Data Science Home work - 5.

Phani. Teja Kesha
KP38691

1. The given Points are

A₁ (2, 10)

A₂ (2, 5)

A₃ (8, 4)

B₁ (5, 8)

B₂ (7, 5)

B₃ (6, 4)

C₁ (1, 2)

C₂ (4, 9)

Here A₁, B₁, C₁ are assigned

as the initial centers of the cluster.

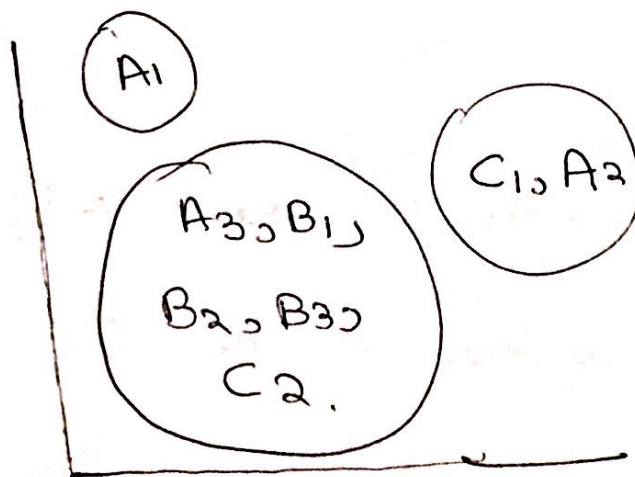→ Eucledian distance gives you the distance

between two Points.

If (x₁, y₁) & (x₂, y₂), The distance is

$$\sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

After a single iteration.

|       | $A_1$ | $A_2$ | $C_1$ $B_1$ $B_2$ $B_3$ $C_1$ $C_2$ |
|-------|-------|-------|-------------------------------------|
| $A_1$ | $\boxed{0}$ | 3.6 | 8.06 |
| $A_2$ | 5 | 4.2 | $\boxed{3.16}$ |
| $A_3$ | 8.4 | $\boxed{5}$ | 7.28 |
| $B_1$ | 3.6 | $\boxed{0}$ | 7.21 |
| $B_2$ | 7.07 | $\boxed{2.6}$ | 6.7 |
| $B_3$ | 7.2 | $\boxed{4.1}$ | 5.3 |
| $C_1$ | 8.06 | 7.2 | $\boxed{0}$ |
| $C_2$ | 2.2 | $\boxed{1.41}$ | 7.6 |

$\left.\begin{array}{c} \\ \\ \end{array}\right\}$ Picking the nearest centers from given values.



→ Let's check for the Second iteration.
For the next iterations centers are updated by doing the

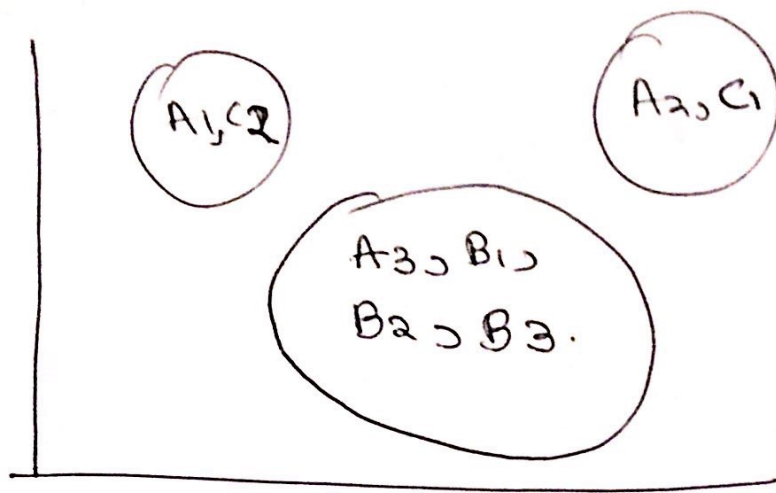average of all the Points w/o the given cluster.

$$C_1 = (2, 10)$$

$$C_2 = (6, 6)$$

$$C_3 = (1.5, 3.5)$$

Second Iteration :-

| | $C_1$ | $C_2$ | $C_3$ |
|---|---|---|---|
| $A_1$ | $\boxed{0}$ | 5.65 | 6.5 |
| $A_2$ | 5 | 4.12 | $\boxed{1.58}$ |
| $A_3$ | 8.4 | $\boxed{2.82}$ | 6.51 |
| $B_1$ | 3.6 | $\boxed{2.23}$ | 5.7 |
| $B_2$ | 7.07 | $\boxed{1.41}$ | 5.7 |
| $B_3$ | 7.2 | $\boxed{2}$ | 4.53 |
| $C_1$ | 8.06 | 6.4 | $\boxed{1.58}$ |
| $C_2$ | $\boxed{2.2}$ | 3.6 | 6.04 |

A cluster diagram showing three groups: circle "A₁,C₂", circle "A₂,C₁", and oval "A₃, B₁, B₂, B₃."

