

Cluster Validation (Cluster Evaluation)

Evaluation of Clustering

⇒ As of now, we have applied various clustering algorithms on data set. Next thing is, how to evaluate whether the clustering results are good or not.

⇒ Major tasks of clustering evaluation includes

① Assessing cluster tendency :-

In this, we assess whether dataset consists of non-random structure in the data. Clustering analysis is meaningful when there is a non-random structure in the data. This can be done using Hopkins' statistic as in (1).
... in a dataset.

Hopkin's statistic :-

- (1) we generate p points that are randomly distributed across the data space, ^{uniformly} i.e. each point has the same probability of being included in the sample.
- (2) Also, sample p actual data points from dataset D .
- (3) Let $u_{i,0}$ be nearest neighbour distance of the artificially generated points and $w_{i,0}$ be nearest neighbour distance of the sample of points from the original dataset.

$$\text{Hopkin statistic}(H) = \frac{\sum_{i=1}^p w_{i,0}}{\sum_{i=1}^p u_{i,0} + \sum_{i=1}^p w_{i,0}}$$

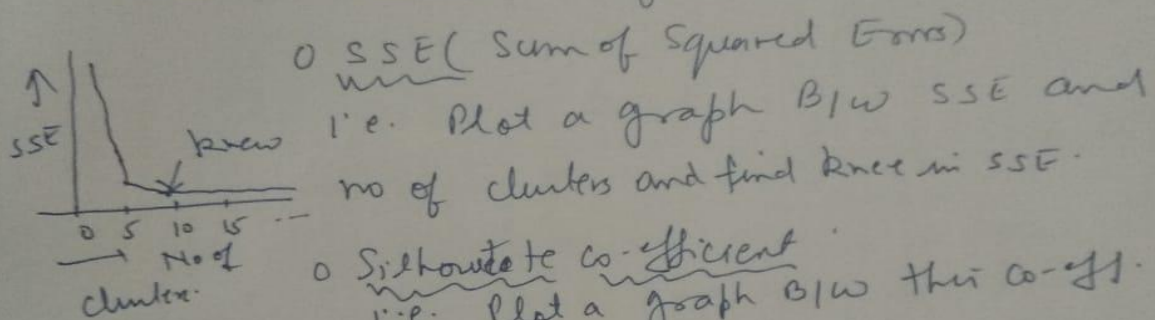
- if both points have same nearest neighbour distances, H is nearly 0.5

•

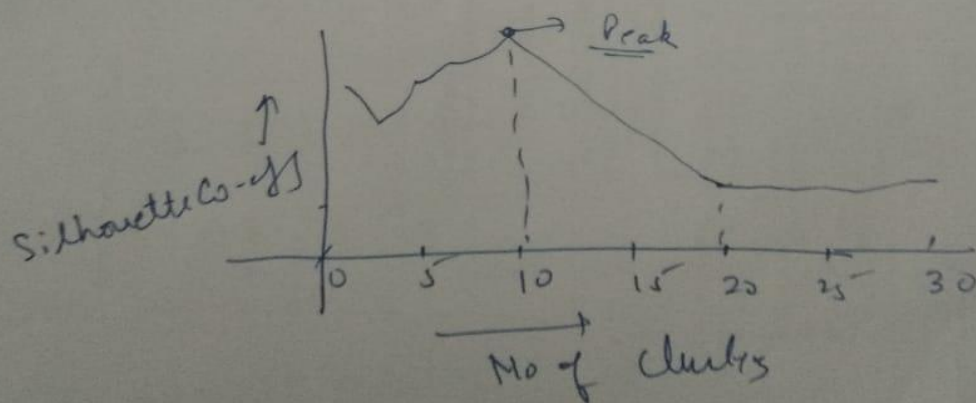
$H=0 \Rightarrow$ High clustering data
 $H=1 \Rightarrow$ randomly distributed

- ② Determining the number of clusters in a dataset.
It is desirable to estimate this number even before a clustering algorithm is used to derive detailed clusters.

