

BIRCH - Balanced Iterative Reducing & Clustering with Hierarchies

cluster the following data points using BIRCH.

B = Branching factor = 3

Threshold T = 5.

Points 2, 10, 12, 4, 25, 3, 30, 20, 11.

Clustering feature CF = $\langle n, LS, SS \rangle$

n = number of points in cluster

LS = linear sum of n points, $\sum_{i=1}^n x_i$

SS = square sum of n points, $\sum_{i=1}^n x_i^2$.

cluster centroid $x_0 = \frac{LS}{n}$.

cluster diameter D = $\sqrt{\frac{2nSS - 2LS^2}{n(n-1)}}$.

Step 1 ph. d.

CF-Tree

CF = $\langle 1, 2, 4 \rangle$

root b=3.

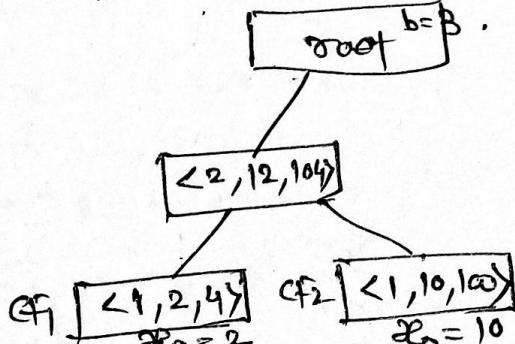
CF₁ $\langle 1, 2, 4 \rangle$

Step 2 ph. 10 Try adding ph. 10 to CF₁. find new CF & diameter

CF₁' = $\langle 2, 2+10, 4+100 \rangle = \langle 2, 12, 104 \rangle$.

$$d = \sqrt{\frac{2 \times 2 \times 104 - 2 \times 12 \times 12}{2(1)}} = \sqrt{60} = 8.$$

d < T, split CF₁.



Step 3. ph. 12 Try adding ph. 12 to CF₁ & CF₂. find new CF & d.

CF₁', ph. 12 is more closer to CF₂ as $|12-10| < |12-2|$

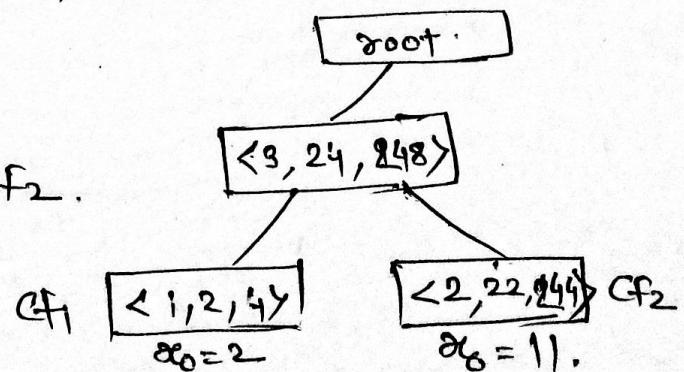
CF₂', ph. 12 is more closer to CF₂ as $|12-2| < 10$.

Try adding ph. 12 to CF₂.

CF₂' = $\langle 2, 10+12, 100+144 \rangle = \langle 2, 22, 244 \rangle$.

$$d = \sqrt{\frac{2 \times 2 \times 244 - 2 \times 22 \times 22}{2(1)}} = 2$$

d < 2 < T, add ph. 12 to CF₂.



$\Delta_0 = 2 < 7$

Try adding pt. 4 to CF_1 ,

$$CF'_1 = \langle 2, 6, 20 \rangle$$

$$d = \sqrt{\frac{2 \times 2 \times 20 - 2 \times 6 \times 6}{2(1)}} = 2$$

$d < T$, add pt. 4 to CF_1 .

