

CLUSTER VALIDITY INDICES

WHY TO EVALUATE THE “GOODNESS” OF THE RESULTING CLUSTERS?

- To avoid finding patterns in noise
- To compare clustering algorithms
- To compare two sets of clusters
- To compare two clusters

DIFFERENT ASPECTS OF CLUSTER VALIDATION

1. Determining the **clustering tendency** of a set of data, i.e., distinguishing whether non-random structure actually exists in the data.
2. Comparing the results of a cluster analysis to externally known results, e.g., to externally given class labels.
3. Evaluating how well the results of a cluster analysis fit the data *without* reference to external information.
 - Use only the data
4. Comparing the results of two different sets of cluster analyses to determine which is better.
5. Determining the 'correct' number of clusters.

For 2, 3, and 4, we can further distinguish whether we want to evaluate the entire clustering or just individual clusters.

FRAMEWORK FOR CLUSTER VALIDITY

- Need a framework to interpret any measure.
 - For example, if our measure of evaluation has the value, 10, is that good, fair, or poor?
- Statistics provide a framework for cluster validity
 - The more “atypical” a clustering result is, the more likely it represents valid structure in the data
 - Can compare the values of an index that result from random data or clusterings to those of a clustering result.
 - If the value of the index is unlikely, then the cluster results are valid
 - These approaches are more complicated and harder to understand.
- For comparing the results of two different sets of cluster analyses, a framework is less necessary.
 - However, there is the question of whether the difference between two index values is significant

MEASURES OF CLUSTER VALIDITY

- Numerical measures that are applied to judge various aspects of cluster validity, are classified into the following three types.
 - **External Index:** Used to measure the extent to which cluster labels match externally supplied class labels.
 - Entropy
 - **Internal Index:** Used to measure the goodness of a clustering structure *without* respect to external information.
 - Sum of Squared Error (SSE)
 - **Relative Index:** Used to compare two different clusterings or clusters.
 - Often an external or internal index is used for this function, e.g., SSE or entropy
- Sometimes these are referred to as **criteria** instead of **indices**
 - However, sometimes criterion is the general strategy and index is the numerical measure that implements the criterion.

EXTERNAL MEASURES

- The correct or ground truth clustering is known priori.
- Given a clustering partition C and ground truth partitioning T, we redefine TP, TN, FP, FN in the context of clustering.
- Given the number of pairs N

$$N = TP + FP + FN + TN$$

EXTERNAL MEASURES ...

- True Positives (TP): X_i and X_j are a true positive pair if they belong to the same partition in T , and they are also in the same cluster in C . TP is defined as the number of true positive pairs.
- False Negatives (FN): X_i and X_j are a false negative pair if they belong to the same partition in T , but they do not belong to the same cluster in C . FN is defined as the number of false negative pairs.
- • False Positives (FP): X_i and X_j are a false positive pair if they do not belong to the same partition in T , but belong to the same cluster in C . FP is the number of false positive pairs.
- True Negatives (TN): X_i and X_j are a true negative pair if they do not belong to the same partition in T , nor to the same cluster in C . TN is the number of true negative pairs.

JACCARD COEFFICIENT

- Measures the fraction of true positive point pairs, but after ignoring the true negatives as,

$$\text{Jaccard} = \text{TP} / (\text{TP} + \text{FP} + \text{FN})$$

- For a perfect clustering C , the coefficient is one, that is, there are no false positives nor false negatives.
- Note that the Jaccard coefficient is asymmetric in that it ignores the true negatives

RAND STATISTIC

- Measures the fraction of true positives and true negatives over all pairs as

$$\text{Rand} = (\text{TP} + \text{TN}) / N$$

- The Rand statistic measures the fraction of point pairs where both the clustering C and the ground truth T agree.
- A perfect clustering has a value of 1 for the statistic.
- The adjusted rand index is the extension of the rand statistic corrected for chance.