Various unsupervised evaluation measures can be used to approximately determine the correct or natural no of clusters.



o Silhouette Co-efficient
i.e. Plot a graph B/w this co-eff. and no of clusters and find a distinct peak in the silhouette co-eff.



③ Measuring Cluster Quality :-

→ How good clustering method resulted into clusters?

→ Two methods :-

- o Extrinsic Methods (when ground truth is known)
- o Intrinsic Methods (when no ground truth is known)

Intrinsic Methods :-

- o Ground truth is not known
- o Hence, this is also known as unsupervised cluster evaluation.

In general, overall cluster validity for a set of K clusters is a weighted sum of the validity of individual clusters.

$$\text{overall validity} = \sum_{i=1}^K w_i \cdot \text{validity}(C_i)$$

Validity function can be cohesion, separation or some combination of them. Weights may vary on validity measure.

In some cases, weights are simple 1 or

→ Size of the cluster or

→ Square root of cohesion and others.

Higher value of cohesion and lower value of separation are better.

(i) Cohesion (Validity function)

→ Let weight = 1

$$\rightarrow \text{Cohesion}(C_i) = \sum_{x \in C_i} \text{Proximity}(x, c_i)$$

c_i is the centroid of cluster C_i

x belongs to cluster C_i

$$\therefore \text{Cohesion}(C_i) = \text{dist}(x_1, c_i) + \text{dist}(x_2, c_i) + \text{dist}(x_3, c_i) + \dots + \text{dist}(x_n, c_i)$$

$$(x_1, x_2, \dots, x_n) \in C_i$$

$$\therefore \boxed{\text{Overall Validity} = \sum_{i=1}^K w_i \text{Cohesion}(C_i)}$$

$$= \text{Cohesion}(C_1) + \text{Cohesion}(C_2) + \dots + \text{Cohesion}(C_K)$$

K = No of clusters.

+ Note:- if proximity measure is Euclidean distance, cohesion is the cluster SSE (Sum of squared errors).

(ii) Separation (Validity function)

$$\text{Separation}(\mathcal{C}_i, \mathcal{C}_j) = \text{Proximity}(c_i, c_j) \text{ --- (1)}$$

c_i & c_j are centroid of \mathcal{C}_i and \mathcal{C}_j (\mathcal{C}_i and \mathcal{C}_j are clusters)

$$\text{Separation}(\mathcal{C}_i) = \text{Proximity}(c_i, c) \text{ --- (2)}$$

c_i = centroid of cluster \mathcal{C}_i

c = overall centroid (mean of all points)

(iii) Silhouette Co-eff: (Combination of (i) & (ii))

The following steps are taken to compute this Co-eff for an individual point.

(a) For the i^{th} object, calculate the average distance to all other objects in its cluster. Call this value a_i .

(b) For the i^{th} object and any cluster not containing the object, calculate object's average distance to all the objects in the given cluster. Find the minimum such value with respect to all clusters, call this value b_i .

(c) for the i^{th} object, silhouette co-eff is defined as

$$S_{i^0} = \frac{(b_{i^0} - a_{i^0})}{\max(a_{i^0}, b_{i^0})} \quad \text{--- (3)}$$

$$\boxed{-1 \leq S_{i^0} \leq 1}$$

→ negative value is undesirable

→ positive value is desirable.

→ overall measure of goodness of clustering can be obtained by compute the average silhouette co-eff of all points.

$$\text{overall silhouette co-eff (S)} = \frac{1}{n} \sum_{i=1}^n S_{i^0}$$

n = no of data points

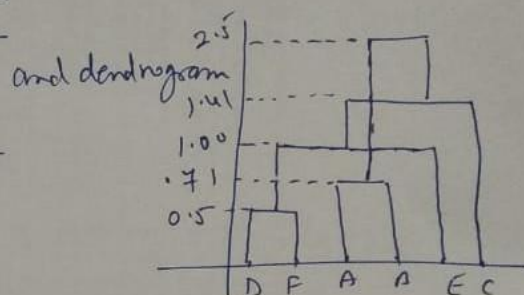
S_{i^0} = calculated from equation 3.

