Representative_based clustering
   ⟶ k means
   ⟶ EM

Hierarchical

Density_based
   ⟶ kernel density estimation

graph clustering
   ⟶ Spectral
   ⟶ Markov chain

---

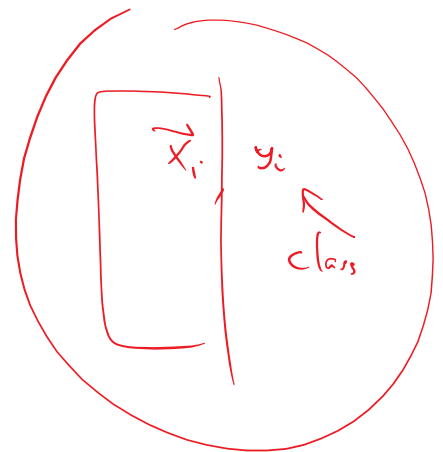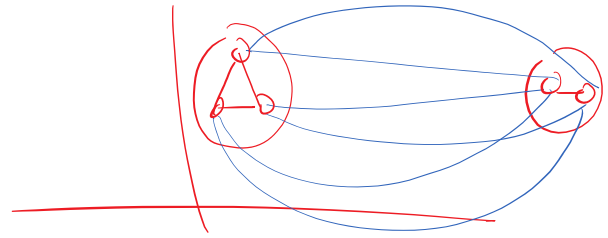Validation / Evaluation of clustering results ?
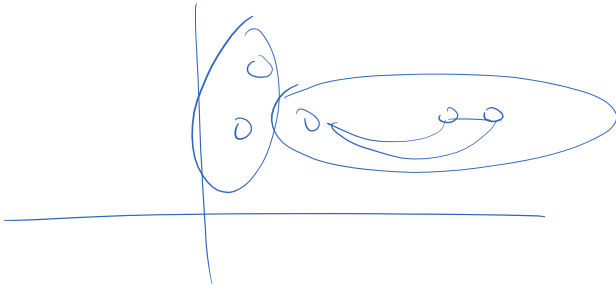
  ① External Measure

    where we have ground-truth.
     ( any classification dataset )

  ② Internal Measures

    ⟶ pair-wise distances
      +
    cluster output

$\vec{X_i}$ | $y_i$
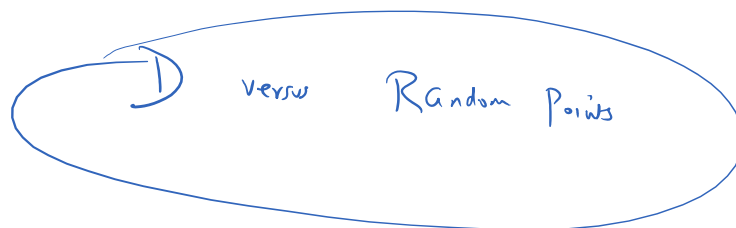
class

③ Relative Measures

Selecting $\underline{k}$  # of clusters

Cluster stability

→ how stable/ robust are the clusters under
small perturbations

Clustering tendency

→ Are there clusters?



versus   Random Points

---

External

$$T = \{T_1, T_2, \ldots, T_k\} \longleftarrow \text{ground truth}$$

↑
True cluster 1

---

$\text{from some algo} = \{c_1, c_2, \ldots, c_k\}$

K b given

algo

extracted clusters

## Contingency table / Confusion matrix

$T_1 \quad T_2 \quad \dots \quad T_k$

|  | $T_1$ | $T_2$ | $\dots$ |  |  |
|---|---|---|---|---|---|
| $C_1$ | $n_{11}$ | $n_{12}$ | $\dots$ |  | $n_1$ |
| $C_2$ |  |  |  |  | $n_2$ |
| $\vdots$ |  |  |  |  | $\vdots$ |
| $C_k$ |  |  |  |  | $n_k$ |
|  | $m_1$ | $m_2$ |  | $m_k$ |  |

