

Clustering Validation & Assessment

10/2

3 types of tasks:

- evaluation: assessing the quality of the clustering
- stability: understand sensitivity of clustering
- tendency: suitability for clustering (generally hard)

not going
to cover
much

Types of measures:

- external: use criterion not in the dataset
e.g. labels provided by an expert
- internal: use criterion from the data
e.g. "compactness" of cluster
- relative: compare algs
e.g. 1 param set vs another

External measures

Assume correctness is already known

- labeled data
- synthetic data



Let $D = \{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n\}$ $\vec{x}_i \in \mathbb{R}^d$

\exists points in D are partitioned into k clusters

$y_i \in \{1, \dots, k\}$ is the ^{label} ground truth for point \vec{x}_i

$T = \{T_1, T_2, \dots, T_k\}$ are ground truth partitioning

e.g. $T_j = \{\vec{x}_i \in D \mid y_i = j\}$ each T_j is a partition

$$m_j = |T_j|$$

And $C = \{C_1, C_2, \dots, C_r\}$ be a clustering produced by some algo

C_i is a cluster

from clustering

$\hat{y}_i \in \{1, \dots, r\}$ be the label for \vec{x}_i

$$n_i = |C_i|$$

Contingency table $r \times k$ matrix

1	2	...	R
1			
2			
3			
4			
5			

w/ i, j entry containing # of points in C_i and T_j
 $N(i, j) = n_{ij} = |C_i \cap T_j|$

truth

$$N(1,2) = |C_1 \cap T_2|$$

How long to compute contingency table?

$\Theta(kr)$ for int to 0s
 $\Theta(n)$ for incrementing

Matching based measures

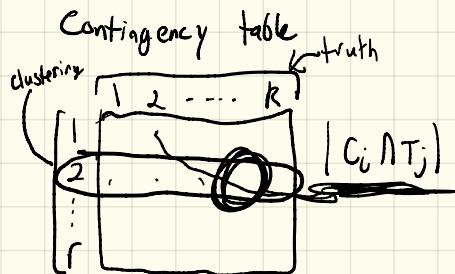
Purity - measure agreement w/ ground truth
value between $[0, 1]$ w/ 1 being all points are in 1 partition

$$\text{purity}_i = \frac{1}{n_i} \max_{j=1}^k \{ n_{ij} \}$$

$$\text{purity} = \sum_{i=1}^r \frac{n_i}{n} \text{purity}_i$$

$$= \frac{1}{n} \sum_{i=1}^r \max_{j=1}^k \{ n_{ij} \}$$

$[0, 1]$ w/ 1 \Rightarrow each cluster
contains points from only 1 partition



Internal measures

No labels so use intrinsic props:-

- Similarity
- Compactness
- Separation

Given a pairwise dist matrix

$$W = \{ \| \vec{x}_i - \vec{x}_j \| \}_{i,j=1}^n$$

Note 1: $\| \vec{x}_i - \vec{x}_j \| = \| \vec{x}_j - \vec{x}_i \|$ } upper Δ w/ diag
 $\| \vec{x}_i - \vec{x}_i \| = 0$

Note 2: W can also be thought of as
a complete undirected graph
w/ edge weights as dist between pts

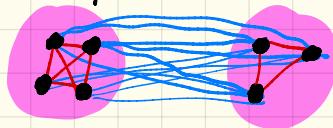
Consider as a k-way cut

Let $V = \bigcup_i C_i$

Given subsets $S, R \subset V$

define $w(S, R) = \sum_{x_i \in S} \sum_{x_j \in R} w_{ij}$

E.g. pick S & R to be clusters



- edges between clusters
intercluster edges
 N_{out} # of intercluster edges

Intracluster weight

$$w_{in} = \frac{1}{2} \sum_{l=1}^k w(C_l, C_l)$$

- edges w/in cluster
intracluster edges
 N_{in} # of intracluster edges

Intercluster weight

Let $\bar{S} = V \setminus S$

$$w_{\text{out}} = \frac{1}{2} \sum_{i=1}^k w(C_i, \bar{C}_i)$$