(II) Hierarchical Clustering :-

- These algorithms can be compared by calculating Cophenetic Correlation Co-efficient
- Cophenetic Coefficient gives the correlation B/w the distance matrix (dissimilarity matrix) and Cophenetic distance matrix.
- Cophenetic matrix is computed using Cophenetic distance which is the proximity at which an agglomerative hierarchical clustering technique puts the objects in same cluster for the first time.
- Eg Let the smallest distance B/w two objects (clusters) that are merged is 0.1, then all points in one cluster have a Cophenetic distance of 0.1 with respect to points in other cluster.
- In Cophenetic distance matrix, entries are the Cophenetic distance B/w each pair of objects.

Eg Let the distance matrix (dissimilarity matrix)

	A	B	C	D	E	F
A	0					
B	0.71	0				
C	5.66	4.95	0			
D	3.61	2.92	2.24	0		
E	4.24	3.54	1.41	1.00	0	
F	3.2	2.5	2.5	0.5	1.12	0



Now compute the cophenetic matrix as follows

(i) @ Firstly as (D, F) are merged with value 0.5,

hence write the entry at (D, F) as 0.5

(ii) Similarly then, (A, B) are merged so place 0.71 for (B, A)

(iii) (D, F) and E
at 1.0 hence
DE & FE = 1.0

	A	B	C	D	E	F
A	0					
B	0.71	0				
C	2.5	2.5	0			
D	2.5	2.5	1.41	0		
E	2.5	2.5	1.41	1	0	
F	2.5	2.5	1.41	0.5	1	0

(iv) ((D, F), E) and C
at 1.41, hence
DC = FC = EC = 1.41

(v) (((D, F), E), C) and (A, B) at
2.5, hence

$$DA = DB = 2.5$$

$$FA = FB = 2.5$$

$$EA = EB = 2.5$$

$$\& \begin{matrix} CA = CB \\ = 2.5 \end{matrix}$$

(ii) Find the correlation co-eff B/w
this matrix (Cophenetic Matrix) and
dissimilarity matrix.

(X) dissimilarity	CP (Y)
① 0.71	0.71
② 5.66	2.50
③ 3.61	2.50
④ 4.24	2.50
⑤ 3.20	2.50
⑥ 4.95	2.50
⑦ 2.92	2.50
⑧ 3.54	2.50
⑨ 2.50	2.50
⑩ 2.24	2.50 1.41
⑪ 1.41	1.41
⑫ 2.50	1.00
⑬ 1.00	0.50
⑭ 0.50	1.00

$$\text{Pearson Co-eff (r)} = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{\sum (x - \bar{x})^2} \sqrt{\sum (y - \bar{y})^2}}$$

\bar{x} = mean of
X variable

\bar{y} = mean of
Y variable.

(iii) Higher the value, better the clustering is.

Extrinsic Methods for cluster evaluation :-

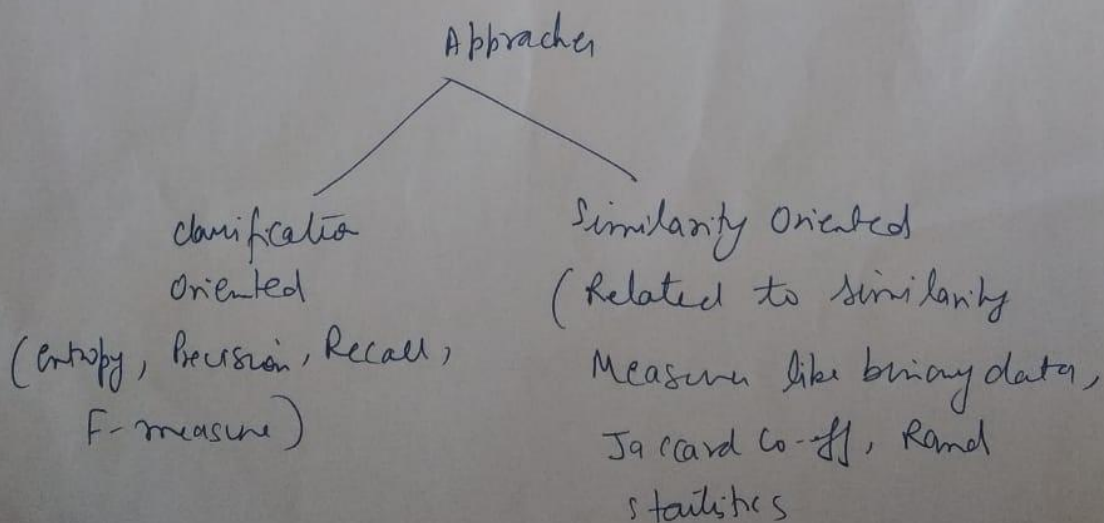
⇒ Ground truth is known (i.e. we know the class labels for data objects)