Step 5 - pt. 25 \rightarrow pt. 25 is more closer to CF_2 as $|25-3| > |25-11|$

Try adding pt. 25 to CF_2

$$CF'_2 = \langle 2, 22, 244 \rangle, \langle 3, 47, 869 \rangle$$

$$d = \sqrt{\frac{2 \times 2 \times 244 - 2 \times 47 \times 47}{2(1)}} =$$

$$d = \sqrt{\frac{2 \times 2 \times 869 - 2 \times 47 \times 47}{2(1)}} = \text{undefined} - \text{ve}$$

Hence add this as separate CF

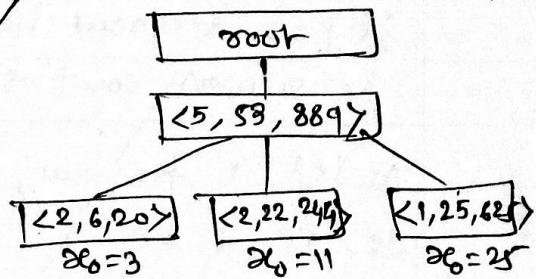
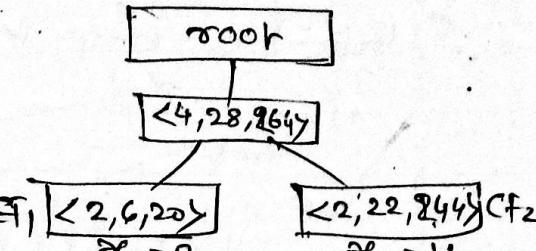
Step 6 - pt. 3 - pt. 3 is more closer to CF_1 as $|3-3| < |3-11| < |3-25|$

Try adding pt. 3 to CF_1 .

$$CF'_1 = \langle 3, 6+3, 20+9 \rangle = \langle 3, 9, 29 \rangle$$

$$d = \sqrt{\frac{2 \times 3 \times 29 - 2 \times 9 \times 9}{3(2)}} = 4.8$$

$d < T$, add pt. 3 to CF_1 .



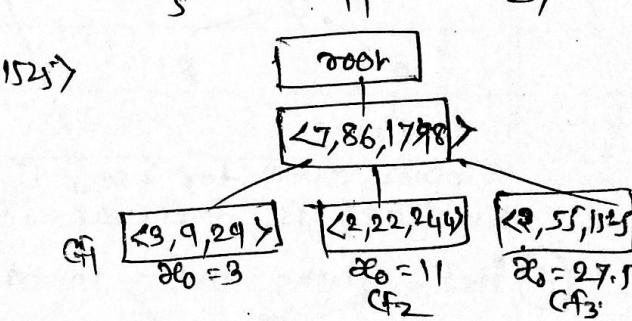
Step 7 - pt. 30 - pt. 30 is more closer to CF_3 as $|30-25| < |30-11| < |30-31|$

Try adding pt. 30 to CF_3 .

$$CF'_3 = \langle 2, 25+30, 625+900 \rangle = \langle 2, 55, 1525 \rangle$$

$$d = \sqrt{\frac{2 \times 2 \times 1525 - 2 \times 55 \times 55}{2(1)}} = 5.$$

$d \leq T$, add pt. 30 to CF_3 .



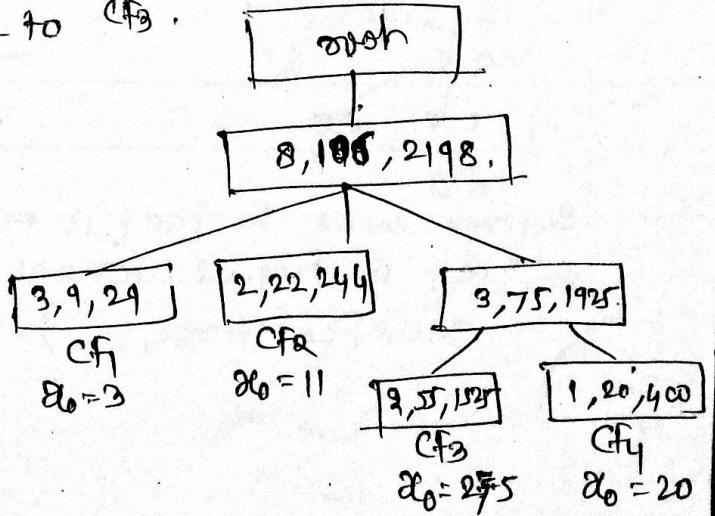
Step 8 - pt. 20 \rightarrow pt. 20 is more closer to CF_3 .

