

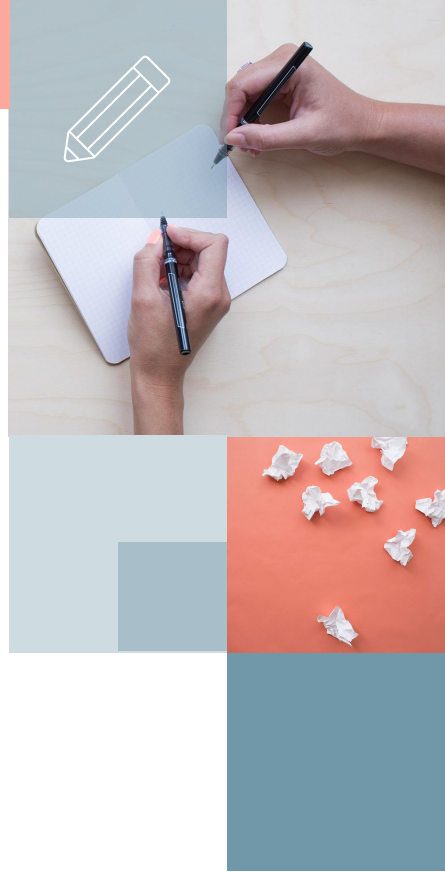


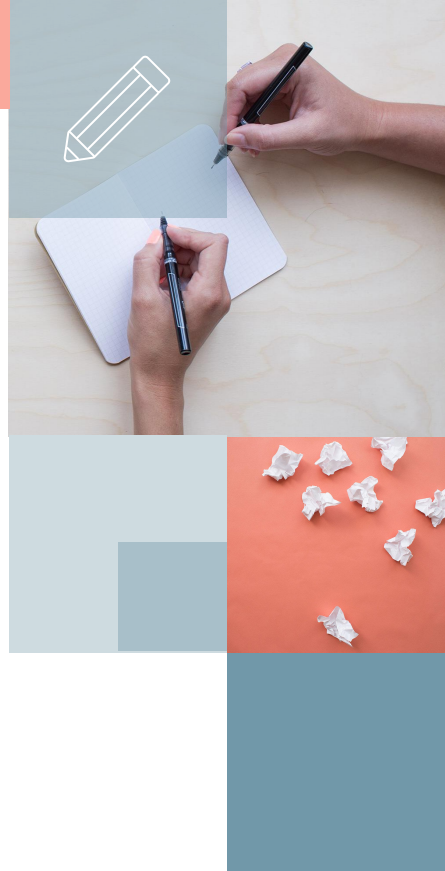
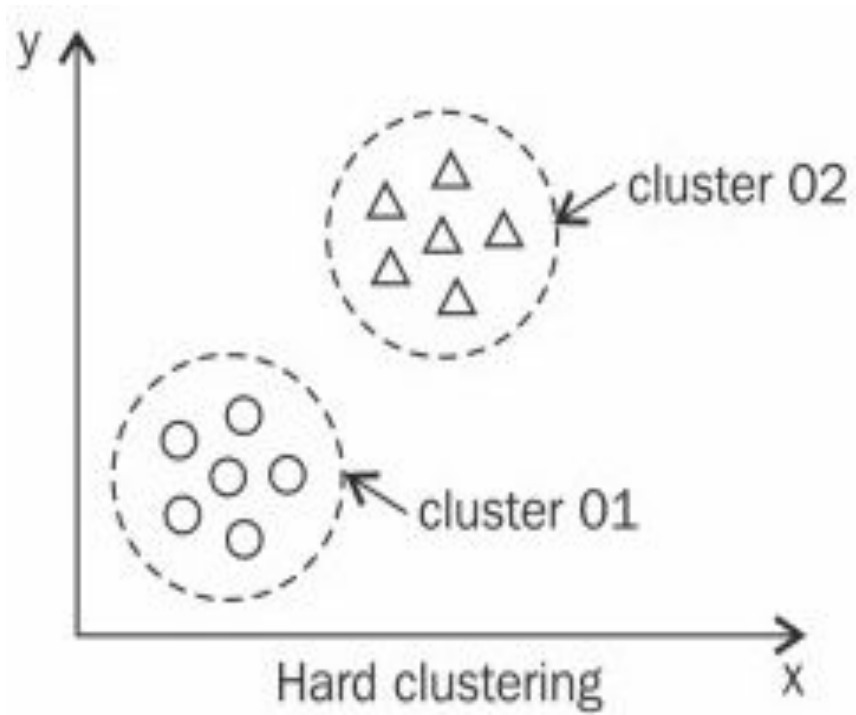
# K-Means