→ Again updating the centers.

$$C_1 = (3, 9.5)$$
$$C_2 = (6.5, 5.25)$$
$$C_3 = (1.5, 3.5)$$

Third Iteration:

|     | $C_1$ | $C_2$ | $C_3$ |
|-----|-------|-------|-------|
| A₁  | 1.11  | 6.5   | 6.5   |
| A₂  | 4.6   | 4.5   | 1.58  |
| A₃  | 7.4   | 1.95  | 6.51  |
| B₁  | 2.5   | 3.13  | 5.7   |
| B₂  | 6.02  | 0.5   | 5.7   |
| B₃  | 6.26  | 1.34  | 4.52  |
| C₁  | 7.7   | 6.3   | 1.58  |
| C₂  | 1.11  | 4.5   | 6.04  |

→ Calculating Centers

$$C_1 = (3.66, 9)$$
$$C_2 = (7, 433)$$
$$C_3 = (1.5, 3.5.$$

**Fourth Iteration :-**

| | $C_1$ | $C_2$ | $C_3$ |
|---|---|---|---|
| $A_1$ | **1.93** | 7.55 | 6.5 |
| $A_2$ | 4.33 | 5.04 | **1.58** |
| $A_3$ | 6.63 | **4.05** | 6.5 |
| $B_1$ | **1.67** | 4.7 | 5.7 |
| $B_2$ | 5.2 | **0.67** | 5.7 |
| $B_3$ | 5.52 | **1.05** | 4.52 |
| $C_1$ | 7.4 | 6.43 | **1.58** |
| $C_2$ | **6.34** | 5.5 | 6.04 |

// centers are same, & not any Change.

→ Since the centers are not changed, it has reached the best possible cluster.

1a)

Center 1
$A_1$

Center 2
$A_3, B_1, B_2$
$B_3, C_2$

Center 3
$C_1, A_2$.

1b)

Center 1
$A_1, B_1, C_2$

Center 2
$A_3, B_2, B_3$

Center 3
$A_2, C_1$

2)
10.10) why does BIRCH encounter difficulty in finding clusters of arbitary shape but OPTICS does not? Propose modifications to BIRCH to help it find clusters of arbitary shape?

BIRCH uses the euclidean distance to find distance between the points, Due to this the shape is nearly spherical, the resulting cluster

is not in the arbitary shape.
whereas OPTICS uses density of
the Points as the metric for distance
all the Points which are very close
enough or within a minimum density
measure form a cluster which can be
of any arbitary shapes which is
based on the original Position of
the Points.

As a modification to BIRCHs it
can be modified to use dyram density measure
to form arbitory shaped Clusters by
clustering very low level B+ trees, which
have very closely Positioned Points. This
will form take density measure and not
distance measures, So that form a CF tree
and will result in arbitary shaped Clusters.

③

→ Partition based and hierarchial clustering uses distance as a measure for creating clusters of the given points.

→ Due to this distance measures they end up having spherical shaped clusters.

→ But in the case of density based clustering, takes uses of the fact that the every cluster formed has a density different than other formed clusters.

→ Density clustering method takes the clusters as the dense region of the data points

→ In this way they arbitarily forms arbitary shaped cluster dense regions and are very suitable than other clustering methods in a way that they can separate those from less density regions in the data point space.

4) Basically, we are given the datapoints and need to build the clustering depending on the Contstraints and Kmeans or other clustering techniques.

→ Lets choose some random ATM's as Centroids. and for each point we assign it with the nearest Centroid. satisfying the given Constraints which we have enforced on the data Points.

→ If a datapoint there is no Kentroid satisfying the Constraint. then the datapoint is not assigned to any Centroid. and the datapoint is then updated to the Centroid list.

→ This Process is repeated ontill any Convergence is found.

→ The algorithm goes as follows

1) Define Constraints function which has

Constraints like → 10000 households per cluster

→ No obstacle objects

2) Randomly choose Centroid Datapoints & assign it to all Datapoints

3) If a Datapoint has no nearest centroid with given requirements then add the datapoint to Centroid list.

4) Repeat This untill Convergence.

5) Then the Points with no clusters assigned by ATM are to by dropping all the Constraints.