⇒ When we have class labels, why to do clustering?

Ans :- Clustering analysis helps to automatically do the classification process that was done manually.

⇒ Hence, we do the cluster evaluation using supervised measures.

⇒ These measures calculate the degree of correspondence b/w cluster labels and class labels (ground truth).



A) Similarity-Oriented Measure:

- In this method, ideal cluster similarity matrix and ideal class similarity matrix (as shown in Table 8.10 and Table 8.11) are found and compared.
- Then, a simple matching coefficient such as **Rand Statistic**, **Jaccard coefficient** is computed.

$$\text{Rand statistic} = \frac{f_{00} + f_{11}}{f_{00} + f_{01} + f_{10} + f_{11}}$$

$$\text{Jaccard coefficient} = \frac{f_{11}}{f_{01} + f_{10} + f_{11}}$$

f_{00} = number of pairs of objects having a different class and a different cluster

f_{01} = number of pairs of objects having a different class and the same cluster

f_{10} = number of pairs of objects having the same class and a different cluster

f_{11} = number of pairs of objects having the same class and the same cluster

	Same Cluster	Different Cluster
Same Class	f_{11}	f_{10}
Different Class	f_{01}	f_{00}

- Example of Matrices:

Table 8.10. Ideal cluster similarity matrix.

Point	p1	p2	p3	p4	p5
p1	1	1	1	0	0
p2	1	1	1	0	0
p3	1	1	1	0	0
p4	0	0	0	1	1
p5	0	0	0	1	1

Table 8.11. Ideal class similarity matrix.

Point	p1	p2	p3	p4	p5
p1	1	1	0	0	0
p2	1	1	0	0	0
p3	0	0	1	1	1
p4	0	0	1	1	1
p5	0	0	1	1	1

(i) Ideal cluster similarity matrix has 1 in i th entry if two objects, i and j , are in the same cluster, else "0".

(ii) Ideal class similarity matrix is defined with respect to class labels, which has 1 in i th entry if two objects, i & j , belong to the same class, and "0" otherwise.

as shown in Table 8.10 and 8.11 for the following example:-

Let $D = \{P_1, P_2, P_3, P_4, P_5\}$ has two clusters as $C_1 = \{P_1, P_2\}$
 $C_2 = \{P_3, P_4, P_5\}$

Calculation of Rand statistics (Matching Co-eff)

→ It is easy to define two way

contingency table to determine whether pair of objects are in the same class and same clusters.

		1 (same) 0 (different)
	Same cluster	Different cluster
Same class	f_{11}	f_{10}
Different class	f_{01}	f_{00}

Table 8.10. Ideal cluster similarity matrix.

Point	p1	p2	p3	p4	p5
p1	1	1	1	0	0
p2	1	1	1	0	0
p3	1	1	1	0	0
p4	0	0	0	1	1
p5	0	0	0	1	1

Table 8.11. Ideal class similarity matrix.