Try adding to CF_3 .

$$CF'_3 = \langle 3, 75+20, 1525+400 \rangle = \langle 3, 75, 1925 \rangle$$

$$d = \sqrt{\frac{2 \times 3 \times 1925 - 2 \times 75 \times 75}{3(2)}} = 7.07$$

$d > T$, split CF_3' into 2 nodes.



J - If, 11 and 13 is more closer to C_2 ,

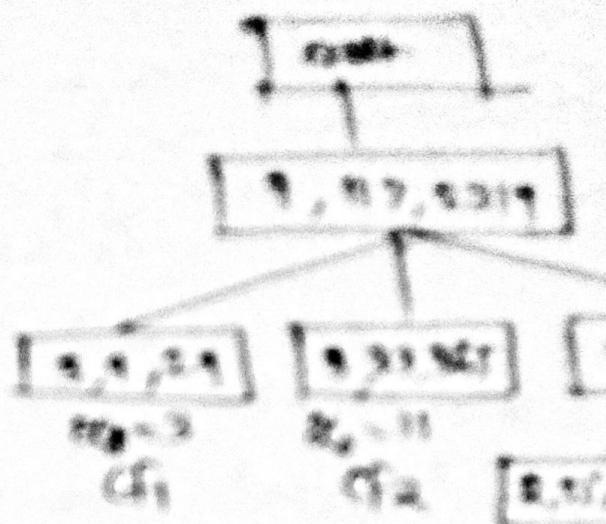
\Rightarrow Assign 11 to C_2 .

$$C_2 = \langle 3, 11, 13, 14, 15, 16 \rangle$$

$$\therefore \langle 9, 10, 12 \rangle.$$

$$\frac{2+3+9+10+12+13+15+16}{8(3)} = 14.1.$$

$\therefore C_1$, and 11 to C_2 .



Final Clusters are -

$$C_1 = \langle 9, 10, 12 \rangle$$

$$C_2 = \langle 13, 15, 16 \rangle$$

$$C_3 = \langle 11, 12, 14, 16 \rangle$$

$$C_4 = \langle 20 \rangle.$$

DSCAN

Use DSCAN to cluster data points $\{2, 10, 12, 4, 25, 3, 30, 20, 11\}$ with Eps radius $\epsilon = 2$ & minPts=2. Mark points clearly as core & noise points. If $\epsilon=5$ is used, whether number of clusters are changed?

Step 0 Unvisited = $\{2, 10, 12, 4, 25, 3, 30, 20, 11\}$

Visited = $\{\}$

Noise = $\{\}$.

Step 1 - Compute ϵ neighbourhood of each point.

$$N(2) = \{4, 3\}$$

$$d(2) = 2$$

$$N(3) = \{2, 4\} \approx$$

$$d(3) = 3$$

$$N(10) = \{12, 11\}$$

$$d(10) = 3$$

$$N(12) = \{10, 11\}$$

$$d(10) = 1$$

$$N(12) = \{2, 3\}$$

$$d(12) = 3$$

$$N(20) = \{3\}$$

$$d(20) = 1$$

$$N(4) = \{2, 3\}$$

$$d(4) = 3$$

$$N(11) = \{10, 12\}$$

$$d(11) = 3$$

$$N(25) = \{\}$$

$$d(25) = 1$$

$$N(\cdot) = \{\}$$

$$\text{Core pts} = \{2, 10, 12, 4, 3, 11\}$$

$$\text{Noise pts} = \{25, 30, 20\}$$

Step 2 - Iterate over each core pt.

pt. 2 Visited = $\{2\}$. $C_1 = \{2\}$.

$N(2) = \{4, 3\}$. As 4 is core pr.

$p' = 4$ Visited = $\{2, 4\}$ $N(2) = N(2) \cup N(4)$,

$N(4) = \{2, 3\} \cup \{2, 3\} = \{2, 3, 4\}$, $C_1 = \{2, 4\}$.

