Homework-4

PART-A

1) **K-Means Algorithm:**

{ (2,10) (2,5) (8,4) (5,8) (7,5) (6,4) (1,2) (4,9) }

a)



b) **3 clusters**

Cluster 1 → (1,2) (2,5)

Cluster 2 → (2,10) (4,9) (5,8)

Cluster 3 → (6,4) (7,5) (8,4)

c) initial centers { (2,5) (5,8) (4,9) }

Cluster 1 → (2,5) (1,2)

Cluster 2 → (5,8) (8,4) (7,5) (6,4)

Cluster 3 → (4,9) (2,10)

d) center 1 ⟹ $\left(\frac{2+1}{2}, \frac{5+2}{2}\right) = (1.5, 3.5)$

center 2 ⟹ $\left(\frac{5+8+7+6}{4}, \frac{8+4+5+4}{4}\right) = (6.5, 5.25)$

center s $\Rightarrow$ $\left(\dfrac{4+2}{2}, \dfrac{9+10}{2}\right) = \left(3, 9.5\right)$.

e) Center's after second iteration?

clustering using $\{ (1.5, 3.5) \ (6.5, 5.25)$

$(3, 9.5) \}$

cluster 1 $\Rightarrow$ $(2,5) \ (1,2), P$

cluster 2 $\Rightarrow$ $(8,4) \ (7,5) \ (6,4)$.

cluster 3 $\Rightarrow$ $(2,10) \ (5,8) \ (4,9)$.

center 1 $\Rightarrow$ $\left(3/2, \ 7/2\right) \Rightarrow \left(1.5, 3.5\right)$

center 2 $\Rightarrow$ $\left(21/3, \ 13/3\right) \Rightarrow \left(7, 4.3\right)$

center 3 $\Rightarrow$ $\left(11/3, \ 27/3\right) \Rightarrow \left(3.6, 9\right)$.

b) Centers after third iteration?

clustering using $\{ (1.5, 3.5), \ (7, 4.3) \ (3.6, 9) \}$

cluster 1 $\Rightarrow$ $(2,5) \ (1,2)$

cluster 2 $\Rightarrow$ $(8,4) \ (7,5) \ (6,4)$

cluster 3 $\Rightarrow$ $(2,10), \ (5,8) \ (4,9)$

center 1 $\Rightarrow$ $(1.5, 3.5)$.

center 2 $\Rightarrow$ $(7, 4.3)$.

center 3 $\Rightarrow$ $(3.6, 9)$.

g) cluster 1 ⇒ $\left[(1,2)(2,5)\right]$ $\left[(1,2)(2,5)\right]$

cluster 2 ⇒ $\left[(2,10)(4,9)(5,8)\right]\left[(8,4)(4,9)(8,8)\right]$

cluster 3 ⇒ $\left[(8,4)(7,5)(6,4)\right.$ $\left[(8,4)(7,5)(6,4)\right]$

Same.

h) No of iterations for the clusters to converge
is 2. The clusters did not change after
that.

i) Resulting centers : $\{(1.5, 3.5)\ (7, 4.3)\ (3.6, 9)\}$
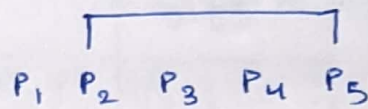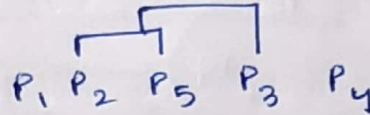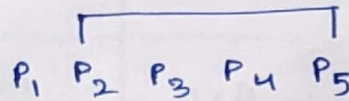
Resulting clusters :

cluster 1 ⇒ $(2,5)(1,2)$.

cluster 2 ⇒ $(8,4)(7,5)(6,4)$.

cluster 3 ⇒ $(2,10)(5,8)(4,9)$.