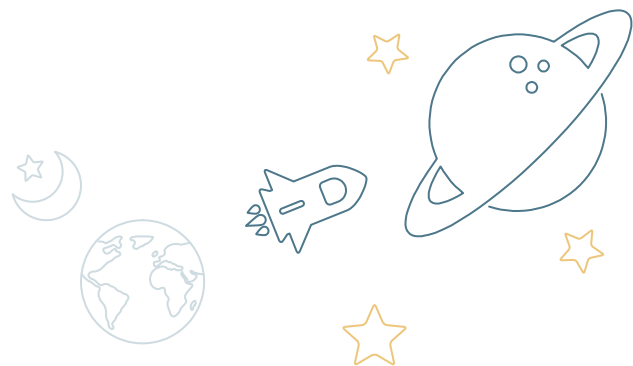
## Fuzzy C-Means

# Hard Clustering

Hard clustering is about grouping the data items such that each item is only assigned to one cluster. For example, we want the algorithm to read all of the tweets and determine if a tweet is a positive or a negative tweet.







# K means

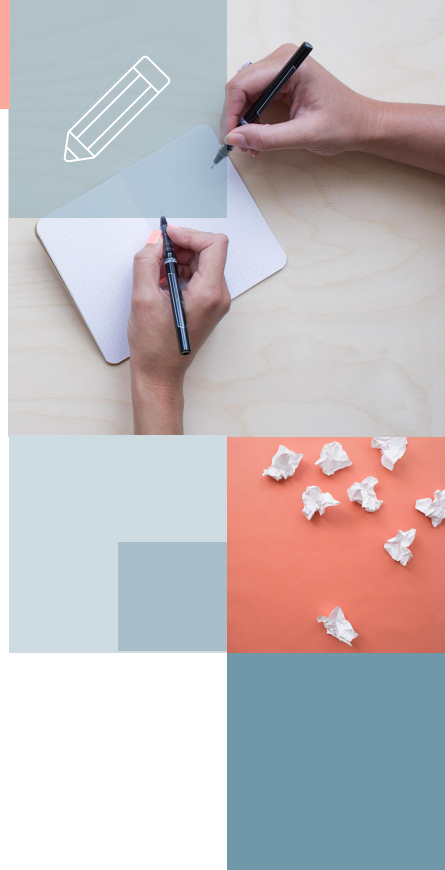
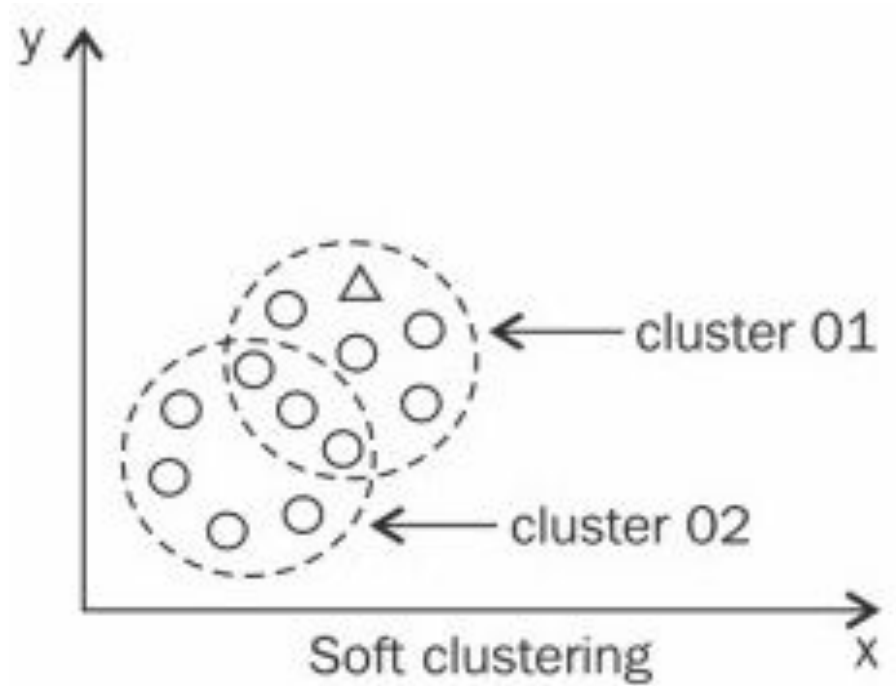