$p' = 3$ Visited = $\{2, 4, 3\}$.

As 3 is core pr.

$N(2) = N(2) + N(3) = \{2, 3, 4\} \cup \{2, 3\} = \{2, 3, 4\}$

$C_1 = \{2, 4, 3\}$.

$C_1 = \{2, 4, 3\}$.

pt. 10 Visited = $\{2, 4, 3, 10\}$.

$C_2 = \{10\}$ $N(10) = \{12, 11\}$.

$p' = 12$ Visited = $\{2, 4, 3, 10, 12\}$

As 12 is core pr. $N(10) = N(10) \cup N(12) = \{12, 11, 10\} \cup \{10, 11\} = \{12, 11, 10\}$

$C_2 = \{10, 12\}$.

$p' = 11$ Visited = $\{2, 4, 3, 10, 12, 11\}$.

As 11 is core pr. $N(10) = N(10) \cup N(11) = \{12, 11, 10\} \cup \{10, 11\} = \{12, 11, 10\}$

$C_2 = \{10, 12, 11\}$.

$p' = 10$ 10 is already visited.

$C_2 = \{10, 12, 11\}$.

$$\text{Visited} = \{2, 4, 3, 10, 12, 11\}$$

$$\text{Noise} = \{25, 30, 20\}$$

$$\text{Clusters} \quad C_1 = \{2, 4, 3\}$$

$$C_2 = \{10, 12, 11\}$$

Step 1 $t=5$, lets repeat the process.

$$N(2) = \{4, 3\} \quad d(2) = 3$$

$$N(12) = \{10, 11\} \quad d(12) = 3$$

$$N(25) = \{30, 20\} \quad d(25) = 3$$

$$N(3) = \{2, 4\} \quad d(3) = 3$$

$$N(20) = \{25, 30\} \quad d(20) = 3$$

$$\text{Corept} = \{2, 10, 12, 4, 25, 3, 30, 20, 11\}$$

$$\text{Noise} = \{25\}$$

$$\text{Visited} = \{2, 4, 3\} \quad \text{Unvisited} = \{10, 12, 11\}$$

Step 2 for 2, Visited = {2}. $C_1 = \{2\}$. As 2 is core pt form new cluster.

$$N(2) = \{4, 3\}$$

$$p=4 \quad \text{visited} = \{2, 4\} \quad \text{As } 4 \text{ is core pt} \quad N(2) = N(2) \cup N(4)$$

$$= \{4, 3\} \cup \{2, 3\}$$

$$C_1 = \{2, 4\} \quad = \{4, 3, 2\}$$

$$p=3 \quad \text{visited} = \{2, 4, 3\} \quad \text{As } 3 \text{ is core pt} \quad N(2) = N(2) \cup N(3)$$

$$= \{4, 3, 2\} \cup \{2, 4\}$$

$$C_1 = \{2, 4, 3\} \quad = \{4, 3, 2\}$$

$$C_1 = \{2, 4, 3\} \quad \text{As all points in } N(2) \text{ are visited now.}$$

Step 3 for 10 Visited = {2, 4, 3, 10} $C_2 = \{10\}$ as 10 is core pt.

$$N(10) = \{12, 11\}$$

$$p=12 \quad \text{visited} = \{2, 4, 3, 10, 12\}$$

$$\text{As } 12 \text{ is core pt} \quad N(10) = N(10) \cup N(12) = \{12, 11, 4\} \cup \{10, 11\}$$

$$C_2 = \{10, 12\} \quad = \{12, 11, 10\}$$

$$p=11 \quad \text{visited} = \{2, 4, 3, 10, 12, 11\}$$

$$\text{As } 11 \text{ is core pt, } N(10) = N(10) \cup N(11) = \{12, 11, 10\} \cup \{10, 12\}$$

$$= \{12, 11, 10\}$$

$$C_2 = \{10, 12, 11\} \quad \text{As all points in } N(12) \text{ are visited now}$$

Step 4 for 25 Visited = {2, 4, 3, 10, 12, 11, 25} $C_3 = \{25\}$ as 25 is core pt.

