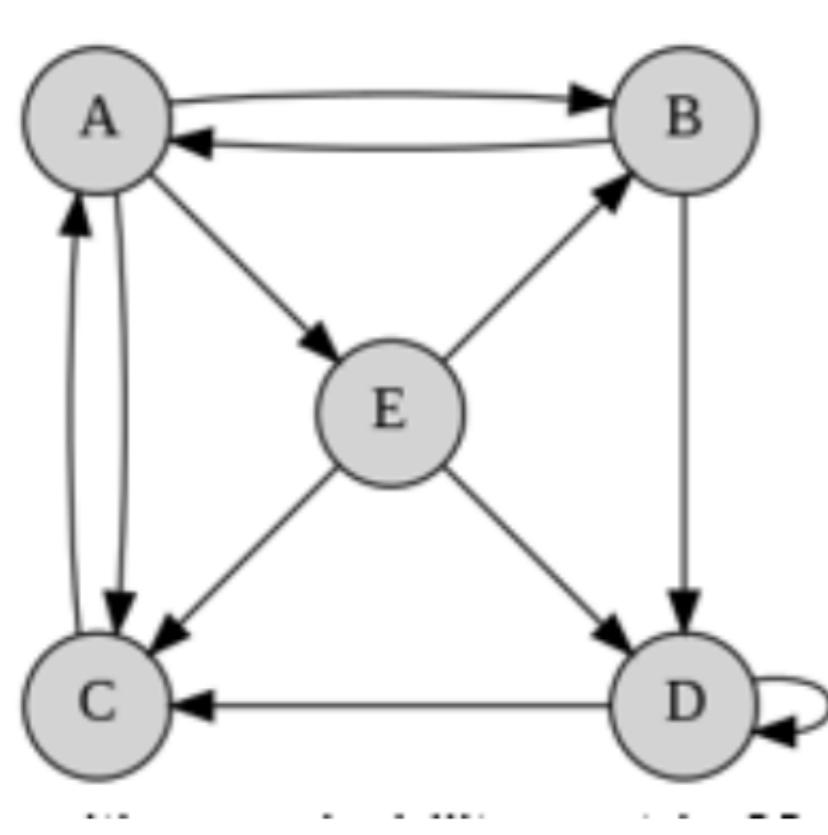


TOPIC 2 PAGERANK

- a) (5 points) Compute the transition probability matrix M of going from each node to any other adjacent node, such that the probabilities sum up to one in a column-wise way. Fill in the placeholders (?) with the appropriate values in the transition probability M matrix shown below.



Transition probability matrix M					
	A	B	C	D	E
A	0	1/2	1	0	0
B	1/3	0	0	0	1/3
C	1/3	0	0	1/2	1/3
D	0	1/2	0	1/2	1/3
E	1/3	0	0	0	0

(+) ↓

Iteration 2

Transition probability matrix M

	A	B	C	D	E
A	0	1/2	1	0	0
B	1/3	0	0	0	1/3
C	1/3	0	0	1/2	1/3
D	0	1/2	0	1/2	1/3
E	1/3	0	0	0	0

$$M \cdot V_1 = \begin{bmatrix} 0 & 1/2 & 1 & 0 & 0 \\ 1/3 & 0 & 0 & 0 & 1/3 \\ 1/3 & 0 & 0 & 1/2 & 1/3 \\ 0 & 1/2 & 0 & 1/2 & 1/3 \\ 1/3 & 0 & 0 & 0 & 0 \end{bmatrix} \times \begin{bmatrix} 0,285 \\ 0,143 \\ 0,228 \\ 0,228 \\ 0,086 \end{bmatrix} = \begin{bmatrix} 0,2995 \\ 0,1236 \\ 0,2516 \\ 0,22876 \\ 0,095 \end{bmatrix}$$

$$V_2' = 0,85 \cdot M \cdot V_1' + \begin{bmatrix} 0,03 \\ 0,03 \\ 0,03 \\ 0,03 \\ 0,03 \end{bmatrix} = \begin{bmatrix} 0,284 \\ 0,135 \\ 0,243 \\ 0,223 \\ 0,770 \end{bmatrix} \leftarrow \text{truncated}$$

- c) (4 points) Compute the topic-sensitive PageRank for the initial graph, assuming the teleport set is $S = \{C, D\}$ and $\beta = 0.8$. All nodes (A – E) share the same initial probabilities. Calculate by hand the values of the column vector v' after the 2nd iteration and fill in the placeholders (?) with the correct values. During all calculations truncate the intermediate and final results to 3 decimal digits. (do not round up or down the numbers)

Topic-Sensitive PageRank

teleportation restricted to C and D

$$\text{NEW initial probability vector } v_0 = [0, 0, 1/2, 1/2, 0]^T$$

Iteration 1

Transition probability matrix M

	A	B	C	D	E
A	0	1/2	1	0	0
B	1/3	0	0	0	1/3
C	1/3	0	0	1/2	1/3
D	0	1/2	0	1/2	1/3
E	1/3	0	0	0	0

$$\times \begin{bmatrix} 1/5 \\ 1/5 \\ 1/5 \\ 1/5 \\ 1/5 \end{bmatrix} = \begin{bmatrix} 0,3 \\ 0,133 \\ 0,233 \\ 0,266 \\ 0,066 \end{bmatrix}$$

$$V_1' = 0,8 \cdot \begin{bmatrix} 0,3 \\ 0,133 \\ 0,233 \\ 0,266 \\ 0,066 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0,24 \\ 0,106 \\ 0,286 \\ 0,312 \\ 0,08 \end{bmatrix}$$

Sx5

Sx1

 $\beta = 0,8 \rightarrow 80\% \text{ follows links}$

20% teleports

$$v' = 0,8 \cdot M \cdot v + 0,2 \cdot \begin{bmatrix} 0 \\ 0 \\ 1/2 \\ 1/2 \\ 0 \end{bmatrix} = 0,8 \cdot M \cdot v + \begin{bmatrix} 0 \\ 0 \\ 0,1 \\ 0,1 \\ 0 \end{bmatrix}$$

Iteration 2

Transition probability matrix M

	A	B	C	D	E
A	0	1/2	1	0	0
B	1/3	0	0	0	1/3
C	1/3	0	0	1/2	1/3
D	0	1/2	0	1/2	1/3
E	1/3	0	0	0	0

$$\times \begin{bmatrix} 0,24 \\ 0,106 \\ 0,286 \\ 0,312 \\ 0,08 \end{bmatrix} = \begin{bmatrix} 0,339 \\ 0,0976 \\ 0,25416 \\ 0,28716 \\ 0,08 \end{bmatrix}$$

Sx5

Sx1

$$V_2' = 0,8 \cdot \begin{bmatrix} 0,339 \\ 0,0976 \\ 0,25416 \\ 0,28716 \\ 0,08 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0,271 \\ 0,078 \\ 0,303 \\ 0,281 \\ 0,064 \end{bmatrix}$$

- d) (2 points) Which node presents the highest rank in each implementation?