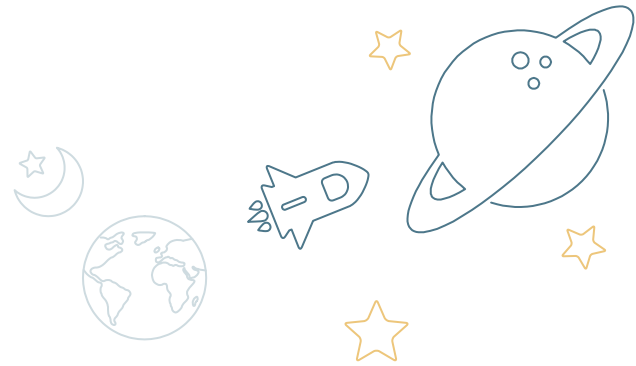
1. Specify the desired number of clusters
2. Randomly assign each data point to a cluster
3. Compute cluster centroids
4. Re-assign each point to the closet cluster centroid
5. Re-compute cluster centroids
6. We repeat steps 4 and 5 until no further improvement is found in the centroid values

# Soft Clustering

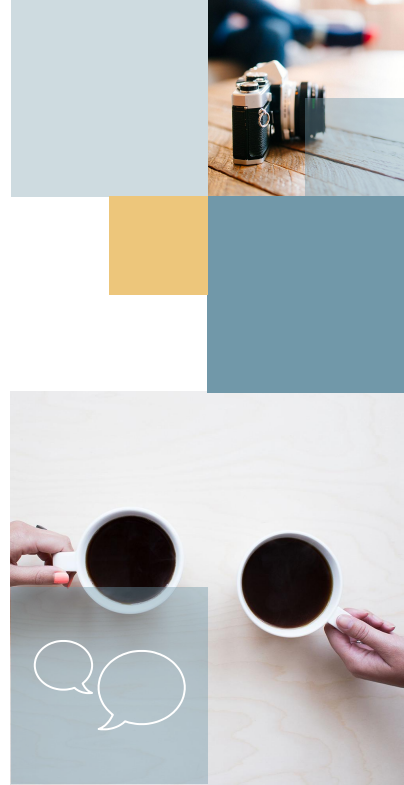
Fuzzy clustering (also referred to as soft clustering or soft k-means) is a form of clustering in which each data point can belong to more than one cluster. For example, preferences of customers in a retail store.







# Fuzzy C means



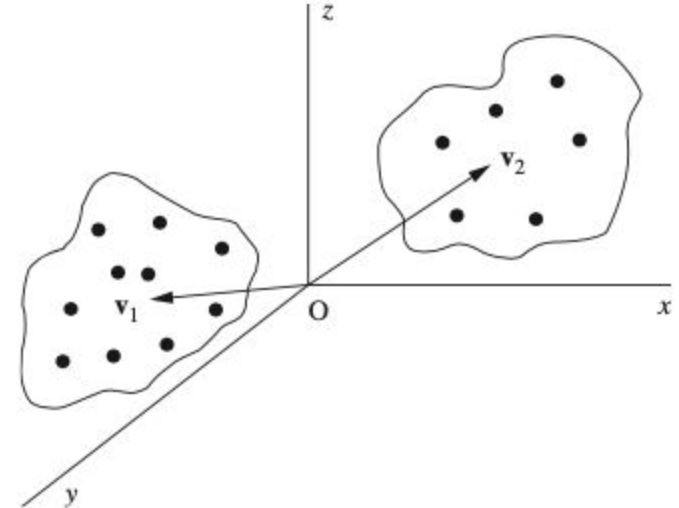


## C Means

Bezdek (1981) suggested using an objective function approach for clustering the data into hyperspherical clusters.

-> minimize the Euclidean distance between each data point in a cluster and its cluster center (a calculated point).

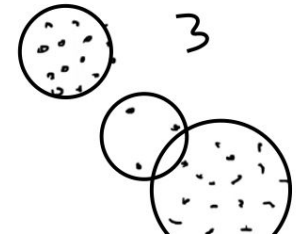
-> maximize the Euclidean distance between cluster centers.



## Fuzzy clustering

Fuzzy clustering is a clustering technique in which each data point belongs to two or more clusters. It is also referred as soft k-means or soft clustering. Clusters are identified using similarity measures. These similarity measures include connectivity, intensity and Distance.

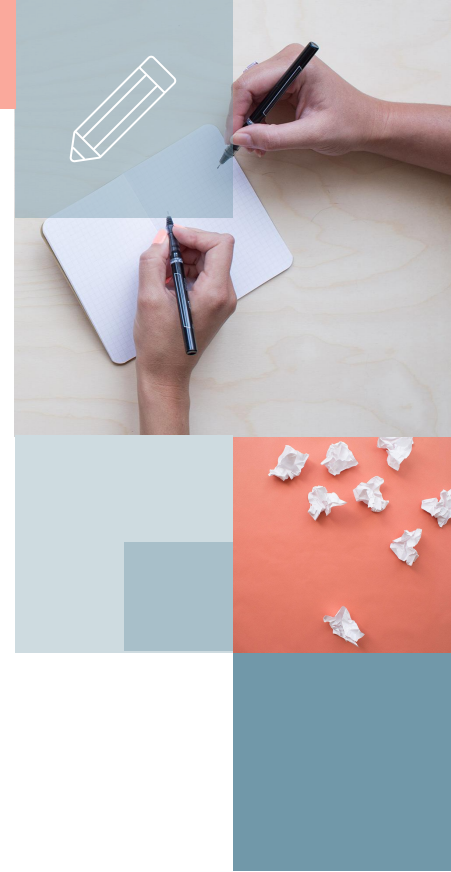
One of the most widely used fuzzy clustering algorithm is Fuzzy C-means Clustering.