$$N(25) = \{30, 20\}$$

$$p=30 \quad \text{visited} = \{2, 4, 3, 10, 12, 11, 25, 30\}$$

$$\text{As } 30 \text{ is core pt. } N(25) = N(25) \cup N(30) = \{30, 20\} \cup \{25, 20\}$$

$$= \{30, 20, 25\}$$

$$C_3 = \{25, 30\}$$

$p = 20$ $\text{visited} = \{2, 4, 3, 10, 12, 11, 25, 30, 20\}$

As 20 is core for, $N(25) = N(25) \cup N(20) = \{30, 20, 25\} \cup \{25, 30\}$
 $C_3 = \{25, 30, 20\}$,
 $= \{30, 20, 25\}$.

As all points in $N(25)$ are visited now.

All points are visited, now.

final clusters, $C_1 = \{2, 4, 3\}$
 $C_2 = \{10, 12, 11\}$,
 $C_3 = \{25, 30, 20\}$,

As ϵ is changed from 2 to 5, number of clusters increased by one.

Evaluation of clustering - Hopkins Stat.
 finding out uniformity in data distribution.
 If Hopkins stat = H ≈ 0.5 , uniformly distributed.
 Otherwise not.

④ Determine whether following dataset has clustering tendency.

$$\{2, 10, 12, 4, 25, 3, 30, 20, 11\} = D.$$

→ Step 1 - Randomly sample 5 points from D.

$$\text{sample} = \{2, 12, 25, 30, 11\}.$$

find out pts which are nearest to each pt in sample & sum all those nearest distances.

$$\begin{aligned}\sum d_{ij} &= d(2, 3) + d(12, 11) + d(25, 20) + d(30, 25) + d(11, 10), \\ &= 1 + 1 + 5 + 5 + 1 \\ &= 18.\end{aligned}$$

Step 2 - Randomly sample 5 points from but don't place them back in D.

$$\begin{array}{lll}S_1 = \{12\} & D = \{2, 10, 4, 25, 3, 30, 20, 11\} & y_1 = d(12, 11) = 1, \\ S_2 = \{30\} & D = \{2, 10, 4, 25, 3, 20, 11\} & y_2 = d(30, 25) = 5, \\ S_3 = \{4\} & D = \{2, 10, 25, 3, 20, 11\} & y_3 = d(4, 3) = 1, \\ S_4 = \{11\} & D = \{2, 10, 25, 3, 20\} & y_4 = d(11, 10) = 1, \\ S_5 = \{3\} & D = \{2, 10, 25, 20\} & y_5 = d(3, 2) = 1.\end{array}$$

$$\sum y_i = 1 + 5 + 1 + 1 + 1 = 9.$$

$$H = \frac{\sum_{i=1}^n y_i}{\sum_{i=1}^n d_{ij} + \sum_{i=1}^n y_i}$$

$$= \frac{9}{18 + 9}$$

$$= \frac{9}{27}$$

$$= 0.4.$$

Looks data is uniformly distributed as per Hopkins stat so it's not much suitable for data clustering.

Evaluation of Clustering - Elbow method

Determining number of clusters.

- ① Consider following dataset, using different values of clusters.
i.e. $k=2, 3, 4, 5, \dots$, determine best value of k using elbow technique.
 $\rightarrow D = \{2, 10, 12, 4, 25, 3, 30, 20, 11\}$.

Assuming Manhattan dist as measure, & initial centre 3, 10.