Answer:

PageRank with teleportation: Node A presents the highest rank (0.284)

Topic Sensitive PageRank: Node C presents the highest rank (0.303)

PAGERANK with teleportation

$$V_2' = \begin{bmatrix} 0,284 \\ 0,135 \\ 0,213 \\ 0,223 \\ 0,770 \end{bmatrix} \leftarrow \text{node A}$$

TOPIC-SENSITIVE PAGERANK

$$V_2' = \begin{bmatrix} 0,271 \\ 0,078 \\ 0,303 \\ 0,281 \\ 0,064 \end{bmatrix} \leftarrow \text{node C}$$

Topic 2 PageRank

(b)

```
[1]: import numpy as np
import pandas as pd

[2]: def truncate(matrix, decimals=3):
    return np.floor(matrix * 10**decimals) / 10**decimals

[3]: m1= np.array([0, 1/3, 1/3, 0 , 1/3])
m2 = np.array([1/2,0,0,1/2,0])
m3 = np.array([1,0,0,0,0])
m4 = np.array([0,0,1/2,1/2,0])
m5 = np.array([0,1/3,1/3,1/3,0])

[4]: v0=np.array([0.2, 0.2, 0.2, 0.2, 0.2])

[5]: M = np.column_stack([m1, m2, m3, m4, m5])
M

[5]: array([[0.          ,  0.5        ,  1.          ,  0.          ,  0.          ],
       [0.33333333,  0.          ,  0.          ,  0.33333333],
       [0.33333333,  0.          ,  0.5        ,  0.33333333],
       [0.          ,  0.5        ,  0.          ,  0.5        ,  0.33333333],
       [0.33333333,  0.          ,  0.          ,  0.          ,  0.        ]])

[6]: Mv0=M @ v0
Mv0

[6]: array([0.3        ,  0.13333333,  0.23333333,  0.26666667,  0.06666667])

[7]: v1=(0.85*Mv0)+np.array([0.03,0.03,0.03,0.03,0.03])
v1_truncated=truncate(v1)
v1_truncated

[7]: array([0.285, 0.143, 0.228, 0.256, 0.086])

[8]: Mv1 = M @ v1_truncated
Mv1

[8]: array([0.2995     ,  0.12366667,  0.25166667,  0.22816667,  0.095      ])

[9]: v2=0.85*Mv1+np.array([0.03,0.03,0.03,0.03,0.03])
v2_truncated = truncate(v2)
v2_truncated

[9]: array([0.284, 0.135, 0.243, 0.223, 0.11  ])
```

(c)

```
[10]: t = 0.2*np.array([0, 0, 1/2, 1/2, 0])
t

[10]: array([0. ,  0. ,  0.1,  0.1,  0. ])

[11]: b=0.8

[12]: Mv0 = M @ v0
Mv0

[12]: array([0.3        ,  0.13333333,  0.23333333,  0.26666667,  0.06666667])

[13]: v1=b*Mv0+t
v1_truncated = truncate(v1)
v1_truncated

[13]: array([0.24        ,  0.106,  0.286,  0.313,  0.053])

[14]: Mv1 = M @ v1_truncated
Mv1

[14]: array([0.339        ,  0.09766667,  0.25416667,  0.22716667,  0.08      ])

[15]: v2 = b*Mv1+t
v2

[15]: array([0.2712     ,  0.07813333,  0.30333333,  0.28173333,  0.064      ])

[16]: truncate(v2)

[16]: array([0.271,  0.078,  0.303,  0.281,  0.063])

[17]: truncate(np.array([0.2712,  0.07813333,  0.30333333,  0.28173333,  0.064]))

[17]: array([0.271,  0.078,  0.303,  0.281,  0.064])

[18]: truncate(v2)

[18]: array([0.271,  0.078,  0.303,  0.281,  0.063])

[19]: print(v2)
print(v2[4])

[0.2712  0.07813333  0.30333333  0.28173333  0.064      ]
0.06399999999999999
```

Interestingly enough, the 5th element in the v2 array, is not 0.064, but 0.0639999999999999. Therefore, when we truncate, it should actually be 0.063.

TOPIC 3

Topic 3: PCY Algorithm

(15 total points) In the following table, a collection of 15 transactions is given. Each transaction consists of a set of items generated from a total of nine (9) distinct items, enumerated from 1 to 9. The minimum support threshold is set to 6.

Transaction	Items	Transaction	Items	Transaction	Items
T1	{1, 2, 5}	T6	{1, 5, 6, 9}	T11	{1, 2, 6, 9}
T2	{2, 3, 6}	T7	{2, 3, 4}	T12	{1, 3, 7}
T3	{3, 4, 5, 7}	T8	{2, 4, 5, 7}	T13	{2, 5, 8}
T4	{1, 3, 5}	T9	{3, 5, 7, 8}	T14	{3, 4, 9}
T5	{2, 4, 7, 8}	T10	{2, 4, 8}	T15	{4, 6, 7, 9}

Answer:

Count how many times they appear

Item	1	2	3	4	5	6	7	8	9
Support	5	8	7	7	7	4	6	4	4

$\sum = 59$

Count how many times both items (i, j) appear together

Support of pairs (i, j)	j								
	2	3	4	5	6	7	8	9	
1	2	2	0	3	2	1	0	2	
2		2	4	3	2	2	3	1	
3			3	3	1	3	1	1	
4				2	1	4	2	2	
5					1	3	2	1	
6						1	0	3	
7							2	1	
8								0	

- a) (2 Points) Compute the support for each item and each pair of items and fill in the placeholders (?) with the correct values in the tables below.

Answer:

Bucket	0	1	2	3	4	5	6	7	8
Support	?	?	6	?	?	7	?	4	?

generally $n \bmod 9 = u$ $\forall n \in \{m \in \mathbb{N} / m < 9\}$

$$h(i,j) = (i \times j) \bmod 9$$

Pair	Support	Bucket
(1,2)	2	$(1 \cdot 2) \bmod 9 = 2 \bmod 9 = 2$
(1,3)	2	$3 \bmod 9 = 3$
(1,4)	0	$4 \bmod 9 = 4$
(1,5)	3	5
(1,6)	2	6
(1,7)	1	7
(1,8)	0	8
(1,9)	2	0

Pair	Support	Bucket
(2,3)	2	$(2 \cdot 3) \bmod 9 = 6 \bmod 9 = 6$
(2,4)	4	$(2 \cdot 4) \bmod 9 = 8 \bmod 9 = 8$
(2,5)	3	$(2 \cdot 5) \bmod 9 = 10 \bmod 9 = 1$
(2,6)	2	$(2 \cdot 6) \bmod 9 = 12 \bmod 9 = 3$
(2,7)	2	$(2 \cdot 7) \bmod 9 = 14 \bmod 9 = 5$
(2,8)	3	$(2 \cdot 8) \bmod 9 = 16 \bmod 9 = 7$
(2,9)	7	$(2 \cdot 9) \bmod 9 = 18 \bmod 9 = 0$

Pair	Support	Bucket
(3,4)	3	$(3 \cdot 4) \bmod 9 = 12 \bmod 9 = 3$
(3,5)	3	$(3 \cdot 5) \bmod 9 = 15 \bmod 9 = 6$
(3,6)	1	$(3 \cdot 6) \bmod 9 = 18 \bmod 9 = 0$
(3,7)	3	$(3 \cdot 7) \bmod 9 = 21 \bmod 9 = 3$
(3,8)	1	$(3 \cdot 8) \bmod 9 = 24 \bmod 9 = 6$
(3,9)	1	$(3 \cdot 9) \bmod 9 = 27 \bmod 9 = 0$
(4,5)	2	$(4 \cdot 5) \bmod 9 = 20 \bmod 9 = 2$
(4,6)	7	$(4 \cdot 6) \bmod 9 = 24 \bmod 9 = 6$

Pair	Support	Bucket
(5,6)	0	1
(5,7)	2	5
(5,8)	1	3
(5,9)	3	8
(6,7)	2	4
(6,8)	0	0
(6,9)	1	6

Pair	Support	Bucket
(7,8)	0	3
(7,9)	1	0
(8,9)	0	0

Bucket	Total Support
0	12
1	7
2	6
3	11
4	9
5	7
6	10
7	4
8	7

Answer:

Bucket	0	1	2	3	4	5	6	7	8
Support	12	7	6	11	2	7	10	4	7

- c) (3 points) Which buckets are frequent? Replace the placeholders (?) with the buckets that are frequent.

