

Dec-28 19/03/21.

Introduction to clustering

✓ Basic concept

✓ distance measures

→ K-means clustering algo.

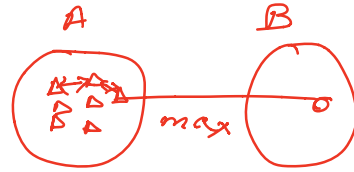
defn :

The process of grouping a set of physical or abstract-objects into classes of similar objects is called grouping.

→ clustering is also called data segmentation because clustering partitions large data set into groups based on their similarity.

→ clustering can be used for outlier detection.

→ In machine learning, clustering is an example of unsupervised learning.



Distance Measures :

(a) Geometric distance measure
(dealing with numeric data)

(b) Percent disagreement -
(dealing with categorical data)

a) Geometric distance measure :

↳ Minkowski distance :

$$d(i,j) = \left(|x_{i1} - x_{j1}|^p + |x_{i2} - x_{j2}|^p + \dots + |x_{in} - x_{jn}|^p \right)^{1/p}$$

where p is a positive integer. Such a distance is called L_p norm.

2) when $p=1$ (i.e. L_1 norm) \rightarrow It represents
Manhattan distance

$$d(i,j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{in} - x_{jn}|$$

3) when $p=2$ (i.e. L_2 norm) \rightarrow It represents
Euclidean distance

$$d(i,j) = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{in} - x_{jn})^2}$$

4) Chebychev distance :

$$dis(A,B) = \max \{ |x_i - y_i| \}$$

5) Power distance :

$$\text{Power dis}(A,B) = \left(\sum_i |x_i - y_i|^p \right)^{1/q}$$

where p and q are parameters defined by user.

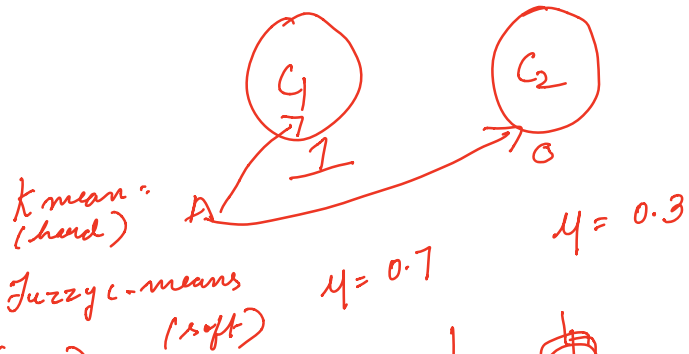
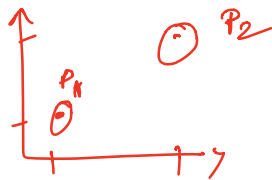
(b) Percent disagreement :

$$\text{Percent disagreement} = \left(\frac{\text{no. of } x_i \neq y_i}{i} \right)$$

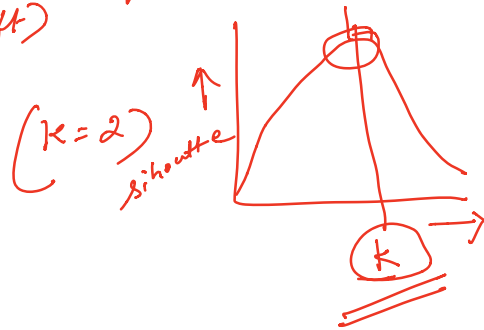
$$\begin{array}{l} G_1: \quad A \begin{array}{|c|} \hline x \\ \hline C \\ \hline \end{array} T \begin{array}{|c|} \hline x \\ \hline G \\ \hline \end{array} T \begin{array}{|c|} \hline x \\ \hline G \\ \hline \end{array} \begin{array}{|c|} \hline x \\ \hline G \\ \hline \end{array} G \begin{array}{|c|} \hline C \\ \hline A \\ \hline \end{array} \\ G_2: \quad A \begin{array}{|c|} \hline T \\ \hline \end{array} T \begin{array}{|c|} \hline C \\ \hline \end{array} T \begin{array}{|c|} \hline A \\ \hline \end{array} \begin{array}{|c|} \hline A \\ \hline \end{array} G \begin{array}{|c|} \hline C \\ \hline A \\ \hline \end{array} \end{array}$$

$$\text{Percent disagreement} = \left(\frac{4}{10} \right) \times 100\% \quad \begin{array}{l} x \Rightarrow 0 \\ = 40\% \end{array}$$

