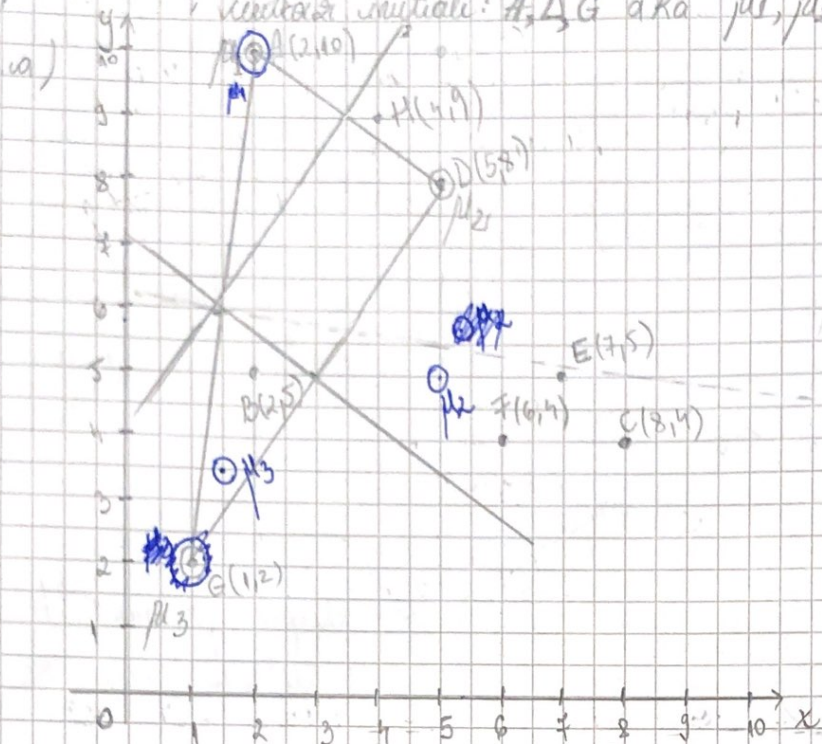


Exerc:

37. $A(2,10)$ $B(2,5)$ $C(8,4)$ $D(5,8)$ $E(7,5)$ $F(6,4)$ $G(1,2)$ $H(4,9)$

prima iterație:

centrați inițiali: A, D, G aka μ_1, μ_2, μ_3



Se observă că punctele E, F, C sunt de aceeași parte a mediatoarei (GD) ca și centrul μ_2 . Iar punctul H este de aceeași parte a mediatoarei (A,D) ca și centrul μ_3 .

\Rightarrow cluster sunt: $\overset{C_1}{\{A(2,10)\}}$; $\overset{C_2}{\{B(2,5), H(4,9), D(5,8), E(7,5), F(6,4), C(8,4)\}}$; $\overset{C_3}{\{G(1,2), B(2,5)\}}$

recalculează centrul:

