Model.

\rightarrow "Entropy" :- The degree to which Each Cluster consists of objects of a single class.

The total Entropy for a set of clusters is calculated as the sum of Entropies of Each cluster weighted by the size of Each Cluster i.e.

$$E = \sum_{i=1}^k \frac{m_i}{m} e_i$$

\rightarrow "purity" :- Another Measure of the extent to which a cluster contains objects of a single class.

The overall purity of a clustering is

$$\text{purity} = \sum_{i=1}^k \frac{m_i}{m} p_i$$

\rightarrow "precision" :- The fraction of a cluster that consists of objects of a specified class.

The precision of cluster i w.r.t class j is

$$\text{precision}(i, j) = p_{ij}$$

\rightarrow "Recall" :- The extent to which a cluster contains all objects of a specified class.

The Recall of cluster i w.r.t class j is

$$\text{Recall}(i, j) = m_{ij} / m_j$$

\rightarrow "F-Measure" :- A combination of both precision & Recall that Measures the extent to which a cluster contains only objects of a particular class & all objects of that class.

The F-Measure of cluster i w.r.t class j is

$$F(i, j) = (2 \times \text{precision}(i, j) \times \text{Recall}(i, j)) / (\text{precision}(i, j) + \text{Recall}(i, j))$$

* Similarity - oriented Measures → In this Approach, the Cluster Evaluation will be done using the comparison of two Matrices

* Ideal Cluster Similarity Matrix

* Ideal class Similarity Matrix.

"Ideal Cluster Similarity Matrix", which has a "1" in the i th Entry if two objects are in same cluster, otherwise "0".

"Ideal class Similarity Matrix", which has a "1" in the i th Entry if two objects are in same class, otherwise "0".

* Relative Methods of Cluster Evaluation :- Compares different Clusterings on Clusters. A relative Cluster Evaluation is a Supervised or Unsupervised Evaluation Measure that is used for the purpose of comparison.

→ Relative Measures are not actually a separate type of Cluster Evaluation Measure, but are instead a specific use of such Measures.

Ex:- Two K-Means Clusterings can be compared using either "SSE" or "Entropy".

* Density - Based Clustering :-

Density Based Clustering Algorithms has played a vital role in finding non linear shape structure based on the density.

→ "Density - Based Spatial Clustering of Application with Noise" (DBSCAN) is a Density Based Algorithm.