## Fuzzy C-Means Clustering (FCM)

Fuzzy c-means clustering (FCM) is a data clustering technique in which the data points is grouped into  $N$  clusters with every data point that belongs to other clusters to a certain degree.

**Example:** The data points which lies close to the centre of a cluster will have a high degree of membership in that cluster, and other data points which lies far away from the centre of the cluster will have a low degree of membership to that cluster.



## Steps for fuzzy c-means clustering:

Step 1: Initialize the data points into desired number of clusters.

Step 2: Find out the centroid of each cluster.

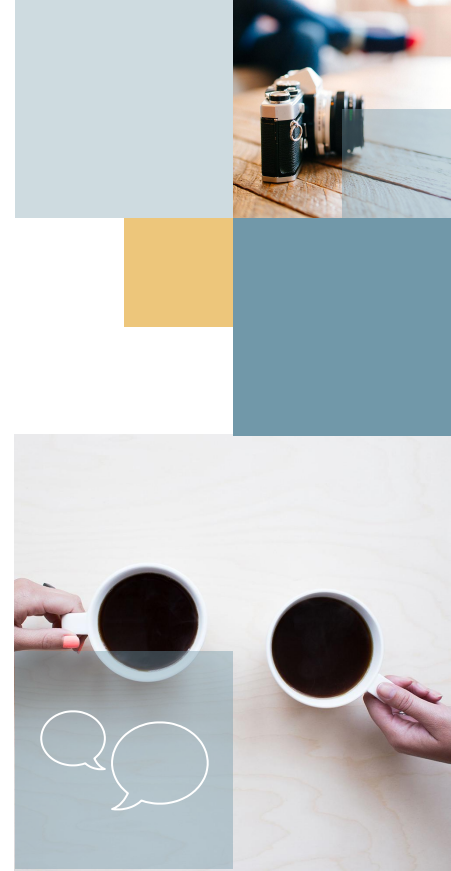
The formula for finding the centroid of cluster

$$V_{ij} = (\sum_1^n (\gamma_{ik}^m * x_k) / \sum_1^n \gamma_{ik}^m$$

Where,

m = fuzziness parameter (generally taken as 2)

$x_k$  = data point.



## Steps for fuzzy c-means clustering:

Step 3: Find out the distance of each point from both the centroids using euclidean distance

Step 4: Update the membership values using this equation.

$$\gamma = \sum_1^n (d_{ki}^2 / d_{kj}^2)^{1/m-1}]^{-1}$$

Step 5: Repeat the steps (2-4), unless centroids are not changing.



# Example:

Consider data points  $\{(1, 3), (2, 5), (6, 8), (7, 9)\}$

**Step 1:** Let's consider no of clusters is 2 in which the data is to be divided

Cluster	(1, 3)	(2, 5)	(4, 8)	(7, 9)
1	0.8	0.7	0.2	0.1
2	0.2	0.3	0.8	0.9

**Step 2:** Find out the centroid of the cluster

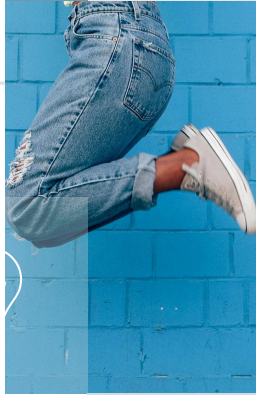
$$V_{11} = (0.82 * 1 + 0.72 * 2 + 0.22 * 4 + 0.12 * 7) / (0.82 + 0.72 + 0.22 + 0.12) = 1.568$$

$$V_{12} = (0.82 * 3 + 0.72 * 5 + 0.22 * 8 + 0.12 * 9) / (0.82 + 0.72 + 0.22 + 0.12) = 4.051$$

$$V_{21} = (0.22 * 1 + 0.32 * 2 + 0.82 * 4 + 0.92 * 7) / (0.22 + 0.32 + 0.82 + 0.92) = 5.35$$

$$V_{22} = (0.22 * 3 + 0.32 * 5 + 0.82 * 8 + 0.92 * 9) / (0.22 + 0.32 + 0.82 + 0.92) = 8.215$$

**Centroids:** (1.568, 4.051) and (5.35, 8.215)



# Contd..

**Step 3:** Find out the distance of each point from both the centroids.

$$D_{11} = ((1 - 1.568)^2 + (3 - 4.051)^2)^{0.5} = 1.2$$

$$D_{12} = ((1 - 5.35)^2 + (3 - 8.215)^2)^{0.5} = 6.79$$

Similarly, the distance of all data points is computed from both the centroids.