Answer:

Frequent buckets are: 0, ?, ?, ?, ?, ?, ?.

minimum support threshold = 6 therefore, frequent buckets: 0, 1, 2, 3, 4, 5, 6, 7, 8

(d) (4 points) Which pairs are finally counted in the 2nd pass of PCY? Replace the placeholders $\{?, ?\}$ with the correct pairs.

Answer:

Candidate pairs are: $\{2, 3\}, \{2, 4\}, \{?, ?\}, \{?, ?\}, \{?, ?\}, \{?, ?\}, \{?, ?\}, \{?, ?\}, \{?, ?\}, \{?, ?\}$

2nd pass candidates are only pairs that: a) have frequent items
b) are bucketed into frequent buckets

Table #1

Item	1	2	3	4	5	6	7	8	9
Support	5	8	7	7	7	4	6	4	4

We rule out infrequent items.

We only consider pairs with 2, 3, 4, 5, and 7

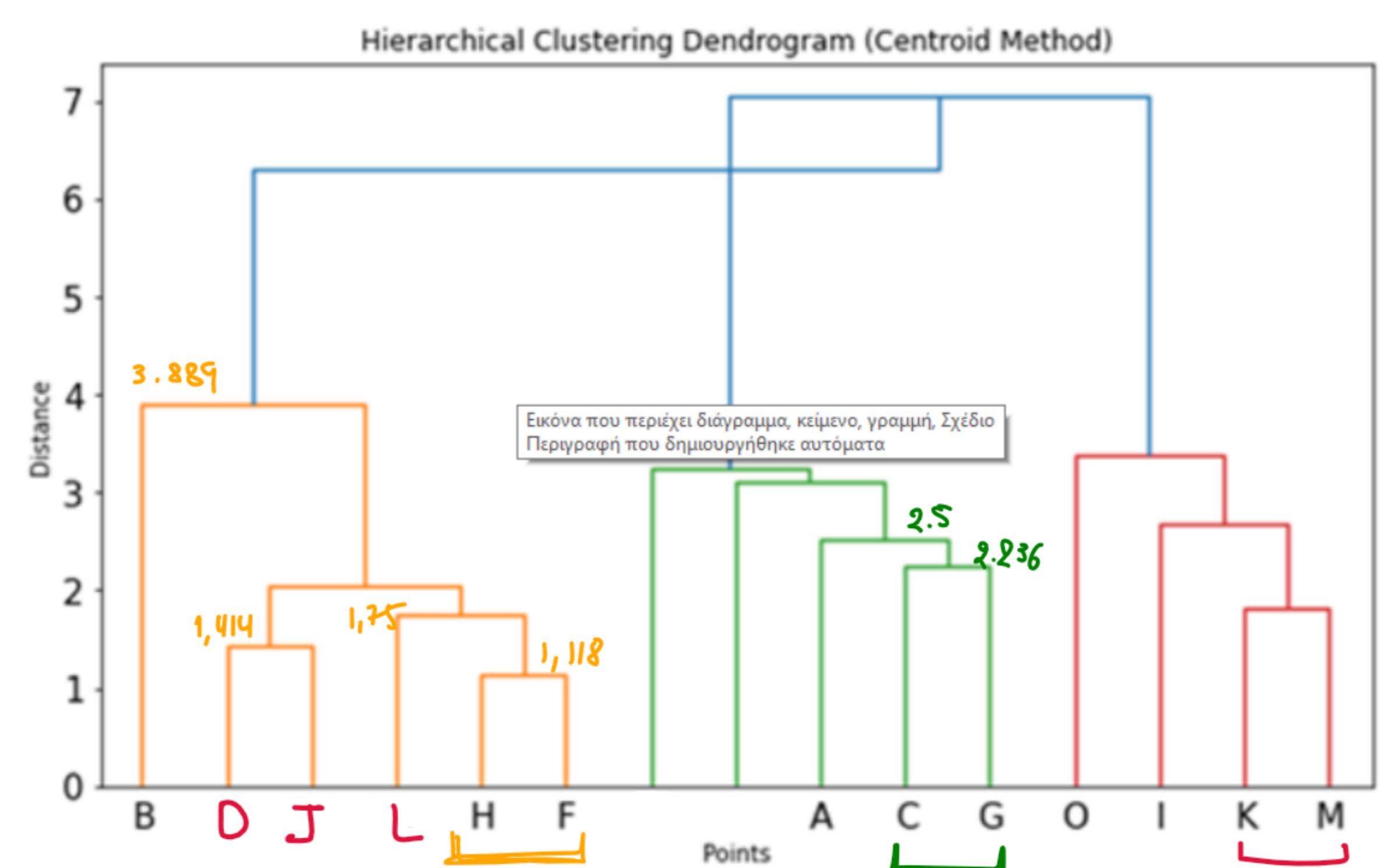
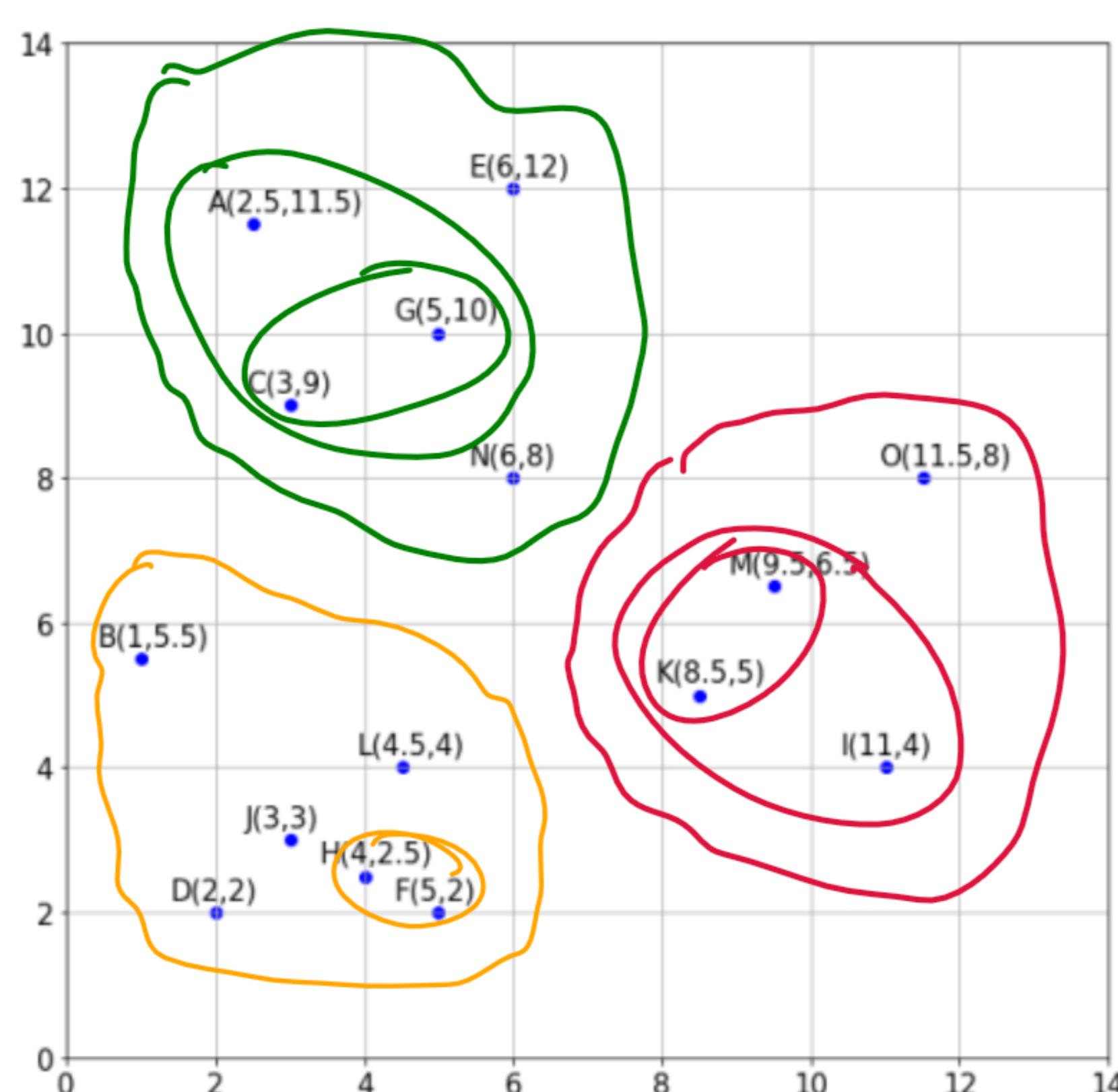
Pair	Bucket	frequent?
(2,3)	6	✓
(2,4)	8	✓
(2,5)	1	✓
(2,7)	5	✓
(3,4)	3	✓
(3,5)	6	✓
(3,7)	3	✓
(4,5)	9	✓
(4,7)	1	✓
(5,7)	8	✓