FOWLKES-MALLOWS MEASURE

- Define precision and recall analogously to what done for classification,

$$\text{Prec} = \text{TP} / (\text{TP} + \text{FP}) \quad \text{and} \quad \text{Recall} = \text{TP} / (\text{TP} + \text{FN})$$

- The Fowlkes–Mallows (FM) measure is defined as the geometric mean of the pairwise precision and recall

$$\text{FM} = \sqrt{(\text{precision} \cdot \text{recall})}$$

- FM is also asymmetric in terms of the true positives and negatives because it ignores the true negatives. Its highest value is also 1, achieved when there are no false positives or negatives.

INTERNAL MEASURES: COHESION AND SEPARATION

- **Cluster Cohesion (Compactness)**: Measures how closely related are objects in a cluster.
- **Cluster Separation (Separation)**: Measure how distinct or well-separated a cluster is from other clusters.
- Example: Squared Error
 - Cohesion is measured by the within cluster sum of squares (SSE)

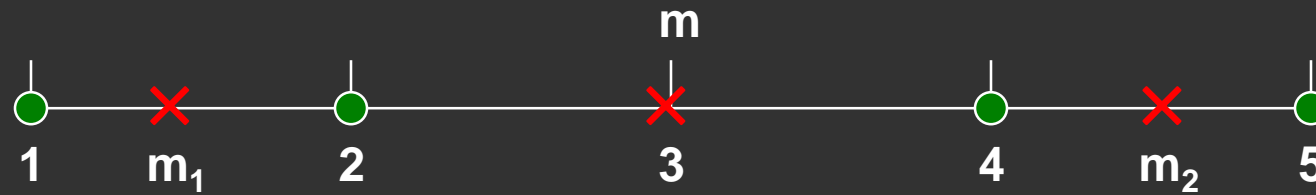
$$WSS = \sum_i \sum_{x \in C_i} (x - m_i)^2$$

- Separation is measured by the between cluster sum of squares

$$BSS = \sum_i |C_i| (m - m_i)^2$$

- Where $|C_i|$ is the size of cluster i

INTERNAL MEASURES: COHESION AND SEPARATION



K=1 cluster:

$$WSS = (1 - 3)^2 + (2 - 3)^2 + (4 - 3)^2 + (5 - 3)^2 = 10$$

$$BSS = 4 \times (3 - 3)^2 = 0$$

$$Total = 10 + 0 = 10$$

K=2 clusters:

$$WSS = (1 - 1.5)^2 + (2 - 1.5)^2 + (4 - 4.5)^2 + (5 - 4.5)^2 = 1$$

$$BSS = 2 \times (3 - 1.5)^2 + 2 \times (4.5 - 3)^2 = 9$$

$$Total = 1 + 9 = 10$$

INTERNAL MEASURES: COHESION AND SEPARATION

- Generally most of the indices used for internal clustering validation combine compactness and separation measures as follow:
- $$\text{index} = \frac{(\alpha \times \text{Separation})}{(\beta \times \text{Compactness})}$$
- Where α and β are weights.

CALINSKI-HARABAZ INDEX:

- Given k clusters, the Calinski-Harabaz score s is given by the ratio of the between-cluster dispersion mean and the within-cluster dispersion

$$CH = \frac{BSS(C)/(k - 1)}{WSS(C)/(N - k)}$$

- The score is higher when clusters are dense and well separated, which relates to a standard concept of a cluster

SILHOUETTE COEFFICIENT

- A measure of how close each point in one cluster is to points in the neighboring clusters.

$$S(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}$$

For each observation i , calculate the average dissimilarity $a(i)$ between i and all other points of the cluster to which i belongs.

For all other clusters C , to which i does not belong, calculate the average dissimilarity $d(i, C)$ of i to all observations of C . The smallest of these $d(i, C)$ is defined as $b(i) = \min_C d(i, C)$. The value of $b(i)$ can be seen as the dissimilarity between i and its “neighbor” cluster, i.e., the nearest one to which it does not belong.

- Observations with a large S_i (almost 1) are very well clustered.
- A small S_i (around 0) means that the observation lies between two clusters.
- Observations with a negative S_i are probably placed in the wrong cluster.