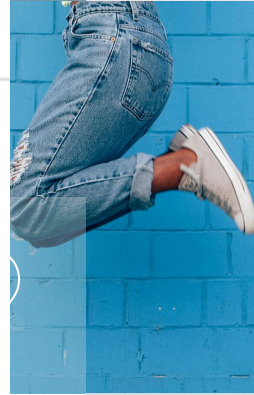
**Step 4:** Update the membership values.

$$\gamma_{11} = \{ [(1.2)^2 / (1.2)^2] + [(1.2)^2 / (6.79)^2] \}^{\{1 / (2 - 1)\}}^{-1} = 0.96$$

$$\gamma_{12} = \{ [(6.79)^2 / (6.79)^2] + [(6.79)^2 / (1.2)^2] \}^{\{1 / (2 - 1)\}}^{-1} = 0.04$$

Similarly, compute all membership values and update the matrix.

**Step 5:** Repeat the steps (2-4), unless centroids are not changing.



# Fuzzy clustering (C-Means)

Data  $X = \{x_1, x_2, \dots, x_n\}$

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

Eg:-

	$c_1$	$c_2$	$c_3$	$c_4$	
$x_1$	.02	.03	.01	.94	$\Rightarrow \sum = 1$
$x_2$	.8	.01	.10	.09	$\Rightarrow \sum = 1$
$x_3$	.01	.01	.97	.01	$\Rightarrow \sum = 1$

$w_{11} = 0.02$   
data  $\leftrightarrow$  cluster

Conditions  $\Rightarrow$

$$\sum_{j=1}^k w_{ij} = 1$$

$$0 < \sum_{i=1}^n w_{ij} < n$$

## Fuzzy c-means (FCM)

Input : Data, K

Output :  $w_{ij}, c_j$  ;  $1 \leq j \leq K$  ;  $i \leq n$  ;

### Steps

$$SSE = \sum_{j=1}^K \sum_{i=1}^n w_{ij}^p \text{dist}(x_i, c_j)^2, \quad p \in [1, \infty)$$

#### Step 1

$\rightarrow$  Initialize data points

#### Step 2

$\rightarrow$  Compute centroids

$$c_j = \frac{\sum_{i=1}^n w_{ij}^p x_i}{\sum_{i=1}^n w_{ij}^p}$$

Here  $p=2$ .

$p=0$ , same as K-means

$p=2, 4, 10, \dots$ , more blurry/clustering



Step 3

- Find distance bt. pts

Step 4

- update fuzzy pseudo-partition

$$w_{ij} = \left( \frac{1}{\text{dist}(x_i, c_j)} \right)^{1/p-1} / \sum_{q=1}^K \left( \frac{1}{\text{dist}(x_i, c_q)} \right)^{1/p-1}$$

If  $p=2$

$$w_{ij} = \frac{1 / \text{dist}(x_i, c_j)}{\sum_{q=1}^K \frac{1}{\text{dist}(x_i, c_q)}}$$



Problem

S1:

Data (1,2) (2,3) (9,4) (10,1)

$x_1$	1	2
$x_2$	2	3
$x_3$	9	4
$x_4$	10	1

	$c_1$	$c_2$
$x_1$	.4	.6
$x_2$	.88	.12
$x_3$	.41	.59
$x_4$	.27	.73
$\sum w_i$	1.18	1.25

S2:

$$\sum_{i=1}^n w_{i1}^2 = (.4)^2 + (.88)^2 + (.41)^2 + (.27)^2 = 1.18$$

$$\sum_{i=2}^n w_{i2}^2 = (.6)^2 + (.12)^2 + (.59)^2 + (.73)^2 = 1.25$$

two dimensions here,  $c_1 = [c_{11}, c_{12}]$   $c_2 = [c_{21}, c_{22}]$

$$c_{11} = \frac{(.4)^2 \times 1 + (.88)^2 \times 2 + (.41)^2 \times 9 + (.27)^2 \times 10}{1.18}$$

$$= 3.97 / 1.18 = \underline{3.38}$$

$$C_{12} = \frac{(-4)^2 \times 2 + (-88)^2 \times 3 + (-41)^2 \times 4 + (-27)^2 \times 1}{1.18}$$