All these pairs are finally counted in the 2nd pass of PCY.

TOPIC 4 CURE ALGORITHM

Topic 4: CURE Algorithm

(15 total points) In this topic, a set of points is clustered using the CURE algorithm. Consider the following 15 points in a two-dimensional Euclidean space which will be clustered using CURE:



(a) (5 points) Perform hierarchical clustering of the given points, assuming clusters are represented by their centroid and at each step the clusters with the closest centroids are merged. Assume that after hierarchically clustering the points, three clusters are created. Fill in the missing points, indicated by the numbered placeholders $1, 2, ?$ etc in the following dendrogram with the correct values, to specify the three formed clusters.

$$\text{New centroid } (H, F) \rightarrow C_{HF} \left(\frac{4+5}{2}, \frac{2.5+5}{2} \right)$$

$$C_{HF} (4.5, 3.75)$$

$$d(L, C_{HF}) = \sqrt{(4.5 - 4.5)^2 + (4 - 3.75)^2} = \sqrt{0 + (1.25)^2} = 1.25$$

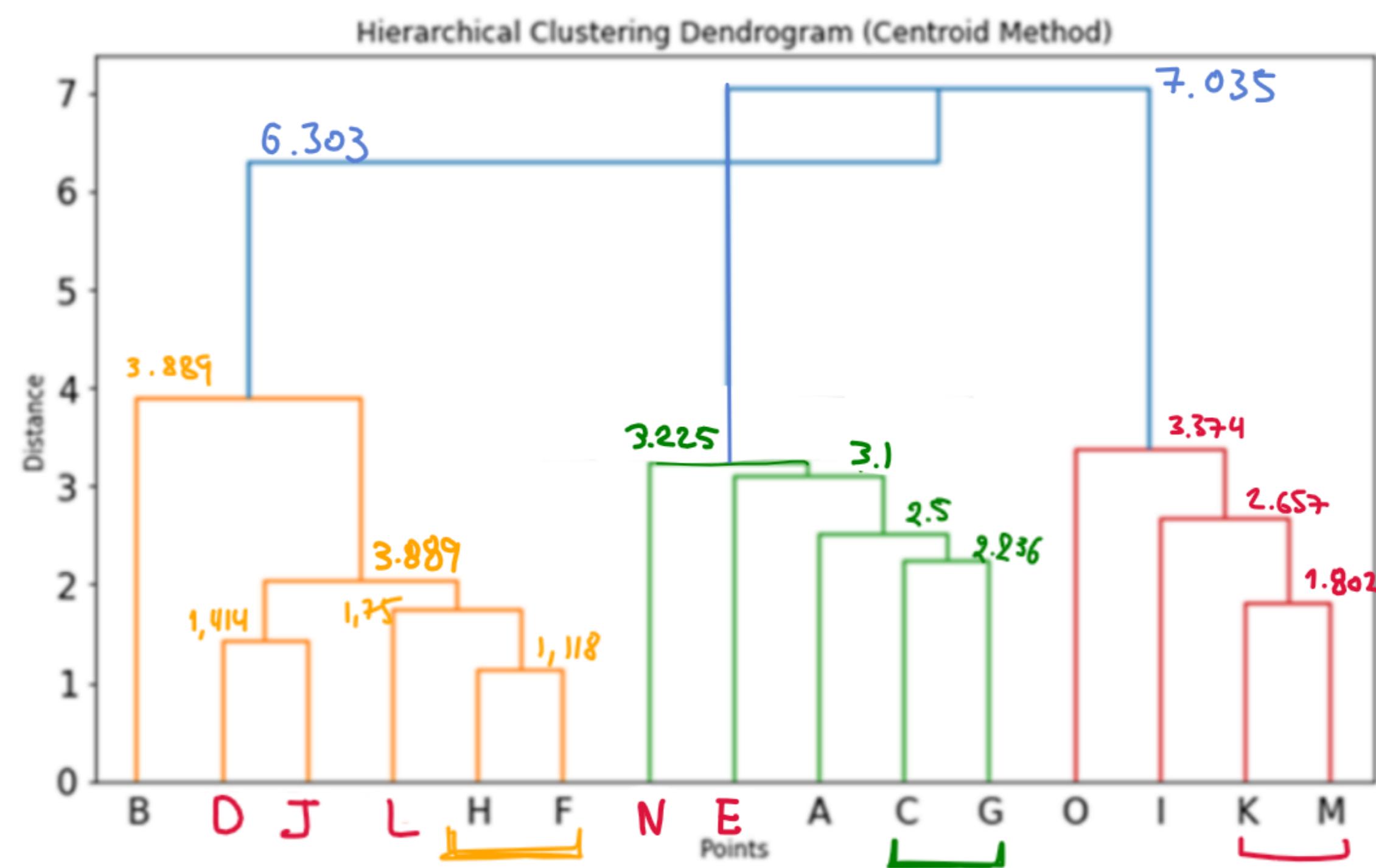
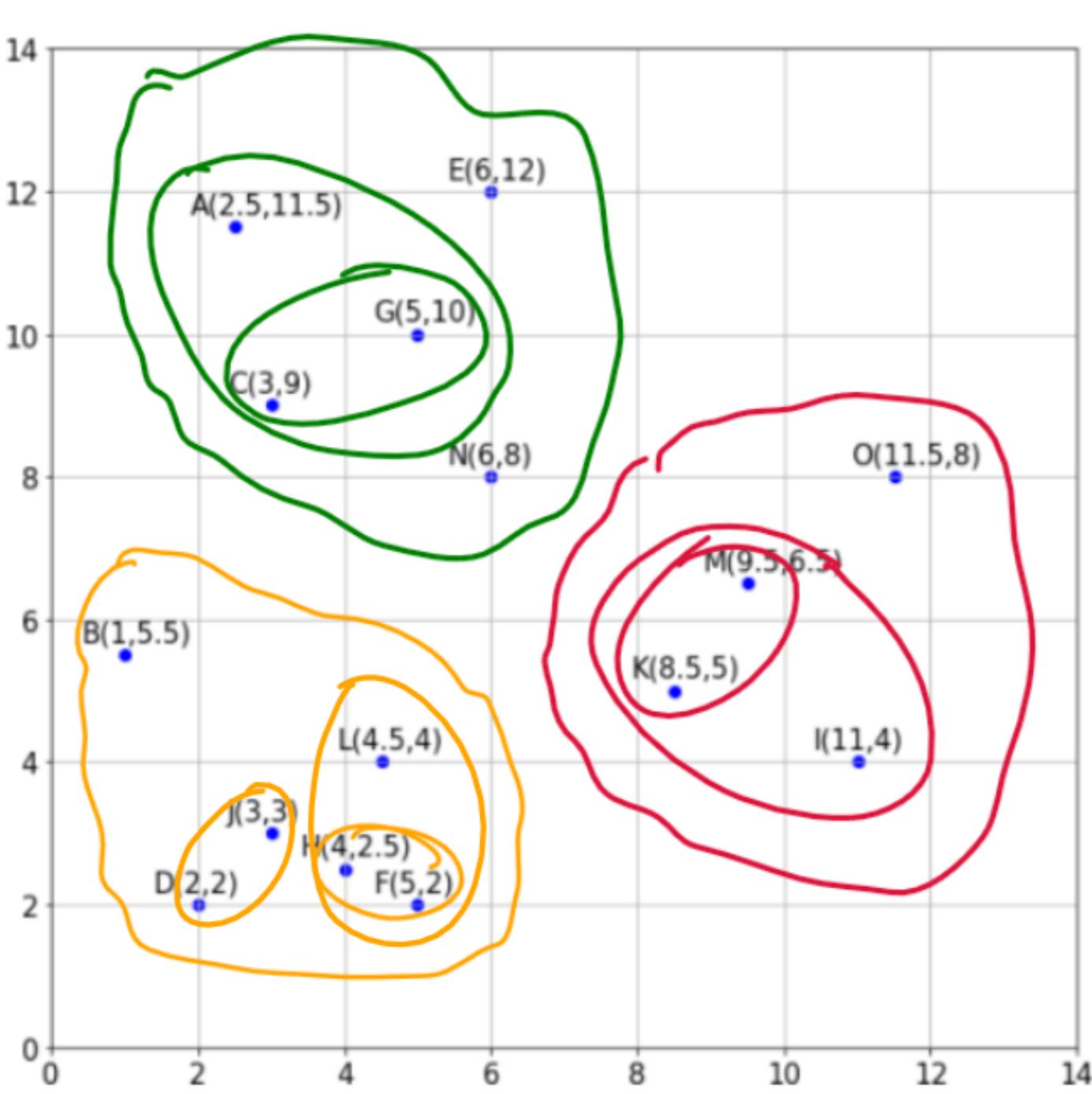
$$d(H, F) = \sqrt{(5 - 4)^2 + (2.5 - 2)^2} = \sqrt{1 + 0.25} = \sqrt{1.25} \approx 1.118$$