$$n_{ij} = |C_i \cap T_j|$$

$$n_i = |C_i|$$

$$m_j = |T_j|$$

① __Purity__ for each cluster $C_i$

$$\text{Purity}_i = \frac{1}{n_i} \left( \max_{j=1}^{k} \{ n_{ij} \} \right)$$

$$\text{Purity} = \sum_{i=1}^{k} \frac{n_i}{n}$$

weighted sum

$T_1 \quad T_2 \quad T_3$

Points

$C_1 \rightarrow T_2$

$C_2 \rightarrow T_3$

$C_3 \rightarrow T$

$\text{Prec}_1 = \frac{90}{94}$

$\text{recall}_1 = \frac{90}{200}$

|  | $T_1$ | $T_2$ | $T_3$ |  |  |
|---|---|---|---|---|---|
| $C_1$ | 2 | 90 | 2 |  | 94 |
| $C_2$ | 1 | 5 | 100 |  | 100 |

lumb

$c_2 \to$ '3

$c_3 \to T_2$

$recall_1 = \dfrac{T_1}{200}$

|       | $m_1$ | $m_2$ | $m_3$ |      |
|-------|-------|-------|-------|------|
|       | (...) | (...) | (...) | (14) |
| $c_2$ | 1     | 5     | (100) | 106  |
| $c_3$ | (99)  | (100) | 1     | 200  |
|       | 100   | 200   | 100   | (400)|

$\overline{\qquad\qquad}$  **Vs**

Maximum
matching

Measure

$c_1 \to T_2$
$c_2 \to T_3$
$c_3 \to T_1$

1-1 matching
with the
highest sum
of edge weights

2
$r_a$

---

$\overline{F}$ — measure:

$$F_i = \frac{2 \cdot Prec_i \cdot Recall_i}{Prec_i + Recall_i}$$

$$\mu_F = \frac{1}{k} \sum F_i$$

F-measure

$$Prec_i = Purity_i = \frac{1}{n_i} \max_j \{n_{ij}\}$$

$$j_i^* = arg\max_j \{n_{ij}\}$$

$$Recall_i = \frac{n_{ij^*}}{m_{j^*}}$$

---

$\gamma$: # of clusters
found

|       | $T_1$ | $T_2$ | $T_3$ |
|-------|-------|-------|-------|
| $c_1$ |       |       |       |

|       | $T_1$ | $T_2$ |
|-------|-------|-------|
|       | 0     | 50    |

$r$: # of clusters found

$r < k$

$c_1$

$c_2$

$r < k$

$r = 4$

| $r_1$ | $r_2$ |
|-----|-----|
| 0 | 50 |
| 50 | 0 |
| 10 | 40 |
| 30 | 20 |

$r > k$

$T_2$

$c_1$ $c_3$    $c_2$ $c_4$

## Information Theory

$$H = -\sum_{i=1}^{k} P_i \log P_i$$

### Conditional Entropy

$H(T \mid c_i)$  (row-wise entropy)

$$H(T \mid c_i) = -\sum^{k} \left(\frac{n_{ij}}{n_i}\right) \log \left(\frac{n_{ij}}{n_i}\right)$$

$$H(T \mid c) = \quad \frac{n_i}{n} H(T \mid c_i)$$

$\underbrace{\qquad}_{\text{weighted sum}}$

Perfect clustering $H(T \mid c) = 0$

$$\boxed{H(T \mid c) = \underbrace{H(T, c)}_{\text{joint entropy}} - \underbrace{H(c)}_{\text{entropy of clustering}}}$$

|       | $T_1$ | $T_2$ | $T_3$ |       |
|-------|-------|-------|-------|-------|
| $n_{11}$ 100 $n_{12}$ | | 0 $n_{13}$ | $n_1$ |
| 0 | 0 | 100 | $n_2$ |
| 100 | 0 | 0 | $n_3$ |
| 100 | 200 | 100 | |

| 0 | 1 | 0 |

|       |     |     |     |
|-------|-----|-----|-----|
| $c_1$ | 50  | 100 | 50  |
| $c_2$ | 40  | -   | -   |
| $c_3$ | 10  |     |     |