$$= 2.88$$

$$C_{21} = \frac{(6)^2 \times 1 + (-12)^2 \times 2 + (-59)^2 \times 9 + (-73)^2 \times 10}{1.25}$$

$$= 7.02$$

$$C_{22} = \frac{(6)^2 \times 2 + (-12)^2 \times 3 + (-59)^2 \times 4 + (-73)^2 \times 1}{1.25}$$

$$= 2.14$$

$$C_1 = [3.38, 2.88] \quad C_2 = [7.02, 2.14]$$



S3:

	cluster centroid	
$C_1$	3.38	2.88
$C_2$	7.02	2.14

Distance

	$C_1$	$C_2$	
$x_1$	2.54	6.03	
$x_2$	1.38	5.10	
$x_3$	5.73	2.71	→ Euclidean distance
$x_4$	6.88	3.19	

S4:

$$w_{11} = \left( \frac{1}{\text{dist}(x_1, C_1)} \right) / \left( \frac{1}{\text{dist}(x_1, C_1)} + \frac{1}{\text{dist}(x_1, C_2)} \right)$$

$$= \frac{1/2.54}{1/2.54 + 1/6.03} = \frac{.39}{.39 + .56} = 0.7$$

$$w_{12} = \left( \frac{1/\text{dist}(x_1, c_2)}{1/\text{dist}(x_1, c_1) + 1/\text{dist}(x_1, c_2)} \right)$$