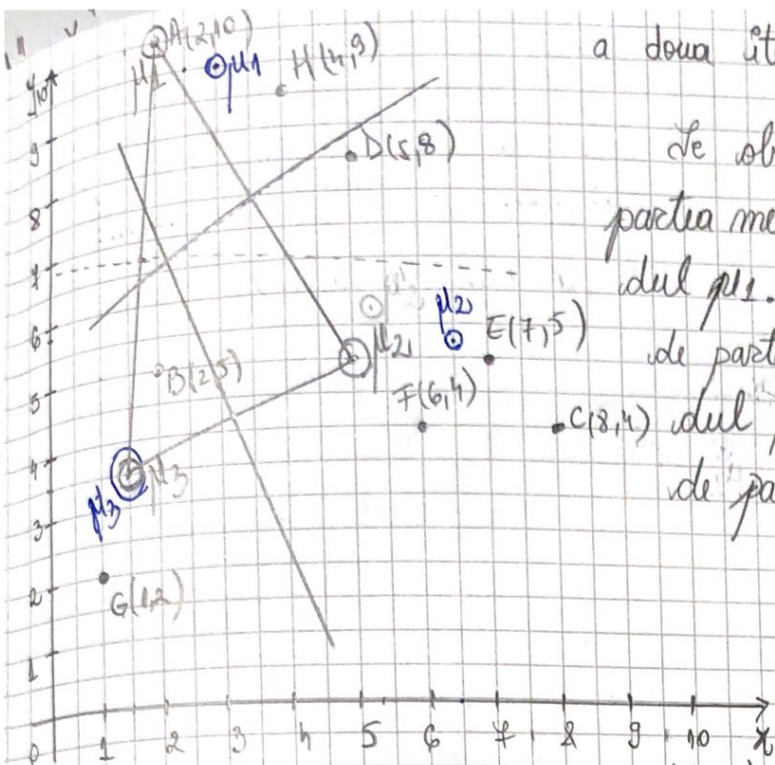
$$\mu_1 = A(2,10)$$

$$\mu_3 = \left(\frac{x_G + x_B}{2}, \frac{y_G + y_B}{2} \right) = \left(\frac{1+2}{2}, \frac{2+5}{2} \right) = \left(\frac{3}{2}, \frac{7}{2} \right) = (1,5; 3,5)$$

$$\mu_2: x_2 = \frac{x_H + x_D + x_E + x_F + x_C}{6} = \frac{4+5+7+6+8}{6} = \frac{30}{6} = 5$$

$$y_2 = \frac{y_H + y_D + y_E + y_F + y_C}{6} = \frac{9+8+5+4+4}{6} = \frac{30}{6} = 5$$

$$\Rightarrow \mu_2 = (10/3; 35/6) \quad \mu_2 = (5,5)$$



a doua iterație:

Se observă că punct H este de partea mediatoarei (μ_1, μ_2) ca și centrul μ_1 . Punctele D, E, F, C sunt de partea mediatoarei ca și centrul μ_2 . Iar punctele B și G sunt de partea centrului μ_3 .

→ clusterelor sunt: $C_1 = \{H(4,9)\}$; $C_2 = \{D(5,8); E(7,5); C(8,4); F(6,4)\}$
 $C_3 = \{B(2,5); G(1,2)\}$

recalculează centrul:

$$\mu_1 = x_1 = \frac{x_A + x_H}{2} = \frac{2+4}{2} = 3 \quad y_1 = \frac{y_A + y_H}{2} = \frac{10+9}{2} = \frac{19}{2} = 9,5$$

$$\mu_1 = (3; 9,5)$$

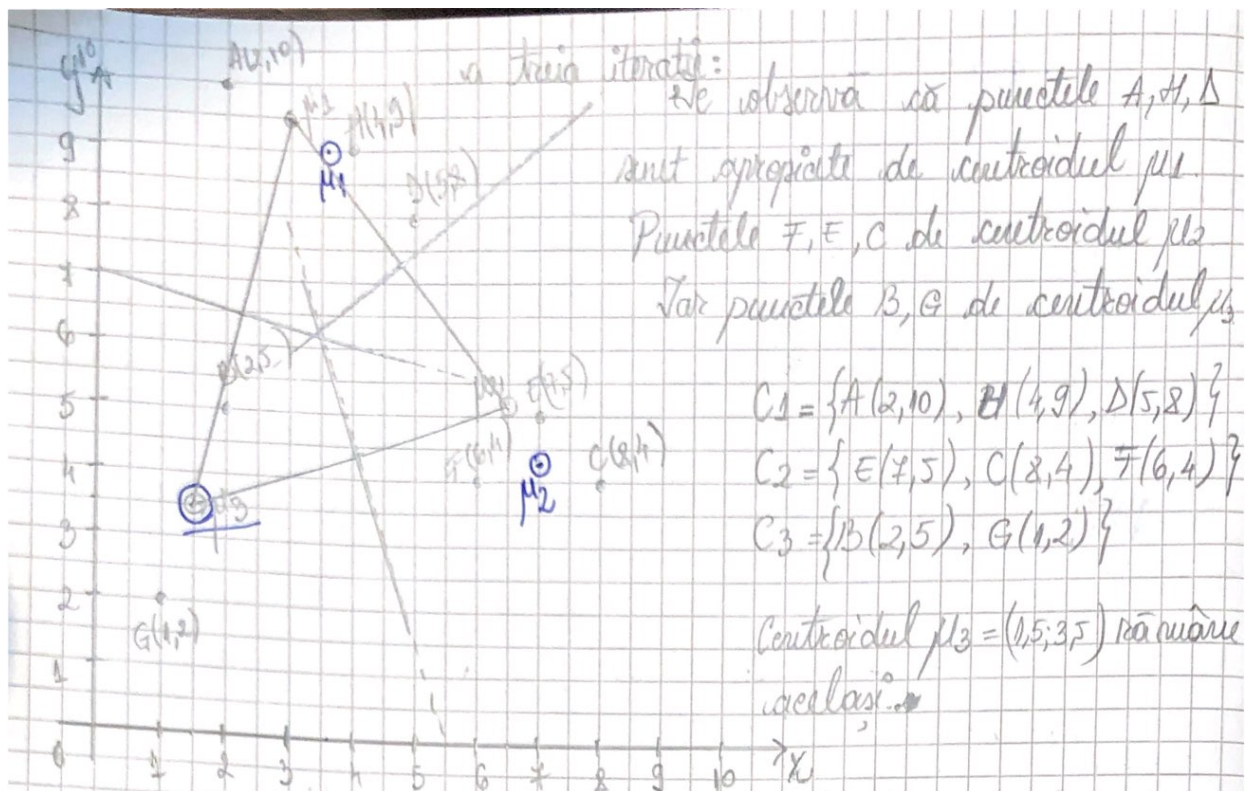
$$\mu_2 = x_2 = \frac{x_D + x_E + x_C + x_F}{4} = \frac{5+7+8+6}{4} = \frac{26}{4} = 6,5$$