$$k=2, C_1 = \{3, 2, 4\} \\ C_2 = \{10, 12, 25, 30, 20, 11\}$$

$VAR(k) = \text{sum of squared dist of each pt from its centroid}$

$$\begin{aligned} &= (3-3)^2 + (3-2)^2 + (3-4)^2 + \\ &\quad (10-10)^2 + (12-10)^2 + (25-10)^2 + (30-10)^2 + (20-10)^2 + (11-10)^2 \\ &= 0+1+1+0+4+225+400+100+1 \\ &= 732 \end{aligned}$$

$$\begin{array}{ll} k=3 & C_1 = \{3, 2, 4\} \\ C_2 = \{10, 12, 11\} \\ C_3 = \{25, 30, 20\} \\ C_4 = 25 \end{array}$$

$$\begin{aligned} VAR(k) &= (3-3)^2 + (3-2)^2 + (3-4)^2 \\ &\quad + (10-10)^2 + (12-10)^2 + (11-10)^2 \\ &\quad + (25-25)^2 + (25-30)^2 + (25-20)^2 \\ &= 0+1+1+0+4+1+0+25+25 = 57 \end{aligned}$$

$$\begin{array}{ll} k=4 & C_1 = \{3, 2, 4\} \\ C_2 = \{10, 12, 11\} \\ C_3 = \{20, 25, 20\} \\ C_4 = \{30, \} \end{array}$$

$$\begin{aligned} VAR(k=4) &= (3-3)^2 + (3-2)^2 + (3-4)^2 \\ &\quad + (10-10)^2 + (10-12)^2 + (10-11)^2 \\ &\quad + (20-20)^2 + (20-25)^2 + (20-25)^2 \\ &\quad + (30-30)^2 \\ &= 0+1+1+0+4+1+0+25+0 = 32 \end{aligned}$$

$$\begin{array}{ll} k=5 & C_1 = \{3, 2, 4\} \\ C_2 = \{10\} \\ C_3 = \{12, 11\} \\ C_4 = \{20, 4\} \\ C_5 = \{25, 30\} \end{array}$$

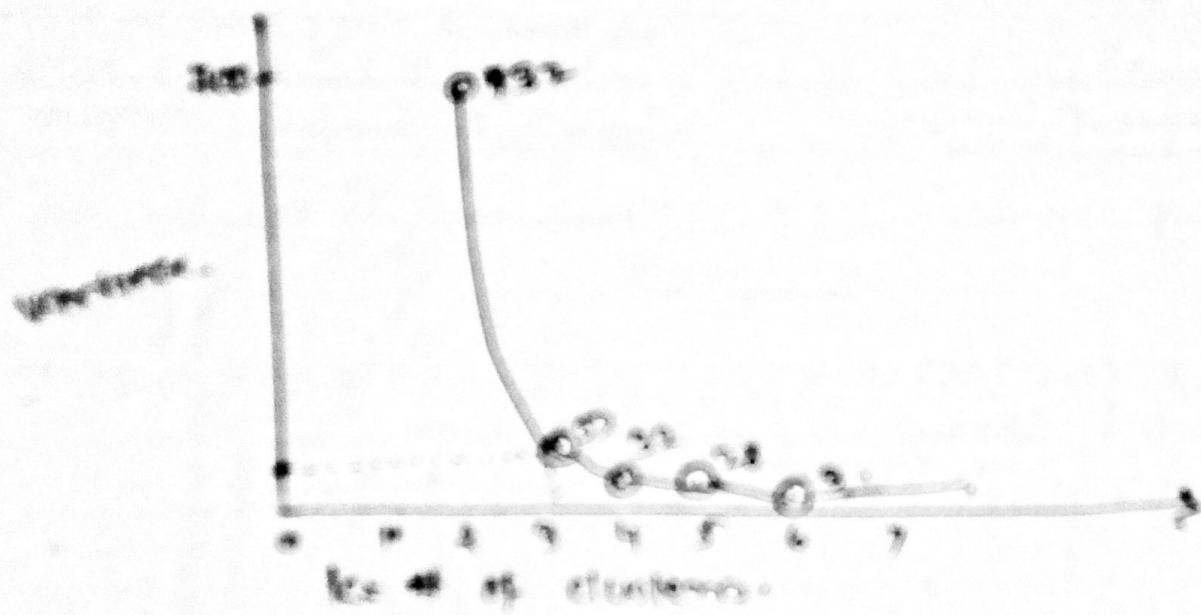
$$\begin{aligned} VAR(k=5) &= (3-3)^2 + (3-2)^2 + (3-4)^2 \\ &\quad + (10-10)^2 \\ &\quad + (12-12)^2 + (12-11)^2 \\ &\quad + (20-20)^2 + (20-25)^2 + (25-30)^2 \\ &= 0+1+1+0+0+1+0+0+25 = 28. \end{aligned}$$

