

Introduction

k-means clustering is a method of cluster analysis which aims to partition n observations into k clusters in which each observation belongs to the cluster with the nearest mean.

□ Notations:

- Patterns: $x_1, x_2, \dots, x_n \in \mathbb{R}^m$
- Cluster Centers: $c_1, c_2, \dots, c_k \in \mathbb{R}^m$
- Euclidean distance: $\|x_i - c_j\|^2$

MacQueen, J. B. (1967). "Some Methods for classification and Analysis of Multivariate Observations", Proceedings of 5th Berkeley Symposium on Mathematical Statistics and Probability. University of California Press. pp. 281–297.

Partition

Based on minimum within cluster measurement.

Let there are $A_1, A_2, A_3, \dots, A_c$ Clusters/partition

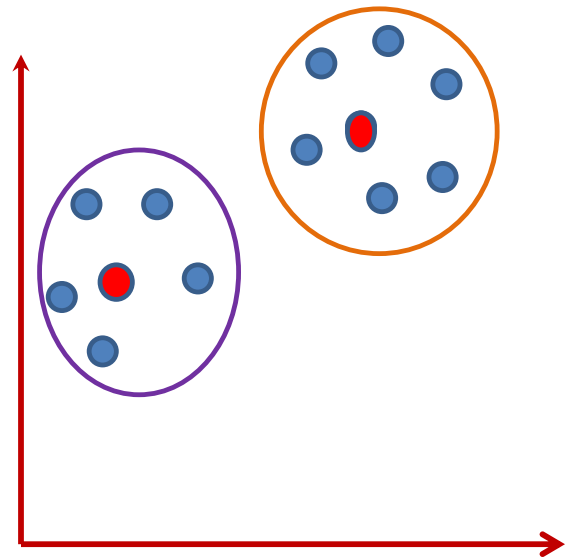
Then they should satisfies these constraints

- $A_i \cap A_j = \emptyset \quad \forall i \neq j$
- $\bigcup_1^c A_i = S \quad s = \{x_1, x_2, x_3, \dots, x_n\}$
- $A_i \neq \emptyset$

Algorithm

There are 4 steps in the algorithm.

- Randomly select k centers from the given data.
- Assign each object to the cluster with the nearest center point. Compute the centers of the clusters of the current partition (i.e., *mean point*, of the cluster).
- Go back to Step 2, stop when no more new assignment.



New Centers will be calculated by

$$\text{Center1} = (x_1 + x_2 + x_3 + x_4 + x_5) / 5$$

$$\text{Center2} = (x_6 + x_7 + x_8 + x_9 + x_{10} + x_{11}) / 6$$

Algorithm

We have given $S = \{x_1, x_2, x_3, \dots, x_n\} \in \mathbb{R}^m$

Step1: Choose c points $y_{11}, y_{12}, y_{13}, \dots, y_{1c} \in \mathbb{R}^m$

Step2: $A_{2i} = \{x \in S: d(\underline{x}, y_{1i}) \leq d(\underline{x}, y_{1j}) \ \forall j \neq i\}$

Step3: $y_{2i} = \text{Mean}(A_{2i}), i=1, 2, 3, \dots, c$

Step4: if $\|y_{1i} - y_{2i}\| \leq \varepsilon$

then

STOP with the output A_{2i}

else

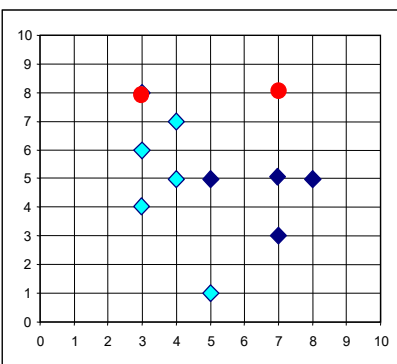
$A_{1i} = A_{2i}$

$A_{2i} = \emptyset$

$y_{1i} = y_{2i}$

GO TO Step2

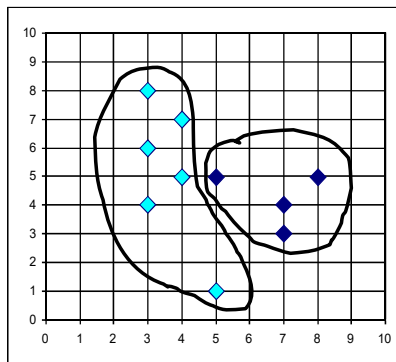
Example-1



K=2

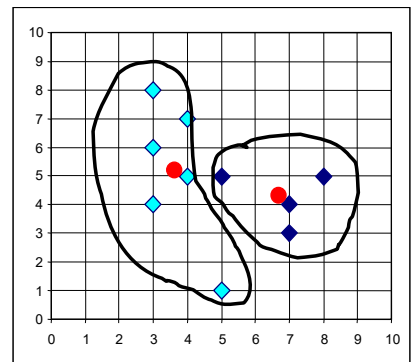
Arbitrarily choose K
object as initial
cluster center