$$d(D, J) = \sqrt{(3 - 2)^2 + (3 - 2)^2} = \sqrt{2} \approx 1.414$$

$$d(J, L) = \sqrt{(4.5 - 3)^2 + (4 - 3)^2} = \sqrt{2.25 + 1} = \sqrt{3.25} \approx 1.80$$

$$d(D, L) > d(J, L) > d(D, J)$$

Therefore $1 \rightarrow D$ $2 \rightarrow J$ $3 \rightarrow L$



$$\text{New Centroid } (A, C, G) \rightarrow C_{\text{Acc}} \left(\frac{2.5+5+3}{3}, \frac{11.5+9+10}{3} \right) \\ C_{\text{Acc}} (3.5, 10.2)$$

$$d(N, C_{\text{Acc}}) = \sqrt{(6-3.5)^2 + (8-10.2)^2} = \sqrt{6.25 + 4.84} = \sqrt{11.09} \approx 3.3$$

$$d(E, C_{\text{Acc}}) = \sqrt{6.25 + (12-10.2)^2} = \sqrt{6.25 + 1.8^2} = \sqrt{9.49} \approx 3$$

$$d(E, C_{\text{Acc}}) < d(N, C_{\text{Acc}})$$

$$d(C, G) = 2.24$$

Therefore, $4 \rightarrow N$
 $5 \rightarrow E$

- b) (4 points) Having formed three clusters, select three points from the first cluster (cluster C1) and the three points from the second cluster (cluster C2) that will be the representative points for each cluster. Fill in the placeholders (?) to specify the representative points. These points should be chosen to be as far from one another as possible. For cluster C1 use F as the initial random point; for cluster C2 use A as the initial random point. Assume that for cluster C3 the selected representative points are {I, K, O} (required for the next question).

Answer:

representative points for cluster C1: {F, B, D}

representative points for cluster C2: {A, N, E}

start F

$$d(F, B) = 5.315 \quad \leftarrow d_{\max}$$

$$d(F, D) = 3$$

$$d(F, J) = 2.236$$

$$d(F, H) = 7.118$$

$$d(F, L) = 2.061$$

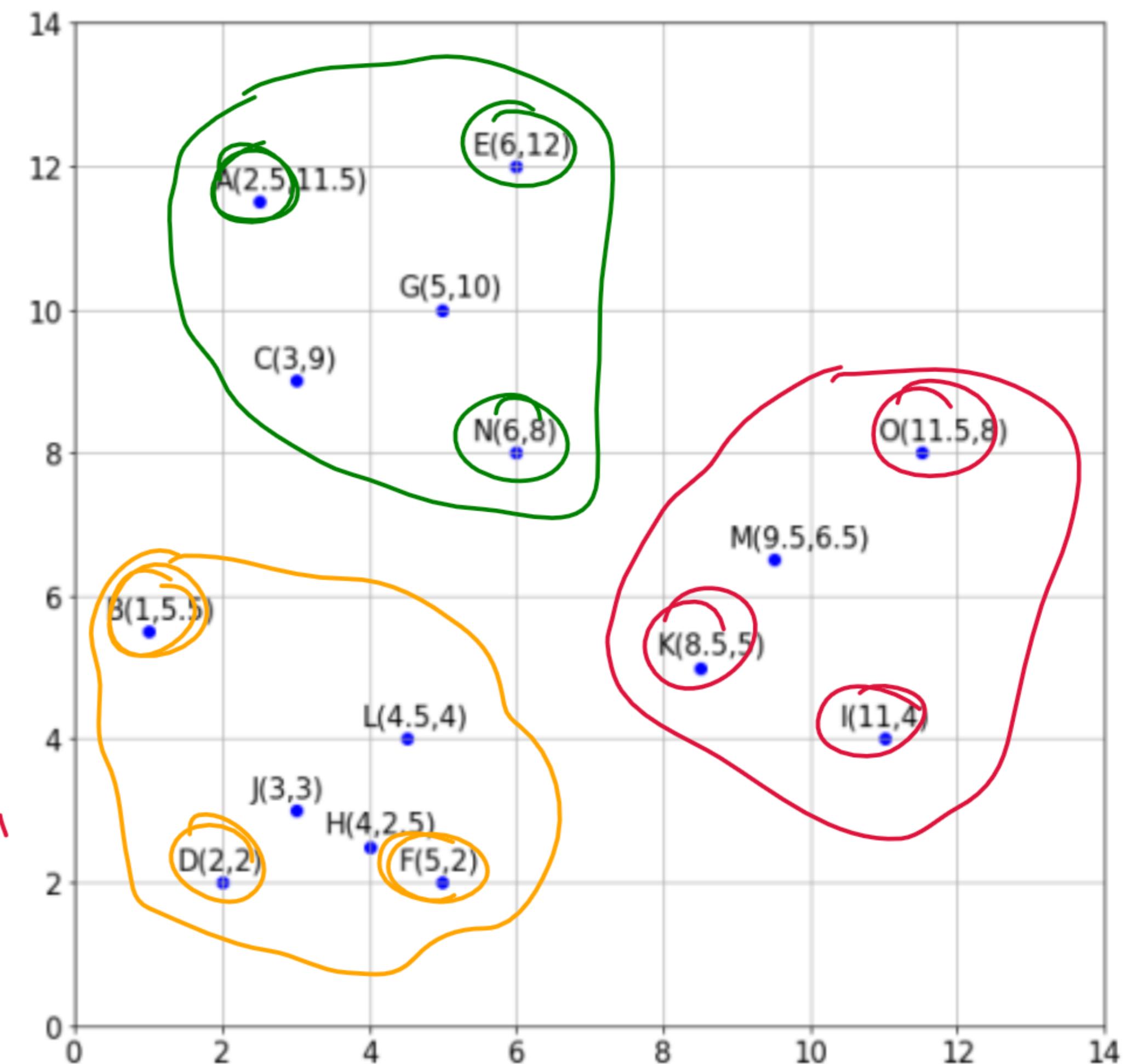
then B

$$d(B, D) = 3.64 \quad \#3$$

$$d(B, J) = 3.201 \quad \text{too close to both}$$

$$d(B, H) = 4.243 \quad \#1$$

$$d(B, L) = 3.807 \quad \#2$$



$$* \frac{d(B, H)}{d(F, H)} = 3.795$$

$$+ \frac{d(B, L)}{d(F, L)} = 1.847$$

$$* \frac{d(B, D)}{d(F, D)} = 1.213$$

D is almost as far away
from F as it is from B

start A

$$d(A, C) = 2.549$$

$$d(A, G) = 2.915$$

$$d(A, E) = 3.536$$

$$d(A, N) = 4.95 \quad \leftarrow d_{\max}$$

then N

$$d(N, C) = 3.162$$

$$d(N, G) = 2.236 \quad \text{too close to both}$$

$$d(N, E) = 4 \quad \leftarrow \text{here the choice is obvious.}$$

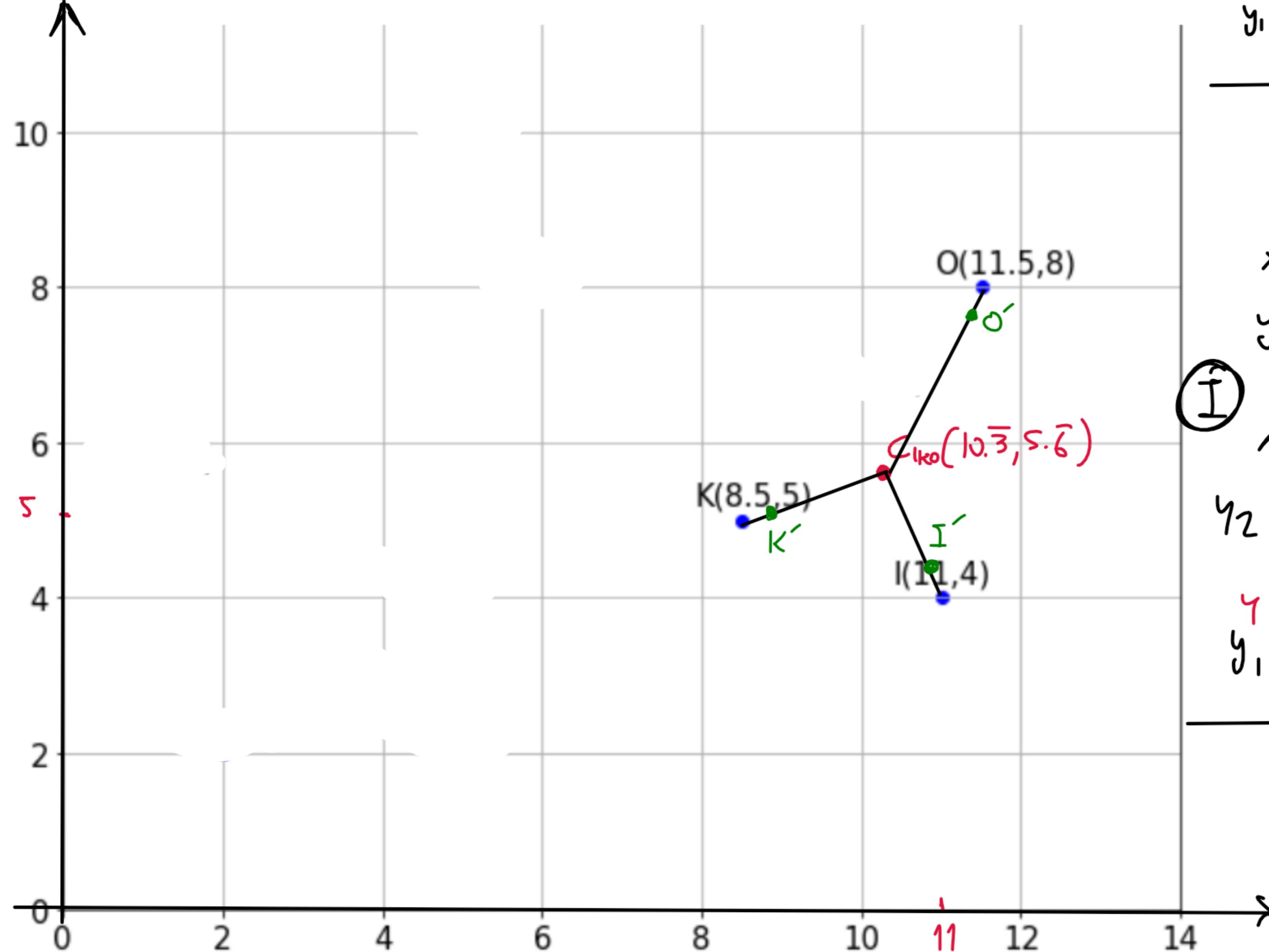
E is further away from A and N
than any other point

(c) (3 points) Move each of the representative points of cluster C3 only 20% closer to the centroid of the cluster it belongs to. Give the new positions of the representative points for cluster C3 and fill in the placeholders (?). During all calculations truncate the intermediate and final results to 3 decimal digits (do not round up or down the numbers). The cluster centroid is computed over the selected representative points.

Answer:

Representative points for cluster C3: I(?, ?), K(?, ?), O(?, ?)