$\frac{1}{4}, \frac{1}{2}, \frac{1}{4}$

$$-\sum_i \sum_j P_{ij} \log P_{ij} \qquad\qquad -\sum_i P_{c_i} \log P_{c_i}$$

## Pairwise

$\boxed{TP:}$    how many points $x_i$ & $x_j$

belong to the same cluster $C_k$ &

they also belong to the same true

cluster $T_a$

$$\boxed{{}^n C_2 = \binom{n}{2} = \frac{n \cdot (n-1)}{2}}$$

all distinct pairs of
points

$$T_1 = \{1 \underbrace{2, 3}\} \quad T_2 = \{4, \underset{\neq}{5\,6}\}$$

$$C_1 = \{1, 4\}, \quad C_2 = \{2, 3, \underbrace{5\,6}\}$$

$C_1' = \{1, 2, 3\}$

$C_2' = \{4, 5, 6\}$

$\boxed{TP:} \quad \underline{1} + 1 = 2$

$T_N:$   # of pairs of points that

belong to diff clusters &

also diff true clusters

|       | $T$  | $\overline{T}$ |
|-------|------|------|
| $C$   | $TP$ | $FP$ |
| $\overline{C}$ | $FN$ | $TN$ |

$$\text{Jaccard}: \quad \frac{TP}{TP + FN + FP} \qquad\qquad \Big|\quad \text{Ignores } \underline{\underline{TN}}$$

FM: Folkes Mallow

geometric mean of Prec & Recall

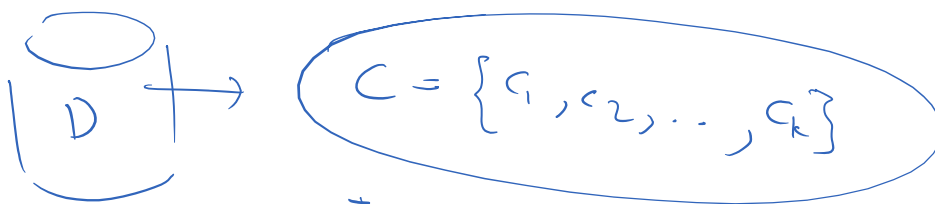$$TP$$

$$recall = \frac{TP}{TP+FN}$$

FM:

$$\sqrt{Prec. \cdot Recall}$$

$$: \sqrt{\frac{TP^2}{(TP+FP)(TP+FN)}}$$

---

## No ground truth !!!

$$D \longrightarrow C = \{c_1, c_2, \ldots, c_k\}$$

$$W = \{w_{ij}\}$$

Pair-wise distance matrix

$$w_{ij} = \|\vec{x_i} - \vec{x_j}\|^2$$
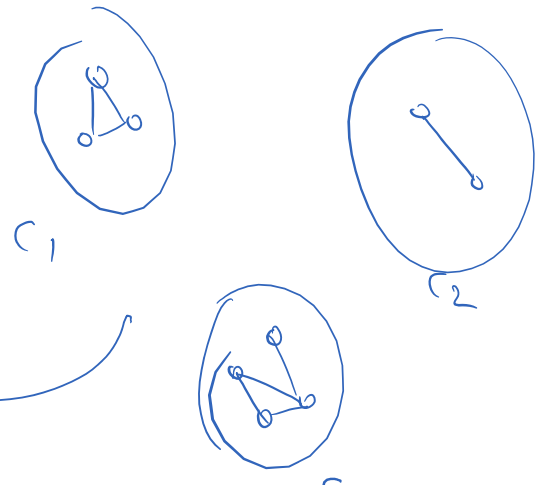
for graph data

$w_{ij} = $ # of hops in the shortest path

$$W(S,T) = \sum_{x_a \in S} \sum_{x_b \in T} w_{ab}$$

weight of the cut

$$W_{in} = \frac{1}{2} \sum^{k} W(C_i, C_i)$$

Internal weights

$C_1$

$C_2$

$C$

N.