$$= \frac{1/6.03}{1/2.54 + 1/6.03} = \frac{0.17}{0.56} = 0.3$$

$$w_{21} = 0.79$$

$$w_{22} = 0.21$$

$$w_{31} = 0.32$$

$$w_{32} = 0.68$$

$$w_{41} = 0.32$$

$$w_{42} = 0.68$$

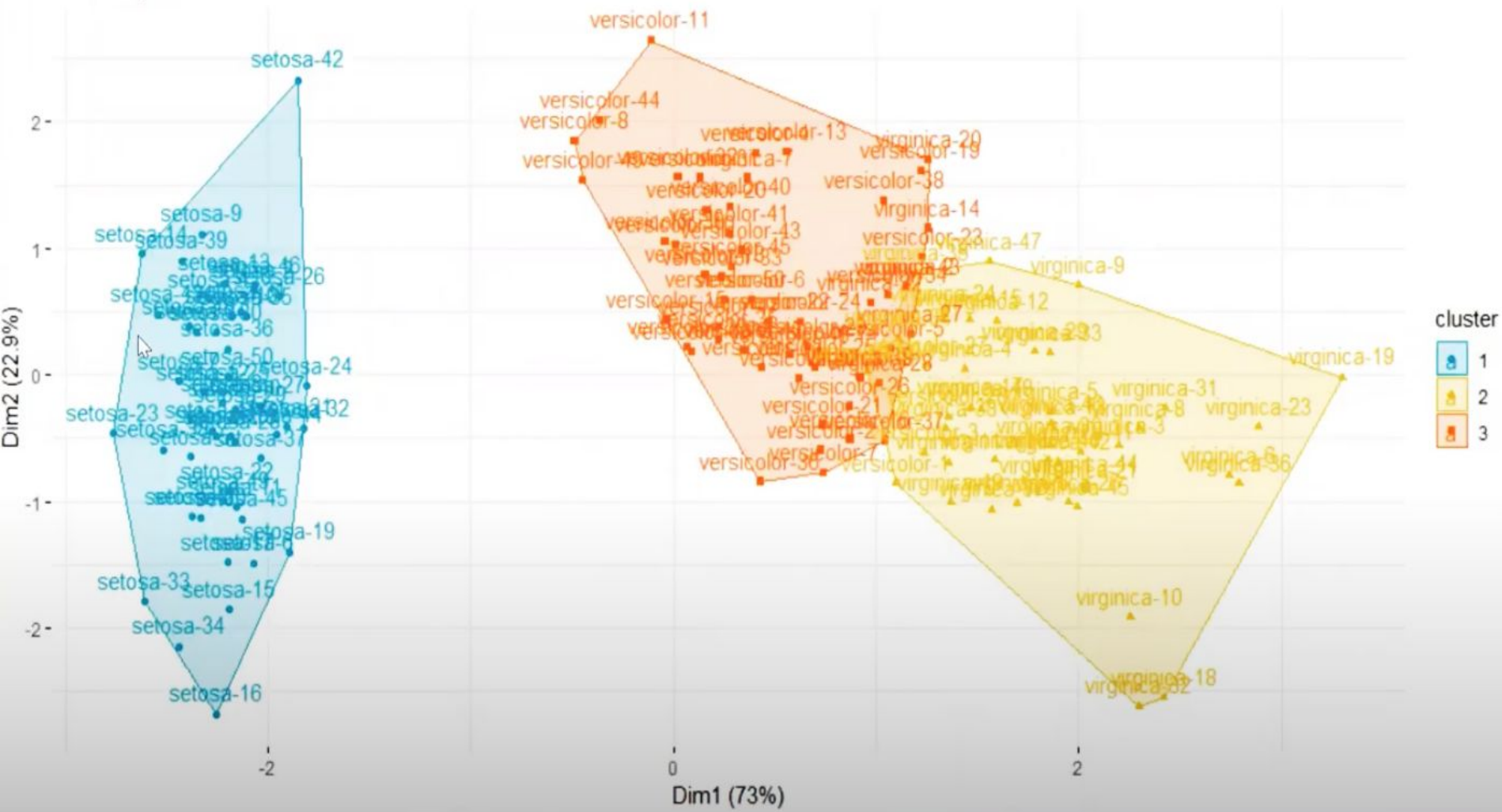
LL →

	$c_1$	$c_2$	
$x_1$	0.7	0.3	= 1
$x_2$	0.79	0.21	= 1
$x_3$	0.32	0.68	= 1
$x_4$	0.32	0.68	= 1

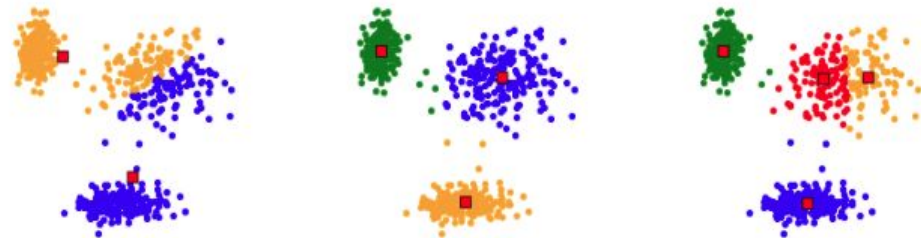
After few itr,

	$c_1$	$c_2$	cluster no
$x_1$	0.88	0.12	→ $c_1$
$x_2$	0.94	0.06	→ $c_1$
$x_3$	0.17	0.83	→ $c_2$
$x_4$	0.18	0.82	→ $c_2$

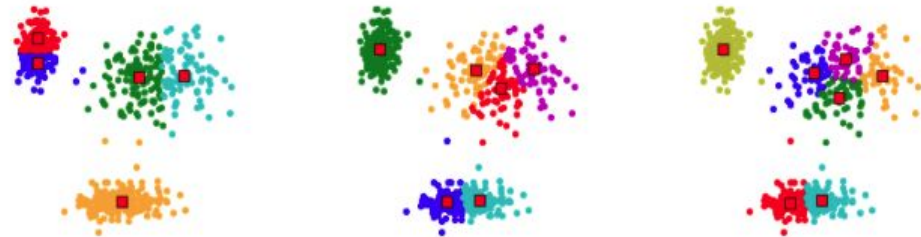
Cluster plot



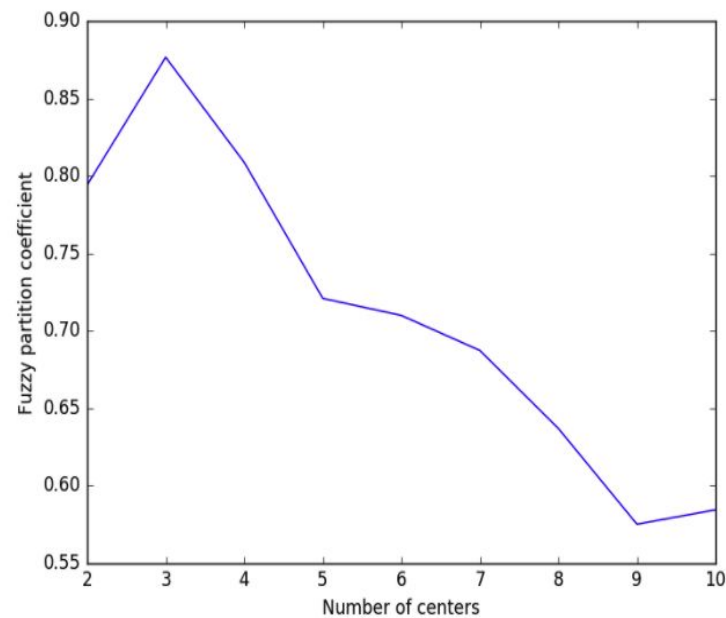
Centers = 2; FPC = 0.79    Centers = 3; FPC = 0.88    Centers = 4; FPC = 0.81