Assign each objects to most similar center



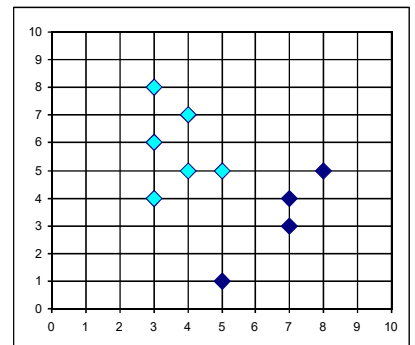
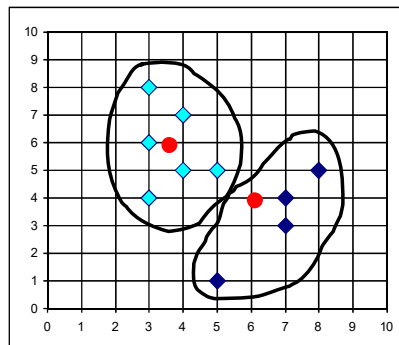
reassign

Update
the
cluster
means



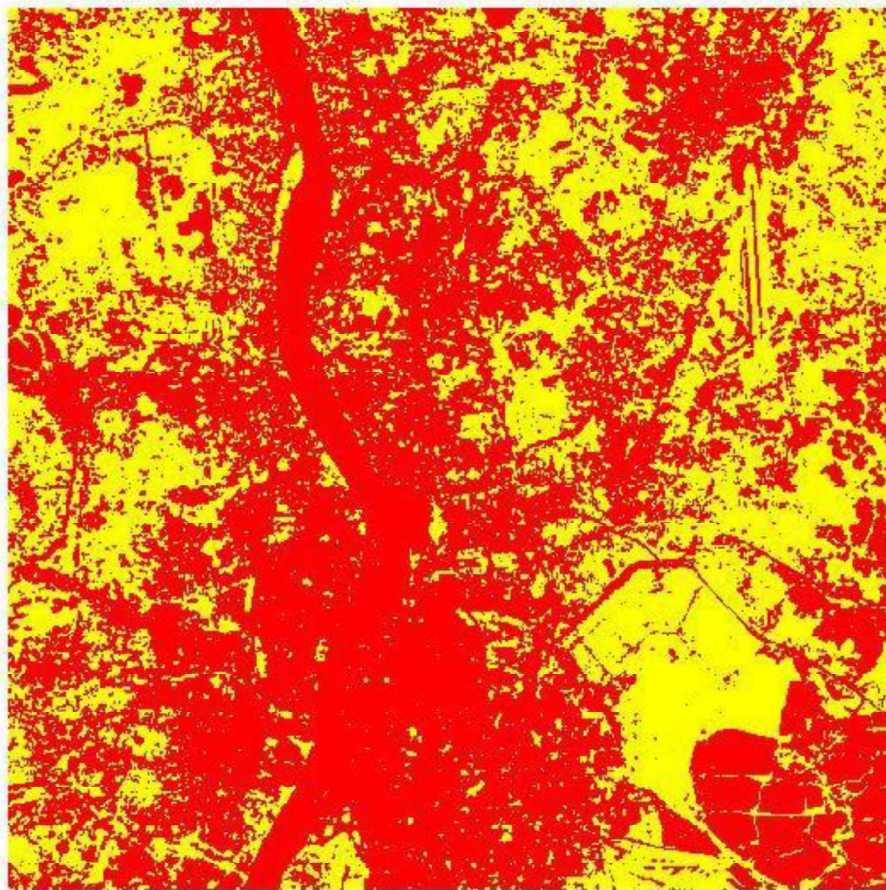
reassign

Update
the
cluster
means



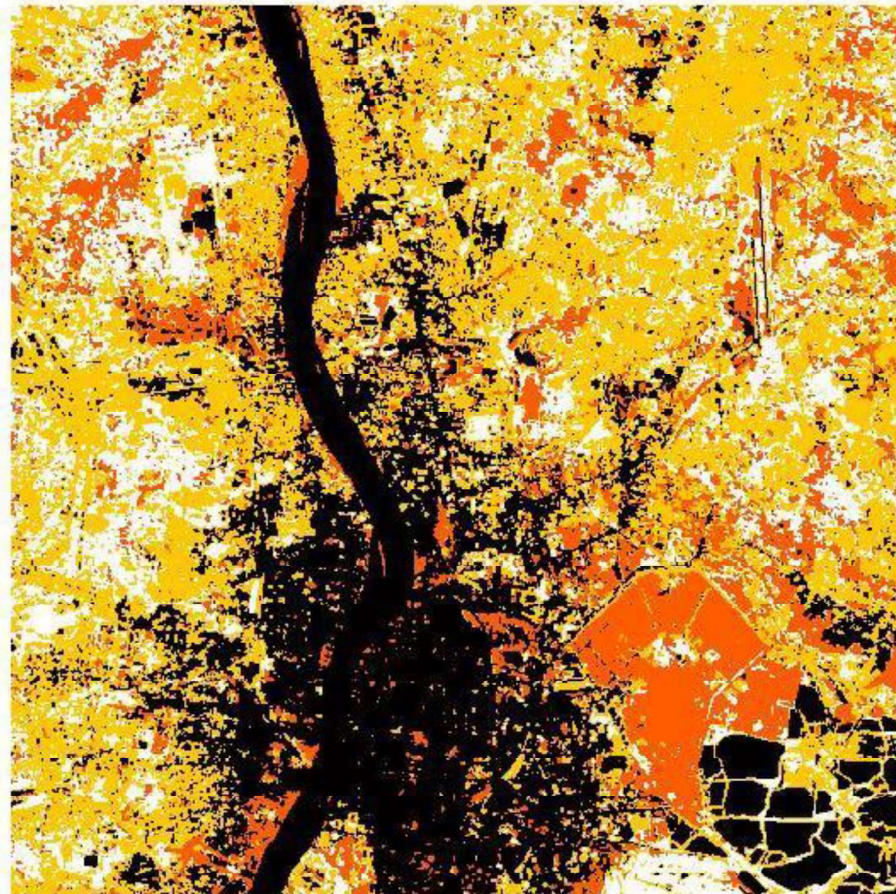
Example-2

C-Mean Classifier When $C=2$



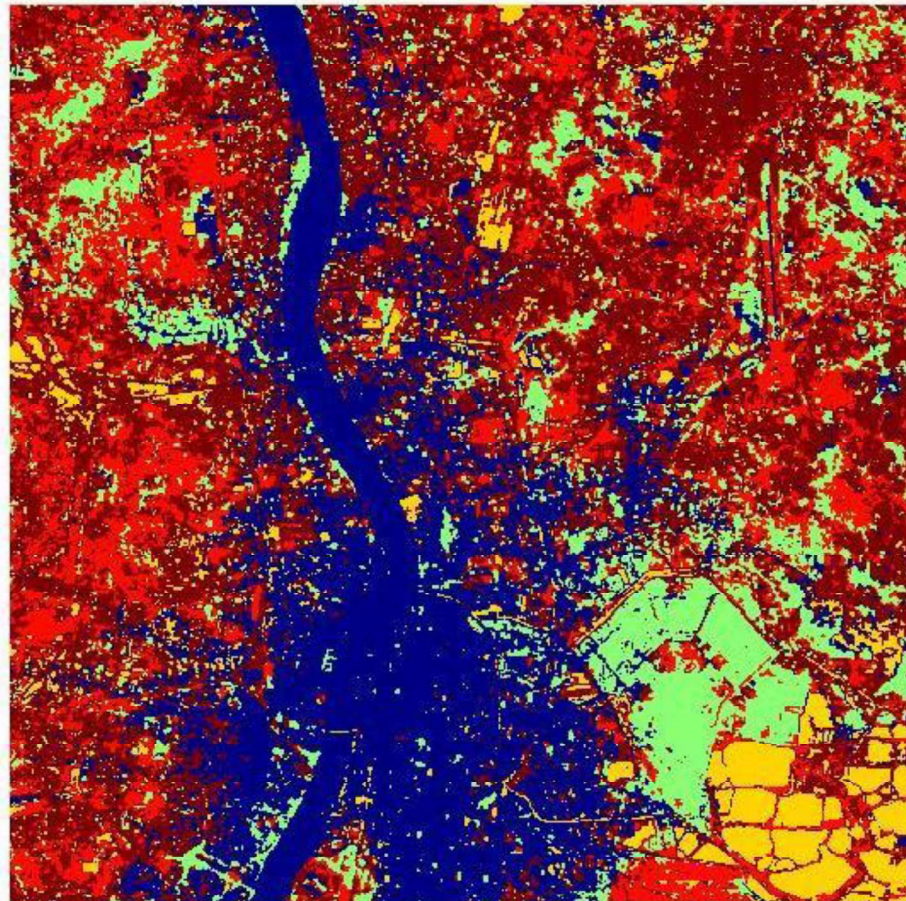
Example-2

C-Mean Classifier When $C=4$



Example-2

C-Mean Classifier When $C=5$



Remarks

- ✓ Algorithm usually converges
- ✓ Two different sets of initial seed point may sometimes give rise to two different final clustering.
- ✓ The algorithm tries to implement the Minimum Within Cluster Distance Criteria.
- ✓ We can start with the initial partition of the dataset (instead of the seed point)

Thanks