$N_{in}$ = how many pairs of internal points

$$\sum_{i=1}$$

$$n_i = |c_i|$$

$$W_{out} : \frac{1}{2} \sum_{i=1}^{k} W(c_i, \overline{c_i})$$

external
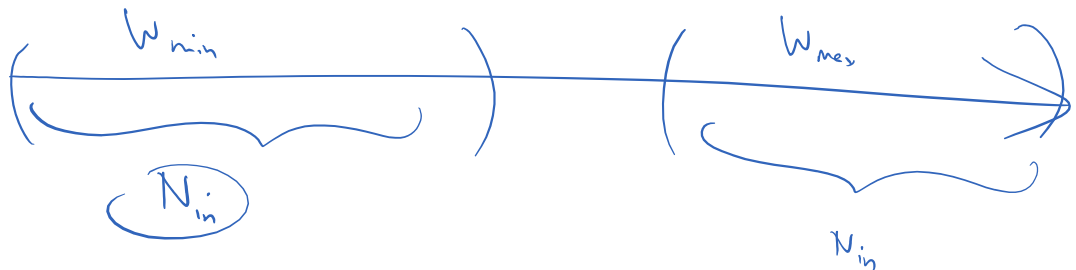weight

$$N_{out} : n - N_{in}$$

$\#$ external
pairs of points

① $\underline{Beta\ CV}$ : $\dfrac{W_{in}/N_{in}}{W_{out}/N_{out}}$     smaller is better!

② $\underline{C\text{-}Index}$     $\underline{\underline{N_{in}}}$ : $\#$ of internal pairs

$$W_{min}(N_{in}) = \quad N_{in} \quad \text{smallest weights in the } \overline{W} \text{ matrix}$$



$W_{min}$          $W_{max}$

$N_{in}$

$N_{in}$

$$C_{index} = \frac{\left(W_{in}\right) - W_{min}\left(N_{in}\right)}{W_{max}\left(N_{in}\right) - W_{min}\left(N_{in}\right)}$$

Smaller the better

# Silhouette Coefficient

$$\underset{\text{Per point } \vec{x}_i}{S_i} = \frac{\mu^{out}(x_i) - \mu^{in}(x_i)}{\max\{\mu^{out}, \mu^{in}\}}$$

$S_i$ Close to $1$ is best!