Centers = 5; FPC = 0.72    Centers = 6; FPC = 0.71    Centers = 7; FPC = 0.69

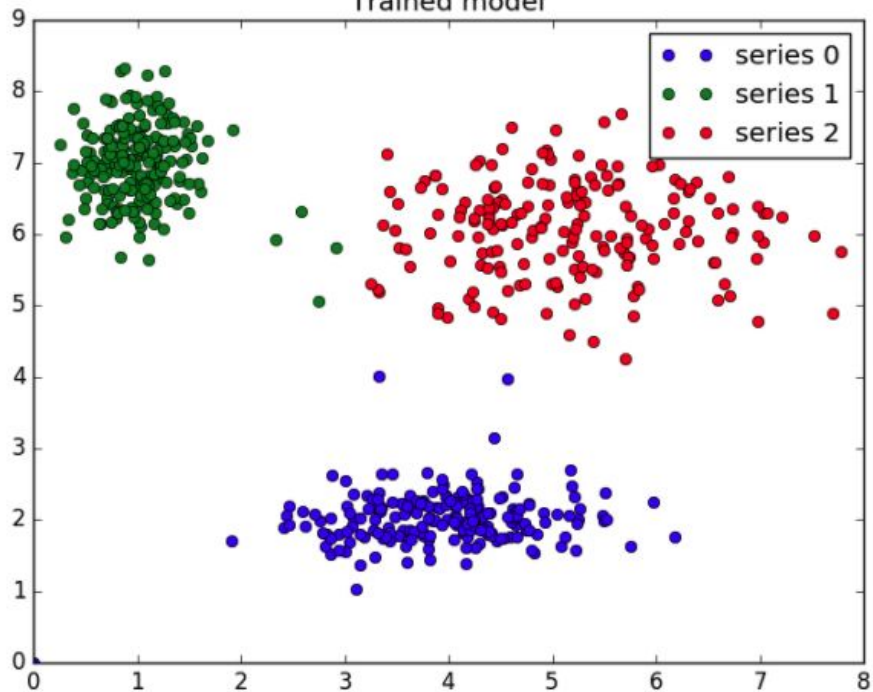


Centers = 8; FPC = 0.64    Centers = 9; FPC = 0.58    Centers = 10; FPC = 0.58

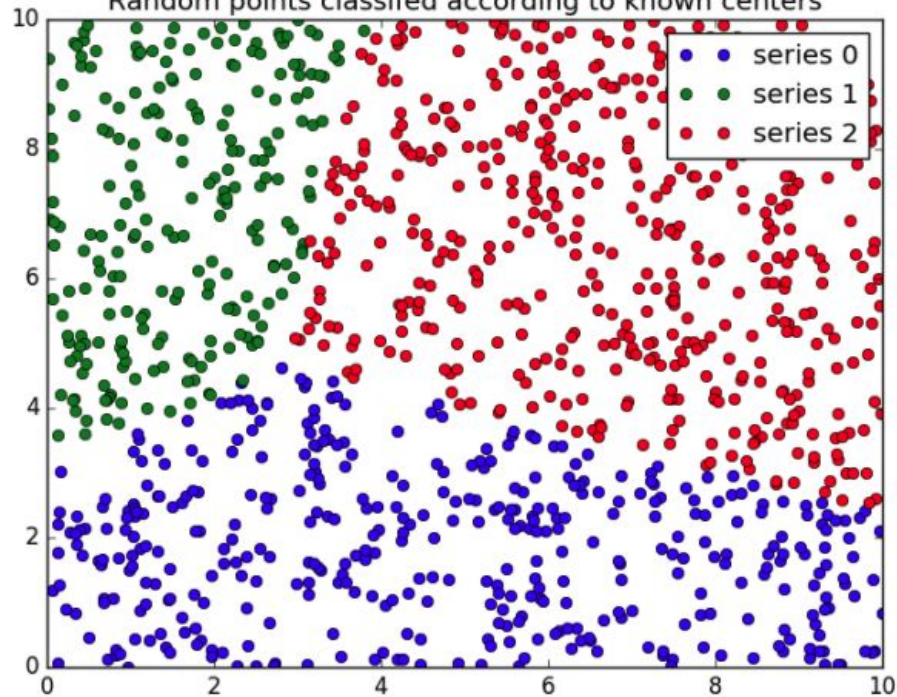




Trained model



Random points classified according to known centers



## Advantages

- 1) Gives best result for overlapped data set and comparatively better than k-means algorithm.
- 2) Unlike k-means where data point must exclusively belong to one cluster center here data point is assigned membership to each cluster center as a result of which data point may belong to more than one cluster center.

## Disadvantages

- 1) Apriori specification of the number of clusters.
- 2) We get the better result but at the expense of more number of iteration.
- 3) Euclidean distance measures can unequally weight underlying factors.



Thanks!

**Any** questions?