✓ k-means clustering (Partitioning method)



$P_1 (1, 1)$	$P_5 (5, 2)$
$P_2 (6, 7)$	$P_6 (2, 3)$
$P_3 (4, 6)$	$P_7 (1, 2)$
$P_4 (5, 7)$	$P_8 (3, 1)$



Itern 1 :
Randomly choose 2 points as cluster centroid .
Here let $P_5 (5, 2)$ and $P_7 (1, 2)$

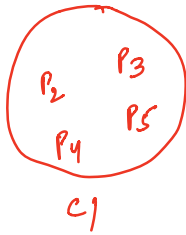
$P_1 \begin{cases} P_5 = 4.123 \\ P_7 = 1 \end{cases} \checkmark$
 $P_2 \begin{cases} P_5 = 5.099 \checkmark \\ P_7 = 7.071 \end{cases}$
 $P_3 \begin{cases} P_5 = 4.123 \checkmark \\ P_7 = 5 \end{cases}$
 $P_4 \begin{cases} P_5 = 5 \checkmark \\ P_7 = 6.403 \end{cases}$

$P_5 \begin{cases} P_5 = 0 \checkmark \\ P_7 = 4 \end{cases}$
 $P_6 \begin{cases} P_5 = 3.162 \\ P_7 = 1.414 \checkmark \end{cases}$
 $P_7 \begin{cases} P_5 = 4.0 \\ P_7 = 0 \checkmark \end{cases}$
 $P_8 \begin{cases} P_5 = 2.236 \\ P_7 = 2.236 \end{cases}$

arbitrary break the tie by assigning P_8 to P_7 cluster

$C_1 (P_5)$	$C_2 (P_7)$
$\{P_2, P_3, P_4, P_5\}$	$\{P_1, P_6, P_7, P_8\}$

Item 2:



$$C_1' = x' = \frac{5+5+6+4}{4} = 5$$

$$y' = \frac{7+6+7+2}{4} = 5.5$$

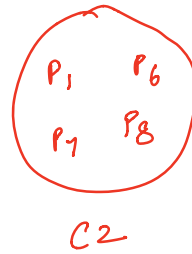
$$P_1 \begin{cases} \rightarrow C_2' = 1.06 \checkmark \\ \rightarrow C_1' = 6.02 \end{cases}$$

$$P_2 \begin{cases} \rightarrow C_2' = 6.75 \\ \rightarrow C_1' = 1.80 \checkmark \end{cases}$$

$$P_3 \begin{cases} \rightarrow C_2' = 4.808 \\ \rightarrow C_1' = 1.118 \checkmark \end{cases}$$

$$P_4 \begin{cases} \rightarrow C_2' = 6.174 \\ \rightarrow C_1' = 1.5 \checkmark \end{cases}$$

$$C_1' = \{ P_2 \ P_3 \ P_4 \}$$



$$C_2' = x' = 1.75$$

$$y' = 1.75$$

$$P_5 \begin{cases} \rightarrow C_1' = 3.5 \\ \rightarrow C_2' = 3.259 \checkmark \end{cases}$$

$$P_6 \begin{cases} \rightarrow C_1' = 3.905 \\ \rightarrow C_2' = 1.274 \checkmark \end{cases}$$

$$P_7 \begin{cases} \rightarrow C_1' = 5.315 \\ \rightarrow C_2' = 0.790 \checkmark \end{cases}$$

$$P_8 \begin{cases} \rightarrow C_1' = 4.924 \\ \rightarrow C_2' = 1.457 \checkmark \end{cases}$$

$$C_2' = \{ P_1 \ P_5 \ P_6 \ P_7 \ P_8 \}$$

Item 3

$$C_1'' \quad x = \frac{5+4+6}{3} = 5$$

$$y = \frac{7+6+7}{3} = 6.66$$

$$C_2'' \quad x = \frac{1+1+2+3+5}{5} = 2.4$$

$$y = \frac{1+2+3+1+2}{5} = 1.8$$

o/p of 3rd iters :

$$C_1'' = \{ P_2 P_3 P_4 \}$$

$$C_2'' = \{ P_1 P_5 P_6 P_7 P_8 \}$$

The points inside cluster don't change in two successive iterations. Thus, K-means is converged.

Exit