I(11, 4)
K(8.5, 5)
O(11.5, 8)
 $C_{1ko} (10.333, 5.666)$

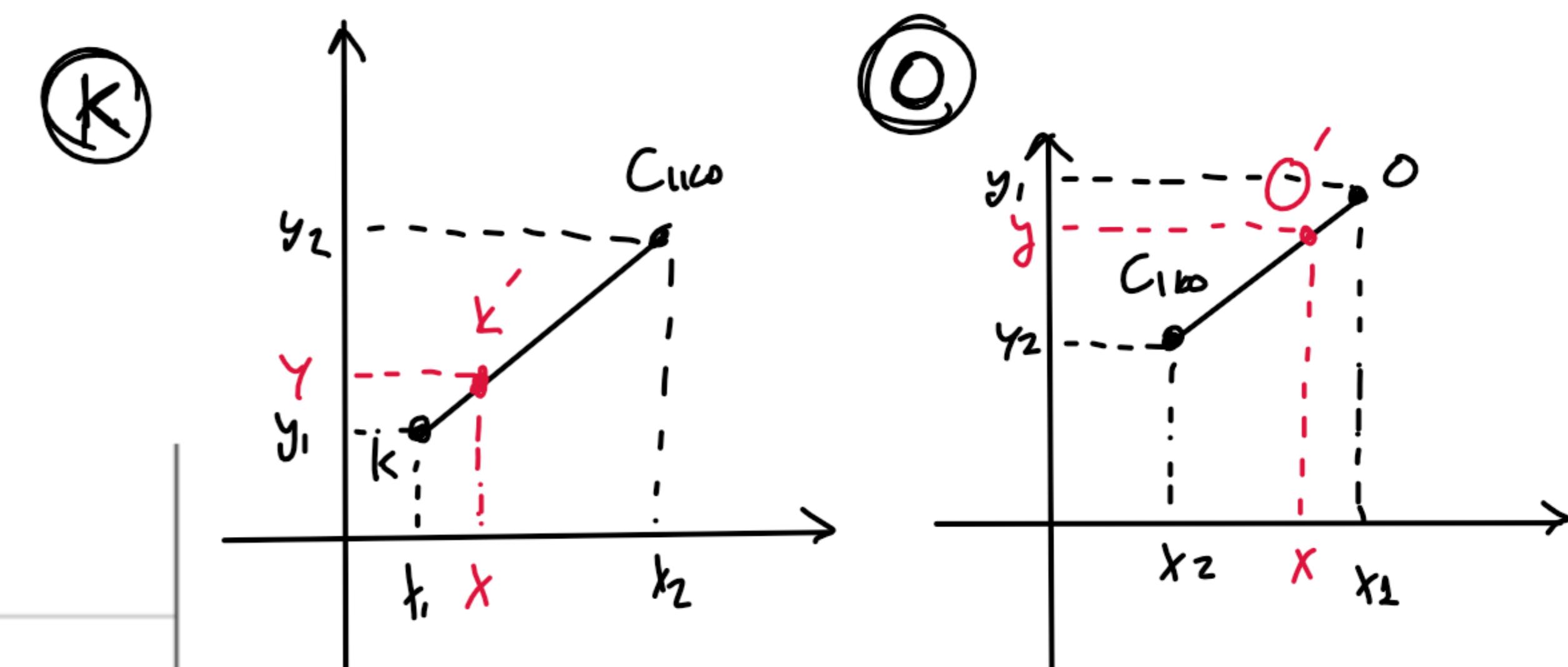


Centroid C_3

$$x = \frac{11 + 8.5 + 11.5}{3} = \frac{31}{3} = 10.333$$

$$y = \frac{4 + 5 + 8}{3} = \frac{17}{3} = 5.666$$

$C_{1ko} (10.333, 5.666)$

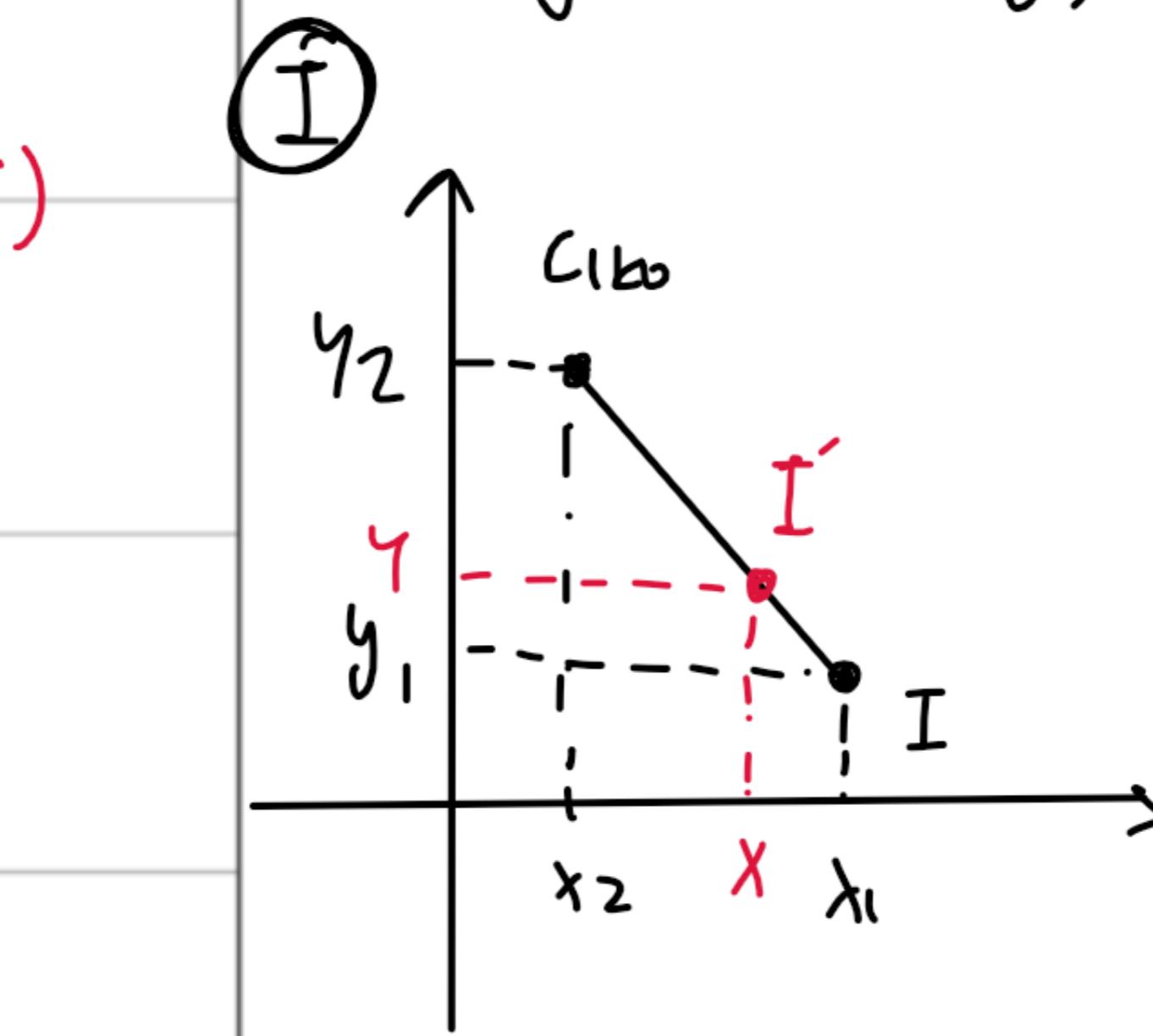


$$x = x_1 + 0.2(x_2 - x_1)$$

$$y = y_1 + 0.2(y_2 - y_1)$$

$$x = x_1 + 0.2(x_2 - x_1)$$

$$y = y_1 + 0.2(y_2 - y_1)$$



$$x = x_1 + 0.2(x_2 - x_1)$$

$$y = y_1 + 0.2(y_2 - y_1)$$

K' $x_{k'} = x_k + 0.2(x_c - x_k) = 8.5 + 0.2 \cdot (10.333 - 8.5) = 8.866$
 $y_{k'} = y_k + 0.2 \cdot (y_c - y_k) = 5 + 0.2 \cdot (5.666 - 5) = 5.133$

O' $x_{o'} = x_o + 0.2(x_c - x_o) = 11.266$
 $y_{o'} = y_o + 0.2(y_c - y_o) = 7.533$

I' $x_{i'} = x_i + 0.2(x_c - x_i) = 10.866$
 $y_{i'} = y_i + 0.2(y_c - y_i) = 4.333$

K'(8.866, 5.133)

O'(11.266, 7.533)

I'(10.866, 4.333)

(d) (3 points) Calculate the minimum possible value for threshold distance d that results in 3 clusters (that is, no pair of clusters are merged), and fill in the placeholder (?). Truncate the distances between representative points to 3 decimal digits (do not round up or down the numbers). Use three decimal digits precision for d .