DUNN'S INDEX:

- If the data set contains compact and well-separated clusters, the diameter of the clusters is expected to be small and the distance between the clusters is expected to be large. Thus, Dunn index should be maximized.

$$\text{Dunn's Index } (D) = \frac{\min\{\text{inter-cluster separation}\}}{\max\{\text{intra-cluster distance or diameter}\}}$$

$$Dunn = \min_{1 \leq i \leq K} \min_{1 \leq j \leq K, j \neq i} \frac{\delta(C_i, C_j)}{\max_{1 \leq k \leq K} \Delta(C_k)}$$

Here,

$$\Delta(S) = \max_{\bar{x}, \bar{y} \in S} d(\bar{x}, \bar{y})$$

max{intra-cluster distance or diameter}

$$\delta(S, T) = \min_{\bar{x} \in S, \bar{y} \in T} d(\bar{x}, \bar{y}),$$

min{inter-cluster separation}

DAVIES–BOULDIN INDEX:

- A function of the ratio of the sum of *within-cluster scatter* to *between-cluster separation*.
- The objective is to minimize the DB index for achieving proper clustering.

$$DB = \frac{1}{K} \sum_{i=1}^K R_{i,q,t}$$

Here,

$$R_{i,q,t} = \max_{j, j \neq i} \frac{S_{i,q} + S_{j,q}}{d_{i,j,t}}.$$

$S_{i,q}$ is the q th moment of the total points in cluster C_i w.r.t. their mean \bar{c}_i .

$$S_{i,q} = \left(\frac{\sum_{\bar{x} \in C_i} \{\|\bar{x} - \bar{c}_i\|_2^q\}}{|C_i|} \right)^{\frac{1}{q}}$$

$d_{i,j,t}$ is the Minkowski distance of order t between the centroids \bar{c}_i and \bar{c}_j

$$d_{i,j,t} = \|\bar{c}_i - \bar{c}_j\|_t$$

XIE-BENI INDEX:

- In the definition of XB-index, the numerator indicates the compactness of the obtained cluster while the denominator indicates the strength of the separation between clusters.
- The objective is to minimize the XB-index for achieving proper clustering.

$$XB = \frac{\sum_{i=1}^K \sum_{j=1}^n \mu_{ij}^2 \|\bar{x}_j - \bar{c}_i\|^2}{n(\min_{i \neq k} \|\bar{c}_i - \bar{c}_k\|^2)}.$$

PS INDEX:

$$PS = \frac{1}{K} \sum_{i=1}^K \left[\frac{1}{n_i} \sum_{j \in S_i} \frac{d_c(\bar{x}_j, \bar{c}_i)}{\min_{m,n=1,\dots,K, m \neq n} (d_c(\bar{c}_m, \bar{c}_n))} \right]$$

Given partition of the data set $X = x_j : j = 1, 2, \dots, n$ and the center of each cluster $c_i (i = 1, \dots, K)$.

Here, S_i is the set whose elements are the data points assigned to the i th cluster, n_i is the number of elements in S_i

- The most desirable partition is obtained by minimizing PS-index

I-INDEX:

- $$I = \left(\frac{1}{K} \times \frac{E_1}{E_k} \times D_k \right)^2$$

Here, E_1 is constant for given dataset.

$$E_k = \sum_{k=1}^K \sum_{i=1}^{n_k} \|\bar{x}_i - \bar{c}_k\|^2$$

$$D_k = \max_{i,j=1}^K \|\bar{c}_i - \bar{c}_j\|^2$$

n_i is the number of points in cluster C_i

\bar{c}_k is the centre of k th cluster.

The best partitioning occurs at the maximum value of the I-index.

CS-INDEX:

- Used for tackling clusters of different densities and/or sizes.

$$CS = \frac{\sum_{i=1}^K K \left[\frac{1}{N_i} \sum_{\bar{X}_i \in C_i} \max_{\bar{X}_q \in C_i} d(\bar{X}_i, \bar{X}_q) \right]}{\sum_{i=1}^K [\min_{j \in K, j \neq i} d(\bar{m}_i, \bar{m}_j)]}$$

where $\bar{m}_i, i = 1, \dots, K$ are the cluster centers.