$$y_2 = \frac{y_D + y_E + y_C + y_F}{4} = \frac{8+5+4+4}{4} = \frac{21}{4} = 5,25$$

$$\mu_2 = (6,5; 5,25)$$

$$\mu_3 = x_3 = \frac{x_B + x_G}{2} = \frac{2+1}{2} = \frac{3}{2} = 1,5; \quad y_3 = \frac{y_B + y_G}{2} = \frac{5+2}{2} = \frac{7}{2} = 3,5$$

$$\mu_3 = (1,5; 3,5)$$



$$\mu_1: x_1 = \frac{2+4+5}{3} = \frac{11}{3} = 3,67 \quad y_1 = \frac{10+9+8}{3} = \frac{27}{3} = 9$$

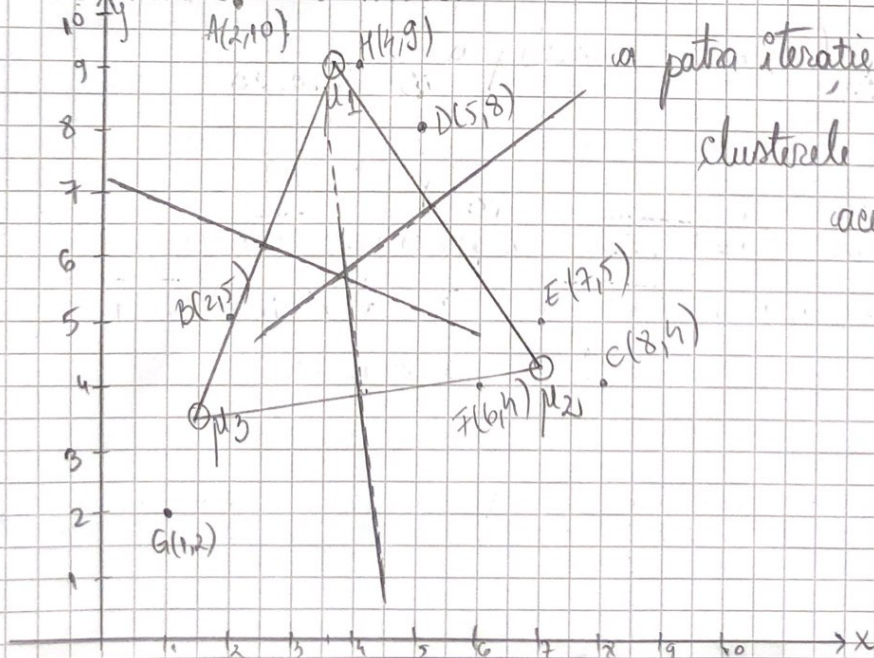
$$\mu_1 = (3,67; 9)$$

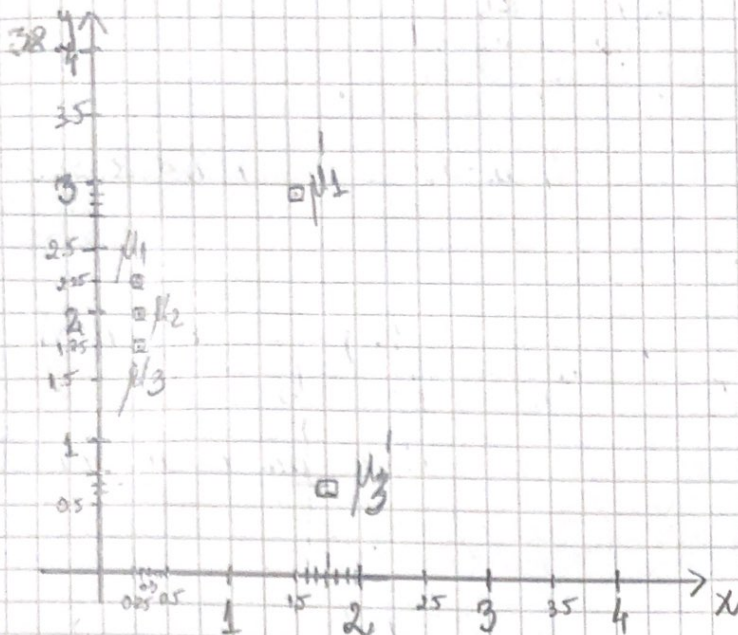
$$\mu_2: x_2 = \frac{7+8+6}{3} = \frac{21}{3} = 7 \quad y_2 = \frac{5+4+4}{3} = \frac{13}{3} = 4,3(3)$$

$$\mu_2 = (7; 4,3(3))$$

Se poate observa ca dupa a patra iteratie algoritmul

K-means a conver.





iteratia 1:

$$\mu_3: \bar{x}_3 = \frac{0.25 + 0.50 + 0.75 + 1 + 1.4 + 0.5 + 0.25 + 0.75 + 0.35 + 0.85 + 3.25 + 3.50 + 3 + 3.25 + 3.45 + 3.75 + 3.25}{17} = \frac{30.05}{17} = 1.768 \Rightarrow \mu_3' = (1.8; 0.7)$$

$$\bar{y}_3 = \frac{0.5 + 0.35 + 0.50 + 0.35 + 0.70 + 0.85 + 2 + 2.5 + 0.85 + 1 + 0.85 + 0.85 + 0.85}{17} = \frac{12.15}{17} = 0.715$$