Answer:

$$d = ?$$

$$C_1 : \{F, B, D\}$$

$$C_2 : \{A, N, E\}$$

$$C_3 : \{I, K, O\}$$

C_1

$$F(5, 2)$$

$$B(1, 5.5)$$

$$D(2, 2)$$

C_2

$$A(2.5, 11.5)$$

$$N(6, 8)$$

$$E(6, 12)$$

Compute the centroids of C_1 and C_2

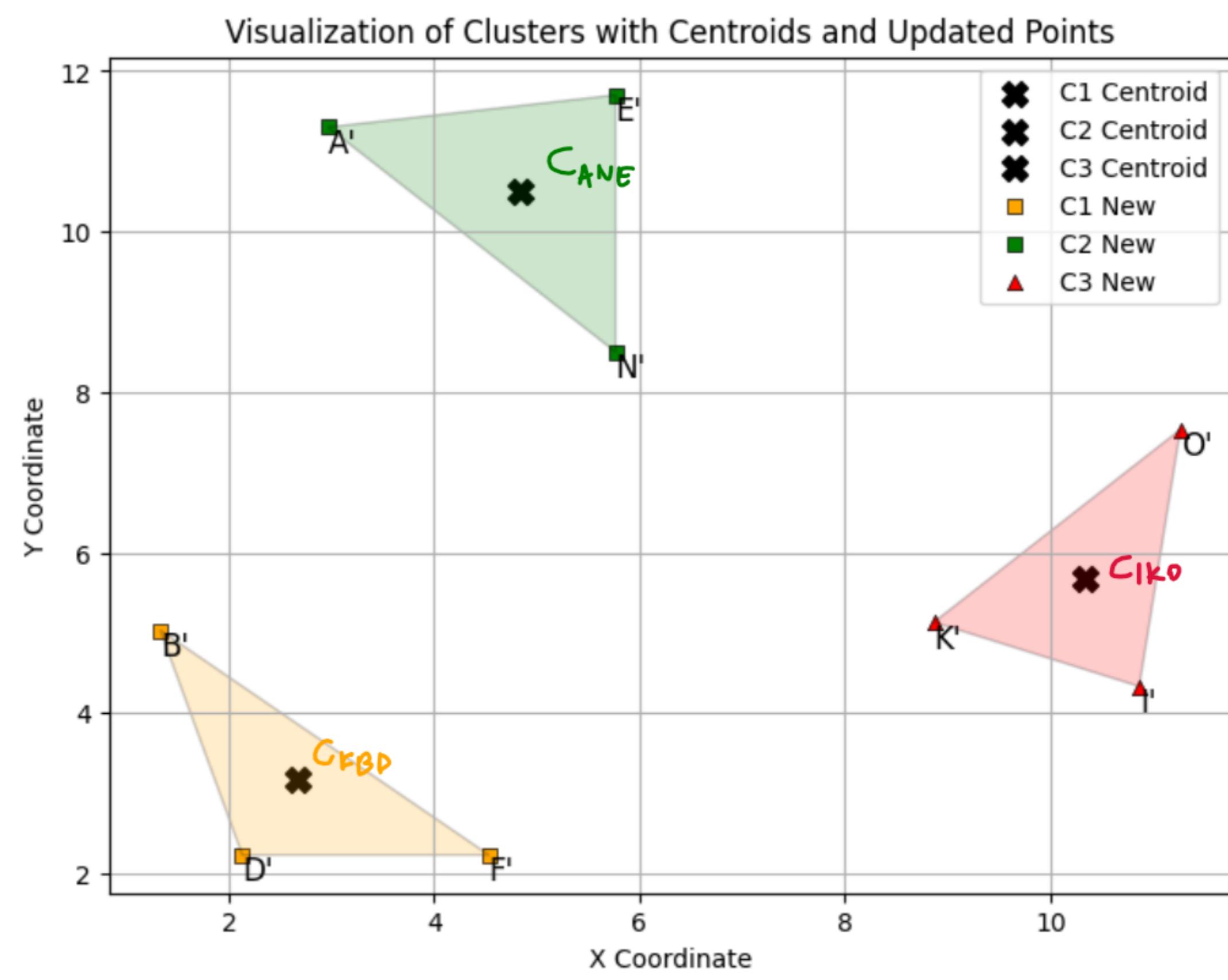
$$C_{FB'D} = \left(\frac{5+1+2}{3}, \frac{2+5.5+2}{3} \right) = (2.666, 3.166)$$

$$C_{ANE} = (4.833, 10.5)$$

Applying $X_{\text{new}} = X_{\text{old}} + 0.2 \cdot (X_{\text{centroid}} - X_{\text{old}})$ we calculate the coordinates of the new points.

$$Y_{\text{new}} = Y_{\text{old}} + 0.2 \cdot (Y_{\text{centroid}} - Y_{\text{old}})$$

$$C_1 \begin{cases} F'(4.533, 2.233) \\ B'(1.373, 5.033) \\ D'(2.133, 2.233) \end{cases} \quad C_2 \begin{cases} A'(2.966, 11.300) \\ N'(5.766, 8.5) \\ E'(5.766, 11.700) \end{cases}$$



Calculate pairwise distances between C_1 and C_2 (I will use F, B, D, A etc instead of F', B', D', A' etc)

$$\begin{array}{lll} d(F, A) = 9.901 & d(B, A) = 6.476 & d(D, A) = 9.105 \\ d(F, N) = 6.387 & d(B, N) = 5.627 & d(D, N) = 7.243 \\ d(F, E) = 9.546 & d(B, E) = 8.006 & d(D, E) = 10.14 \end{array}$$

minimum distance

$$d_{\min}(C_1, C_2) = 5.627$$

Calculate pairwise distances between C_1 and C_3

$$\begin{array}{lll} d(F, I) = 6.672 & d(B, I) = 9.558 & d(D, I) = 8.981 \\ d(F, K) = 5.213 & d(B, K) = 7.533 & d(D, K) = 7.33 \\ d(F, O) = 8.568 & d(B, O) = 10.242 & d(D, O) = 10.559 \end{array}$$

minimum distance

$$d_{\min}(C_1, C_3) = 5.213$$

Calculate pairwise distances between C_2 and C_3

$$\begin{array}{lll} d(A, I) = 10.533 & d(N, I) = 6.585 & d(E, I) = 8.96 \\ d(A, K) = 8.534 & d(N, K) = 4.576 & d(E, K) = 7.261 \\ d(A, O) = 9.114 & d(N, O) = 5.584 & d(E, O) = 6.9 \end{array}$$

minimum distance

$$d_{\min}(C_2, C_3) = 4.576$$

To prevent merging, the smallest distance between any two clusters must be greater than 4.576