Point	p1	p2	p3	p4	p5
p1	1	1	0	0	0
p2	1	1	0	0	0
p3	0	0	1	1	1
p4	0	0	1	1	1
p5	0	0	1	1	1

Note: f_{ij} ($i \Rightarrow$ class, $j \Rightarrow$ cluster)

We need to compute these quantities for all pair of distinct objects. (i.e. total no of pair of objects = $m(m-1)/2$, where m is no of objects in dataset)

$$f_{11} \Rightarrow (p_1, p_2), (p_1, p_3), (p_2, p_3) [2]$$

$$f_{00} \Rightarrow (p_4, p_5), (p_4, p_3), (p_5, p_3), (p_4, p_3) [4]$$

$$f_{01} \Rightarrow (p_1, p_3), (p_2, p_3) [2]$$

$$f_{10} \Rightarrow (p_3, p_4), (p_3, p_5) [2]$$

$$\therefore \text{Rand Statistic} = \frac{(f_{00} + f_{11})}{(f_{00} + f_{11}) + (f_{01} + f_{10})}$$

$$= \frac{6}{6 + (2+2)} = \frac{6}{10} = 0.6$$

$$\therefore \text{Jaccard coefficient} = \frac{f_{11}}{(f_{01} + f_{10}) + f_{11}} = \frac{2}{(2+2) + 2}$$

$$= \frac{2}{6} = \frac{1}{3} = 0.33$$

Classification-Oriented Approaches

These measures evaluate the extent to which two objects that are in the same class. In this context, predicted class labels are the cluster labels.

(i) Entropy :-

→ The degree to which each cluster consists of objects of a single class.

→ entropy of each cluster i^o is :-

$$e_{i^o} = - \sum_{j^o=1}^L p_{i^o j^o} \log_2 p_{i^o j^o}, \text{ where } L \text{ is no of classes.}$$

$$\text{where } p_{i^o j^o} = \left(\frac{m_{i^o j^o}}{m_{i^o}} \right)$$

⇒ m_{i^o} = no of objects in cluster i^o .

⇒ $p_{i^o j^o}$ = the probability that a member of cluster i^o belong to class j^o .

⇒ $m_{i^o j^o}$ = no of objects of class j^o in cluster i^o .

→ Total entropy of a set of clusters is calculated as:

$$C = \sum_{i^o=1}^K \frac{m_{i^o} e_{i^o}}{m} \text{ where } K \text{ is no of clusters.}$$

m is the total no of

points.

e_{i^o} = entropy of i^o th cluster.

(ii) Purity :-

→ The extent to which a cluster contains objects of a single class.

→ Purity of cluster $i^0 = P_{i^0} = \max_{j^0} (P_{i^0 j^0})$

$$\text{Overall purity} = \sum_{i^0=1}^K \left(\frac{m_{i^0}}{m} \right) (P_{i^0})$$

(iii) Precision :-

→ The fraction of a cluster that consists of objects of a specified class.

→ The precision of cluster i^0 with respect to

$$\text{class } j^0 = \text{Precision}(i^0, j^0) = P_{i^0 j^0} = \left(\frac{m_{i^0 j^0}}{m_{i^0}} \right)$$

(iv) Recall :-

→ The extent to which a cluster contains all objects of a specific class.

→ The recall of cluster i^0 with respect to

$$\text{class } j = \text{Recall}(i^0, j) = \left(\frac{m_{i^0 j}}{m_j} \right)$$

m_j = no of objects in class j

(v) F-measure :-

→ Combination of precision and recall.

$$\rightarrow F(i, j) = \frac{2 \times \text{Precision}(i, j) \times \text{Recall}(i, j)}{\text{Precision}(i, j) + \text{Recall}(i, j)}$$

F-measure of cluster i with respect to class j .

Let us take an example shown in following table which is a result of K-means clustering of 3204 news articles consisting of 6 classes from India. No of clusters are 6.

Cluster	Enter-tainment	Financial	Foreign	Metro	National	Sports
1	3	5	40	506	96	27
2	4	7	280	29	39	2
3	1	1	1	7	4	671
4	10	162	3	119	73	2
5	331	22	5	70	13	23
6	5	358	12	212	48	13
Total	354	555	341	943	273	738

- Ideally, each cluster contains documents from only one class.
- But in reality, each cluster contains documents from many classes.
- Many clusters contain documents primarily from just one class.