$S_i \in [-1, 1]$

$S_i = -1$

$\Rightarrow$ Mis-clustered point

$C_2$

$C_3$

$x_i$

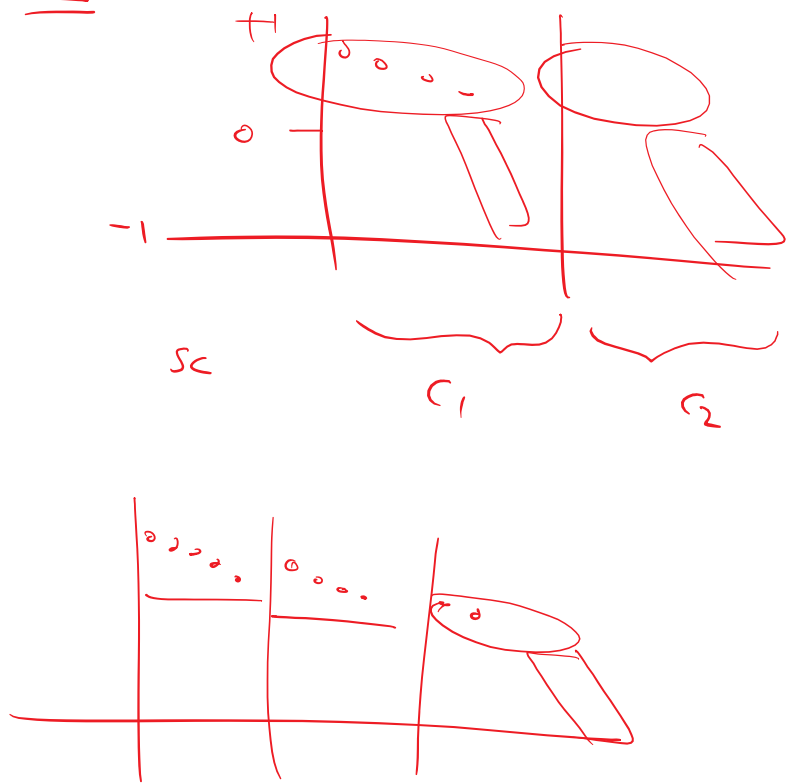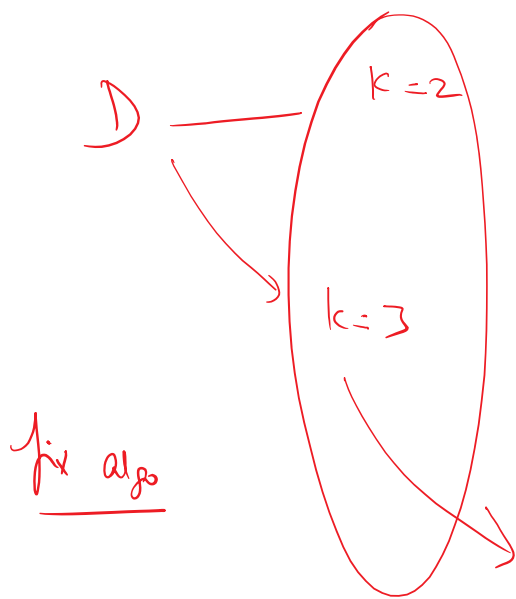$$\mu^{in}(x_i) = \frac{1}{n_k} \sum_{x_j \in C_k} \|\vec{x}_i - \vec{x}_j\|^2$$

$$\mu^{out}(x_i) = \underset{C_k}{\min} \left\{ \text{avg distance to cluster } C_k \right\}$$

$$SC: \frac{1}{n} \sum S_i$$

avg over all points

# Relative Measure to select $k$

$D$ — $\Big($ $k=2$

$k=3$ $\Big)$

$\underline{\text{fix algo}}$



$SC$      $C_1$      $C_2$

---

# gap statistic to select $k$

$$C = \{c_1, c_2 \dots c_k\}$$

$|D| \longrightarrow W_{in}$ : total internal weight

$\uparrow$

is this good or bad?

generate random samples in the same data space as $D$

$q$

random

$\Bigg\{ \begin{array}{l} R_1 \longrightarrow W_{in}^k(R_1) \\ R_2 \longrightarrow W_{in}^k(R_2) \\ \vdots \end{array} \Bigg\}$
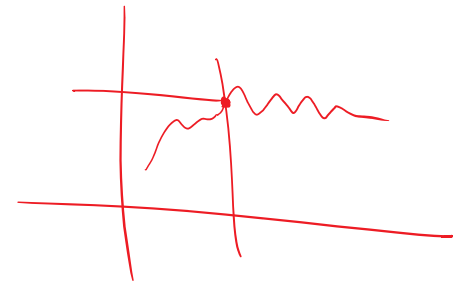
$\mu^k \longleftarrow$ mean internal

$\underline{\quad}^k$

$V$

$V$ random samples

$\left\{ \begin{array}{ccc} R_2 & \longrightarrow & Win(R_2) \\ \vdots & & \vdots \\ R_q & \longrightarrow & Win(R_q) \end{array} \right\}$

$\mu \longleftarrow$ mean internal

$\sigma^k \longleftarrow$ Stdev.

$$gap(k) = \frac{\mu^k}{\underbrace{\phantom{xxx}}_{\text{from } q \text{ random samples}}} - \frac{Win^k_{in}}{\underbrace{\phantom{xxx}}_{\text{from } D}}$$

gap(k)

k

① Pick the max k

② Choose smallest $k$ such that

$$\left[ gap(k) \geq gap(k+1) - \sigma^{k+1} \right]$$

within a deviation of the next value of k