2) Hierarchial algorithm

1) Single link:

$(P_2, P_5) - 0.98$

P₁  P₂  P₃  P₄  P₅

$(P_3, P_5) - 0.85$

P₁  P₂  P₃  P₄  P₅

P₁  P₂  P₅  P₃  P₄

$(P_4, P_5) - 0.75$

P₁  P₂  P₅  P₃  P₄

$(P_2, P_3) - 0.64$   (already exists)

$(P_1, P_4) - 0.55$

P₁  P₂  P₅  P₃  P₄

2) Complete link:

$(P_2, P_5) - 0.98$

P₁  P₂  P₃  P₄  P₅

|  | $P_1$ | $P_2 \cup P_5$ | $P_3$ | $P_4$ |
|---|---|---|---|---|
| $P_1$ | 1.0 | 0.1 | 0.41 | 0.55 |
| $P_2 \cup P_5$ | 0.10 | 1.0 | 0.64 | 0.47 |
| $P_3$ | 0.41 | 0.64 | 1.0 | 0.44 |
| $P_4$ | 0.55 | 0.47 | 0.44 | 1.0 |

P₁  P₂  P₅  P₃  P₄

| | $P_1$ | $P_2 \cup P_3 \cup P_5$ | $P_4$ |
|---|---|---|---|
| $P_1$ | 1.0 | 0.1 | 0.55 |
| $P_2 \cup P_3 \cup P_5$ | 0.1 | 1.0 | 0.44 |
| $P_4$ | 0.55 | 0.44 | 1.0 |



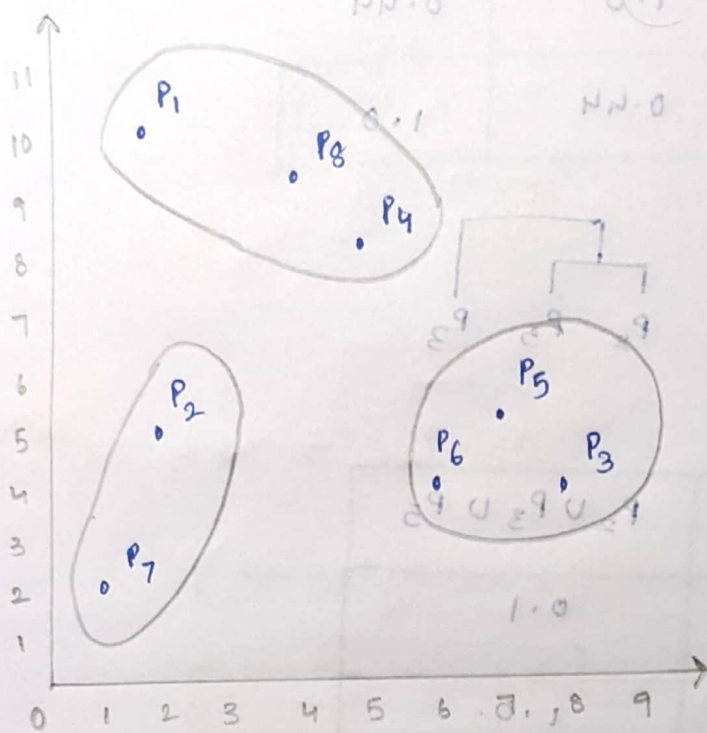| | $P_1 \cup P_4$ | $P_2 \cup P_3 \cup P_5$ |
|---|---|---|
| $P_1 \cup P_4$ | 1.0 | 0.1 |
| $P_2 \cup P_3 \cup P_5$ | 0.1 | 1.0 |



3) DBSCAN Algorithm:

epsilon = 2
min_sample = 2

epsilon = $\sqrt{10}$.

min_sample = 2.



PART-B

4) **Memtable:** Bigtable servers store recent updates in memory in a data structure known as memtable.

Reading from memtable is quick and if the server crashes, memtable can be constructed again using commit log.

5) **SSTable:** A sorted string table is used to store table data in google file system. It is different from relational database as relational databases don't provide partition tolerance.

6) CAP Theorems

For any system that shares data it is impossible to guarentee simultaneously the following properties:

1) Consistency.

2) Availability.

3) Partition-Tolerance.

very large scale system will partition at some point.

1) it is essential to decide between consistency and availability.

2) Traditional DBMS prefer consistency over availability and partition tolerance.

3) Most web applications choose availability.

7) A <u>Tablet</u> is a set of consecutive rows of a table and is a unit of distribution and load balancing within Bigtable.

A <u>tablet server</u> stores and serves tablets to clients. For a given tablet, a tablet server acts as a leader and other servers follow. replicas of the tablet.