$$\mu_1': \bar{x}_1 = \frac{20.44}{13} = 1.575 = 1.6 \Rightarrow \mu_1' = (1.6; 2.9)$$

$$\bar{y}_1 = \frac{37.45}{13} = 2.88 = 2.9 \quad \mu_2' = (0.30, 2.00)$$

iteratia 2:

$$d(\mu_2', A) = d(\mu_2', (0.25, 0.50)) = \sqrt{2.2525} < d(\mu_3', A) = \sqrt{2.44}$$

$$d(\mu_2', D) = d((0.30, 2), (0.75, 1.00)) = \sqrt{0.45^2 + 1^2} = \sqrt{1.20} < d(\mu_3', D) = \sqrt{3.94}$$

$$d(\mu_2', T) = d((0.30, 2), (1.25)) = \sqrt{0.74} \quad d(\mu_1', T) = d((1.6; 2.9); (1.25)) = \sqrt{0.52}$$

recalculăm centrul: $\mu_1'' = (1.7; 2.9)$

$$\bar{x}_1'' = \frac{0.85 + 1.0 + 1.0 + 1.10 + 1.20 + 1.25 + 1.50 + 1.20 + 1.37 + 3.0 + 3.25 + 3.1}{12} = \frac{19.82}{12} = 1.65 = 1.7$$

$$\bar{y}_2'' = \frac{3.35 + 3.0 + 2.5 + 3.20 + 2.40 + 2.50 + 2.50 + 2.75 + 2.75 + 3.25 + 3.0 + 3.50}{12} = \frac{34.7}{12} = 2.9$$

$$39. J(c^{(t)}, \mu^{(t)}) = \sum_{i=1}^n (x_i - \mu_{c^{(t)}(x_i)}^{(t)})^2$$

a) set de date: -9, -8, -7, -6, -5, 5, 6, 7, 8, 9, 9

$$\mu_1^{(0)} = -20$$

$$\mu_2^{(0)} = -10$$

$$\text{dum c\u0103 pt. } t=1 \Rightarrow J(c^{(t)}, \mu^{(t)}) \leq J(c^{(t-1)}, \mu^{(t-1)})$$