$$\begin{array}{ll} k=6 & C_1 = \{3, 2, 4\} \\ C_2 = \{10, 11\} \\ C_3 = \{12, \} \\ C_4 = \{20, \} \\ C_5 = \{25, \} \\ C_6 = \{30, \} \end{array}$$

$$\begin{aligned} VAR(k=6) &= 0+1+1 \\ &\quad + 0+1 \\ &\quad + 0+0+0+0 = 3 \end{aligned}$$



Variance



No. of students.

Hence $R^2 = 0.937$, reduction in variance is not that much.
 Hence full effect variance of student's result
 i.e. It varies $\approx 8\%$.

Q. Using data on distances given below, compute silhouette coefficient for each point, each of two clusters and the overall clustering. Cluster 1 contains {P₁, P₂, P₃}, cluster 2 contains {P₃, P₄}.

	P ₁	P ₂	P ₃	P ₄
P ₁	0			
P ₂	0.1	0		
P ₃	0.65	0.70	0	
P ₄	0.55	0.60	0.30	0

→ for P₁ a = 0.1 avg dist betw P₁ & all other points in C₁.

$$b = \frac{0.65 + 0.55}{2} = 0.6$$

= avg dist betw P₁ & all other points in C₂

$$S = \frac{b - a}{\max(a, b)} = \frac{0.6 - 0.1}{\max(0.1, 0.6)} = \frac{0.5}{0.6} = 0.833$$

for P₂

$$a = 0.1$$

$$b = \frac{0.65 + 0.6}{2} = 0.625$$

$$S = \frac{b - a}{\max(a, b)} = \frac{0.625 - 0.1}{0.625} = 0.846$$

for P₃

$$a = 0.3$$

$$b = \frac{0.65 + 0.7}{2} = 0.675$$

$$S = \frac{0.675 - 0.3}{0.675} = 0.555$$

for P₄

$$a = 0.3$$

$$b = \frac{0.55 + 0.6}{2} = 0.575$$

$$S = \frac{0.575 - 0.3}{0.575} = 0.478$$

for cluster C₁ $S = \frac{0.833 + 0.846}{2} = 0.839$

for cluster C₂ $S = \frac{0.555 + 0.478}{2} = 0.516$

for overall clustering $S = \frac{0.839 + 0.516}{2} = 0.692$

- Q. How does F-score help in quantifying cluster quality ?
 Given clustering results of newspaper articles dataset given below, compute F-score for cluster representing metro & financial articles.

Cluster	Entertain	fin	foreign	Metro	Nation	Sports	Total
# 1	1	1	5	11	4	676	698
# 2	27	89	333	827	253	33	1562
# 3	126	465	8	105	16	29	749
Total	154	555	346	943	273	738	3009

$$f = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

→ for cluster 1

$$\text{class Metro. } f = \frac{2 \times \frac{11}{676} \times \frac{11}{943}}{\frac{11}{676} + \frac{11}{943}} = 0.00195$$

$$\text{class financial } f = \frac{2 \times \frac{1}{676} \times \frac{1}{555}}{\frac{1}{676} + \frac{1}{555}} = 0.00162$$

for cluster 2

$$\text{class Metro } f = \frac{2 \times \frac{827}{1562} \times \frac{827}{943}}{\frac{827}{1562} + \frac{827}{943}} = 0.660$$

$$\text{class financial } f = \frac{2 \times \frac{89}{1562} \times \frac{89}{555}}{\frac{89}{1562} + \frac{89}{555}} = 0.084$$

for cluster 3

$$\text{class Metro } f = \frac{2 \times \frac{105}{749} \times \frac{105}{943}}{\frac{105}{749} + \frac{105}{943}} = 0.124$$

$$\text{class financial } f = \frac{2 \times \frac{465}{749} \times \frac{465}{555}}{\frac{465}{749} + \frac{465}{555}} = 0.713$$