

* Evaluating or clustering

cluster validity

The validation of cluster created is a troublesome task.

→ the problem here is

clusters are "in the eyes of the beholder".

A good cluster will have

- High inter class similarity
- Low intra class similarity.

Aspects of cluster validation

• External → compare your cluster to ground truth.

• Internal → Evaluating the cluster without reference to the external data.

• Reliability → the clusters are not formed by chance (randomly) —

Some statistical framework can be used

External measure

→ no. of objects in the data $P = p_1, p_2, \dots, p_m$
the set of ground truth clusters $C = c_1, c_2, \dots, c_n$
the set of clusters formed by the algorithm.
the Incidence matrix $N \times N$ matrix.

→ $P_{ij} = 1$, if the two points o_i & o_j belong to the same cluster in the ground truth
else, $P_{ij} = 0$.

→ $C_{ij} = 1$, if the two points o_i & o_j belong to the same cluster in the ground truth
else, $C_{ij} = 0$

now there can be following scenarios

1) $C_{ij} = P_{ij} = 1$ → Both the points belong to the same cluster for both our algorithm & Ground truth (Agree) — SD

2) $C_{ij} = P_{ij} = 0$ → Both the points don't belong to same cluster for both our algorithm & Ground truth (Agree) — DD

3) $C_{ij} = 1$ but $P_{ij} = 0$ → the points belong in the same cluster for our algorithm but different cluster in ground truth (Disagree) — SD

4) $C_{ij} = 0$ but $P_{ij} = 1$ → the points don't belong in same cluster for our Algo but same cluster in ground truth (Disagree) — DS

Just like accuracy score

$$\text{Rand Index} = \frac{\text{total agree}}{\text{total disagree}} = \frac{(SS + DD)}{(SS + DD + DS + SD)}$$

→ the disadvantage of this is it can be dominated by DD.

$$\text{Jaccard coefficient} = \frac{SS}{(SS + DS + SD)}$$

∴ A higher value of Rand Index & Jaccard coefficient mean that the clusters generated by our algorithm mostly agree to the ground truth.

confusion matrix

n = no. of points

m_i = points in cluster i

C_j = points in class j

n_{ij} = points in cluster i coming from class j .

$$P_{ij} = \frac{n_{ij}}{m_i}$$

(probability of element from cluster i to be assigned to class j)

| | class 1 | class 2 | class 3 | |
|-----------|--------------|--------------|--------------|-------|
| cluster 1 | n_{11}/P_1 | n_{12}/P_2 | n_{13}/P_3 | m_1 |
| cluster 2 | n_{21}/P_1 | n_{22}/P_2 | n_{23}/P_3 | m_2 |
| cluster 3 | n_{31}/P_1 | n_{32}/P_2 | n_{33}/P_3 | m_3 |
| | C_1 | C_2 | C_3 | n |

→ Entropy of cluster i ,

$$e_i = - \sum P_{ij} (\log P_{ij})$$

→ For evaluating clustering algorithm, entropy can be given as

$$E = \sum \frac{m_i}{n} e_i$$

→ Purity of cluster i , $P_i = \max(P_{ij})$

And for entire cluster it is $P(C) = \sum \frac{m_i}{n} \cdot P_i$

→ Purity, is the total percentage of data points "clustered correctly".

→ A high value of Purity score means that our clustering algorithm performs well against the ground truth.

note

→ In calculating External measure, most of time we don't have ground truth.

① Internal measures

There are the methods we use to measure the quality of clusters without external reference. There are two aspects of it.

Cohesion

How closely the objects in the same cluster are related to each other. It is the within-cluster sum of squared distances. It is the same metric that we used to calculate for K-means Algorithm.

$$WCSS = \sum \sum (x - m_i)^2$$

Separation

How different objects in different clusters are and how different a well separated cluster is from other clusters.

It is the between cluster sum of squared distance

$$BSS = \sum c_i (m - m_i)^2$$

where, c is the size of individual cluster & m is centroid of all data points.

Note

☞ $BSS + WSS$ is always a constant

→ the silhouette can be calculated as

$$s(m) = \frac{b(m) - a(m)}{\max\{a(m), b(m)\}}$$

where,

$a(m)$ is avg. distance of x from all other

$b(m)$ is avg. distance of x from all other point in same cluster.

from all other point in other clusters.

☞ And silhouette coefficient is,

$$SC = \frac{1}{N} \sum s(m)$$

☞ Higher $s(m)$ means that the

inter cluster similarity is less and the

intra cluster dissimilarity is more (Good clustering)