Inegalitatea de mai sus este dat\u0103 din d\u0103ru\u0107 inegalit\u0107i:

$$J(c^t, \mu^t) \stackrel{(1)}{\geq} J(c^t, \mu^{t+1}) \stackrel{(2)}{\geq} J(c^{t+1}, \mu^{t+1})$$

Inegalitatea 1: $J(c^t, \mu^t) \geq J(c^t, \mu^{t+1})$

$$J(c^t, \mu^t) = \sum_{i=1}^n (x_i - \mu_{c^t(x_i)}^t)^2$$

not. $x_{i_1}, x_{i_2}, \dots, x_{i_l}$ instan\u0219ele din comp. clusterului C_j^t ,
unde $l = |C_j^t|$

$$\Rightarrow J(C_j^t, \mu_j^t) = \sum_{p=1}^l (x_{i_p} - \mu_j^t)^2$$

$$\text{i\u0107 } J(c^t, \mu^t) = \sum_{j=1}^K J(C_j^t, \mu_j^t)$$

Se consider\u0103 C_j^t fixat, iar μ_j^t variabil.

! Este normal ca pentru o n\u0103u\u0107are iterativ\u0107 s\u0103
avem un J mai mic sau egal deoarece distan\u0219ele
de la centroid la puncte scad sau r\u0103m\u0107n la fel,
nu se \u0103n\u0103u\u0107esc / dep\u0103\u0107esc.

40. Dacă algoritmul K-means realizează K-partiția respectivă, lucru care este posibil, înseamnă că algoritmul K-means a convergat. Deoarece clusterurile rămân aceleași și presupunem că și centrozii rămân aceiași.

41. Este soluția optimă pentru alegerea centrozilor, deoarece astfel, vom afla punctul în care putem plasa centrul iar distanța până la valori să fie minimă. Așa cum se aleg centrozii și în cazul algoritmului K-means.

48. Single-linkage ne conduce cel mai probabil la formarea unor clusteruri asemănătoare cu cele obținute de K-means deoarece K-means se bazează pe distanța minimă dintre punct și centrul, iar single-linkage se bazează de asemenea pe distanța minimă dintre punctele (x, y) , cu $x \in X$ și $y \in Y$.



Clustriizarea optimă pentru setul de date este: #.

$$C_1 = \{1, 2\} \quad C_2 = \{5\} \quad C_3 = \{7\}$$

Iar valoarea pt. funcția J este

$$J = \sum_{i=1}^n \min_{j \in \{1, \dots, K\}} \|x_i - \mu_j\|^2$$

$$J = 0,5 + 0,5 + 0 + 0 = 1.$$

b) $C_1 = \{1, 2\} \quad C_2 = \{5, 7\} \quad C_3 = \{\emptyset\}$

Valoarea \emptyset din clusterul C_3 înseamnă că centrul aflat la poziția pp. 10, care are fiecare element apropiat lui.

K -m nu poate îmbunătăți aceste clustere deoarece μ_3 rămâne în aceeași poziție, iar C_1 și C_2 vor converge tot spre centrozii lor.

c) Se poate observa că această proprietate este satisfăcută pe parcursul alg. K -means, deoarece centrozii vor fi plasați în centrul unor valori apropiate, cecătoare, nu există cazuri în care

$$C_1 = \{1, 2, 7\}$$

$$C_2 = \{5, 8\}$$

Deoarece centrul pt. C_2 va trage aproape de el și valoarea 7; iar pentru val $\{1, 2\}$ va rămâne centrul pt. C_1 .

d) Ideea mea ar fi să așezăm din cei K centrozii, unul în limita inf, unul în limita

suprafață, după care luăm lungimea spațiului
de id valori și o împărțim cu $(K-2)$ pentru a
împărți centrul pe toate suprafețele spațiului.
După care fiecare valoare va fi asociată unui centroid,
adică celui mai apropiat, iar centroidul vor fi
recalculați pentru a fi în mijlocul valorilor.