# EVALUATING CLUSTERING METHODS

(1) Internal clustering validation:

- No ground-truth available.
- How good is the cluster structure?

(2) External cluster evaluation:

- Take a benchmark with known labels
- How "close" are the clusters to the true labels?

# INTERNAL EVALUATION

$C_1, C_2, \underline{\quad}, C_u$ = output clusters

## DUNN INDEX

(1) Inter-cluster distance

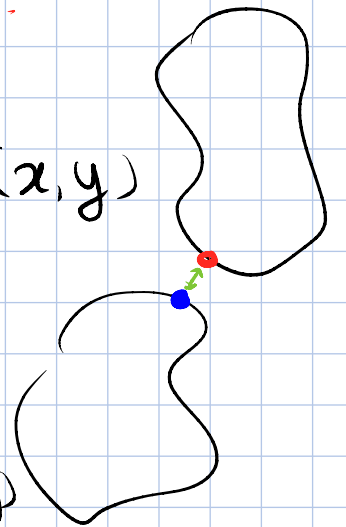= how well-separated the clusters are.

### Examples:

- $\delta(C_i, C_j) = \min_{\substack{x \in C_i \\ y \in C_j}} d(x, y)$

  Single-linkage distance

  distance

- $\delta(C_i, C_j) = \frac{1}{|C_i||C_j|} \sum_{\substack{x \in C_i \\ y \in C_j}} d(x, y)$
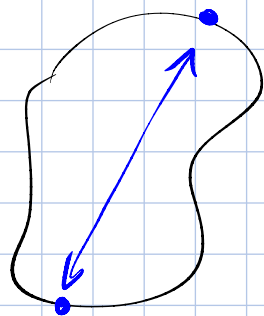
  average linkage distance.

(2) Intra-cluster distance =
　　　　how "concentrated"

Example:

- $\Delta(c_i) = \max\limits_{x,x' \in c_i} d(x,x')$

diameter

- $\Delta(c_i) = \dfrac{1}{|c_i|(|c_i|-1)} \sum\limits_{x,x' \in c_i} d(x,x')$ → average distance.

The Dunn Index:

$$DI = \dfrac{\min\limits_{1 \le i < j \le k} \delta(c_i, c_j)}{\max\limits_{1 \le \ell \le k} \Delta(c_\ell)}$$

worst case

worst case.
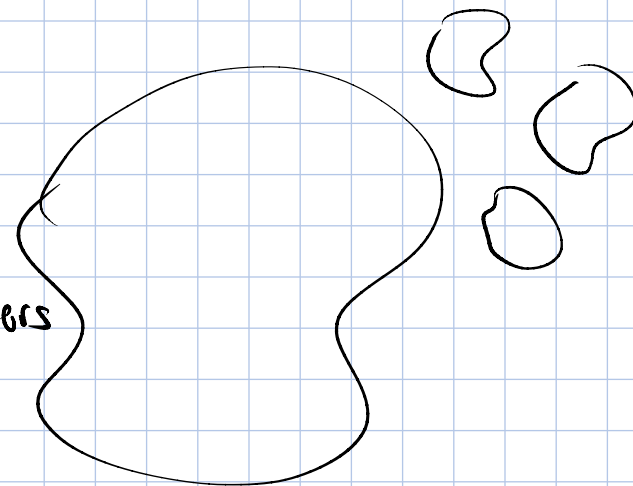
high = good separation

high as possible.

low = good concentration

- $DI \in [0, \infty)$

- High $DI \Rightarrow$ better clustering.

- Issue:

When we have a mix of small & big clusters
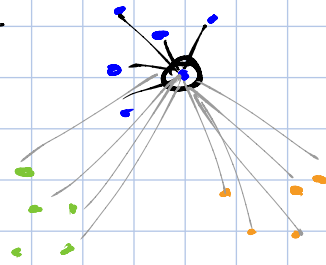
↓

$DI \to \underline{low}$

(even if method works well)

# SILHOUETTE ANALYSIS

- Clusters: $C_1, C_2, \ldots, C_k$.

- Pick $x \in C_i$:

$$a(x) = \frac{1}{|C_i|-1} \sum_{x' \in C_i} d(x, x') \rightarrow \text{within-cluster distance.}$$

$$b(x) = \min_{j \neq i} \frac{1}{|C_j|} \sum_{y \in C_j} d(x, y) \rightarrow \text{between-cluster distance.}$$

## Silhouette Coefficient:

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

If $x \in C_i$
and $|C_i| = 1$
$\downarrow$
$s(x) = 0$.

## Comments:

- $a(x) < b(x) \Rightarrow s(x) = 1 - \frac{a(x)}{b(x)} \in (0, 1]$

- $a(x) > b(x) \Rightarrow s(x) = \frac{b(x)}{a(x)} - 1 \in [-1, 0)$

- $a(x) = b(x) \Rightarrow s(x) = 0$.

$$\rightarrow \quad -1 \leq s(x) \leq 1$$

### Good clustering:
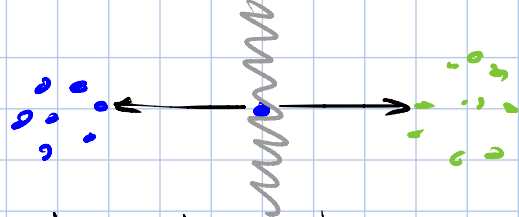
$$a(x) << b(x) \rightarrow s(x) \approx +1$$

### Bad clustering:

$$b(x) << a(x) \rightarrow s(x) \approx -1$$

⊗ $S(x) = 0$

↓

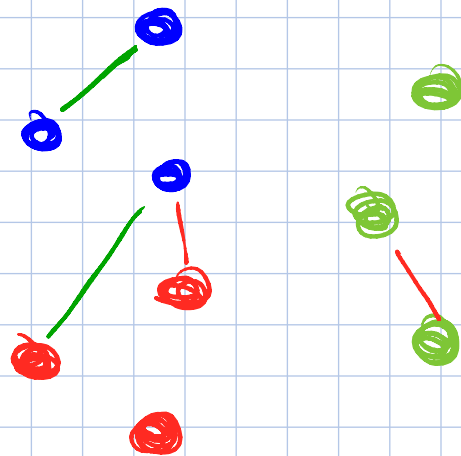$x$ is on the decision boundary.

## MEAN SILHOUETTE COEFFICIENT:

$$SI = \frac{1}{n} \sum_x S(x) \rightarrow$$ overall quality of algorithm.
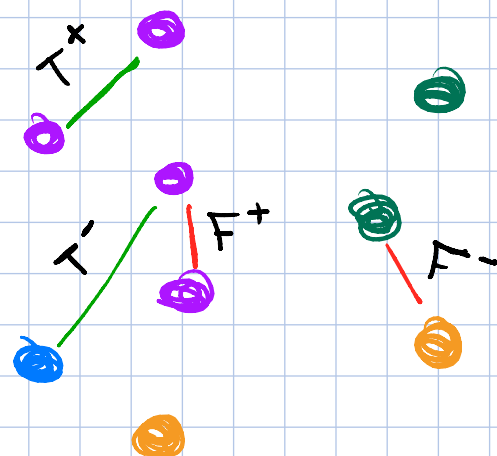
Compare results to some benchmark.

- labelled data.

- data with concensus on the "truth"

## RAND INDEX



"truth"                    our algorithm

Data set : $\{x_1, x_2, \_\_, x_n\}$

Ground truth: $C^* = \{c_1^*, c_2^*, \_\_, c_m^*\}$

Our algorithm: $C = \{c_1, c_2, \_\_, c_k\}$

Denote:

- $x_i \sim x_j$  if $x_i, x_j$ in the same cluster in $C$

- $x_i \overset{*}{\sim} x_j$  if $x_i, x_j$ in the same cluster in $C^*$

Define:

$$(x_i, x_j) \longrightarrow \begin{cases} T^+ & \text{if } x_i \sim x_j \vee x_i \overset{*}{\sim} x_j \quad \overset{\text{and}}{} \\ T^- & \text{if } x_i \nsim x_j \text{ and } x_i \overset{*}{\nsim} x_j \\ F^+ & \text{if } x_i \sim x_j \text{ and } x_i \overset{*}{\nsim} x_j \\ F^- & \text{if } x_i \nsim x_j \text{ and } x_i \overset{*}{\sim} x_j \end{cases}$$

$TP =$ total number of $T^+$'s

$TN = \quad \text{''} \quad \text{''} \quad \text{''} \quad T^-$'s

$FP = \quad \text{''} \quad \text{''} \quad \text{''} \quad F^+$'s

$FN = \quad \text{''} \quad \text{''} \quad \text{''} \quad F^-$'s

RAND INDEX:

$$RI = \frac{TP + TN}{\underbrace{TP + TN + FP + FN}_{\binom{n}{2} = \frac{n \cdot (n-1)}{2}}} \quad \in [0, 1].$$

$RI \approx 1 - \underline{good}$ clustering