

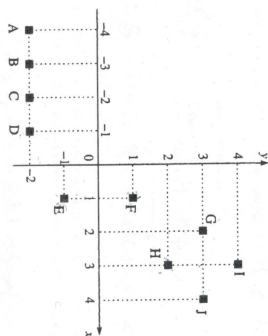
2.

(Clusterizare ierarhică aglomerativă: exemplificare pe date din \mathbb{R}^2 , folosind tipurile de similitudine "single-linkage" și "average-linkage")

Considerăm setul de date din figura alăturată.

a. Realizați clusterizarea ierarhică a datelor în maniera bottom-up, folosind similitudine de tip "single-linkage" și distanța euclidiană între puncte.

b. Realizați clusterizarea ierarhică a datelor, folosind de această dată similitudine de tip "average-linkage".



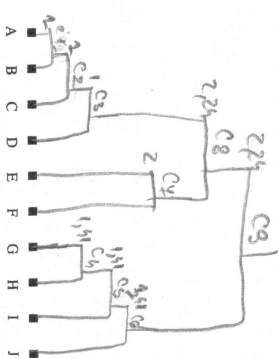
Pentru ușurința calculului, în tabelul alăturat sunt listate distanțele dintre perechile de instanțe date mai sus.

	A	B	C	D	E	F	G	H	I	J
A	0.00	1.00	2.00	3.00	5.10	5.83	7.81	8.06	9.22	9.43
B	1.00	0.00	1.00	2.00	4.12	5.00	7.07	7.21	8.49	8.60
C	2.00	1.00	0.00	1.00	3.16	4.24	6.40	6.40	7.81	7.81
D	3.00	2.00	1.00	0.00	2.24	3.61	5.83	5.66	7.21	7.07
E	5.10	4.12	3.16	2.24	0.00	2.00	4.12	3.61	5.39	5.00
F	5.83	5.00	4.24	3.61	2.00	0.00	2.24	2.24	1.41	2.00
G	7.81	7.07	6.40	5.83	4.12	2.24	0.00	1.41	2.00	1.41
H	8.06	7.21	6.40	5.66	3.61	2.24	1.41	0.00	2.00	1.41
I	9.22	8.49	7.81	7.21	5.39	3.61	2.00	2.00	0.00	1.41
J	9.43	8.60	7.81	7.07	5.00	3.61	2.00	1.41	1.41	0.00

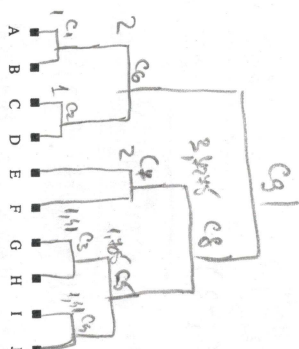
Observație (1): Dacă la o iterație a algoritmului de clusterizare „distanțele” (adică similitudinile) dintre două perechi de clusterizări au aceeași valoare, prioritatea la alcătuirea noului cluster este dictată de ordinea alfabetică.

Observație (2): La construcția dendrogramei, înălțimea [nodului rădăcinii] pentru fiecare cluster C_i va fi dată de distanța (adică măsura de similitudine) dintre clusterelor C_j și C_k din care a fost construit clusterul C_i .

Veți construi cele două dendrograme pornind de la reprezentările de mai jos.



Dendrograma single-linkage



Dendrograma average-linkage

$$\begin{aligned}
 d(A, B, C, D) &= \frac{2 + 3 + 1 + 2}{4} = 2 \\
 d(F, G, H) &= \frac{2 + 2 + 2 + 2}{4} = 2 \\
 d(G, H, I, J) &= \frac{1 + 1 + 2 + 1 + 1 + 2}{4} = \frac{8}{4} = 2 \\
 d(A, B, C, D, E, F) &= \frac{6 + 10 + 4 + 1 + 2 + 3 + 10 + 2 + 2 + 4 + 3 + 6 + 1}{8} = 4.125 \\
 d(E, F, G, H, I, J) &= \frac{4 + 1 + 2 + 3 + 1 + 6 + 1 + 5 + 3 + 3 + 5 + 6 + 1}{8} = 4.125
 \end{aligned}$$

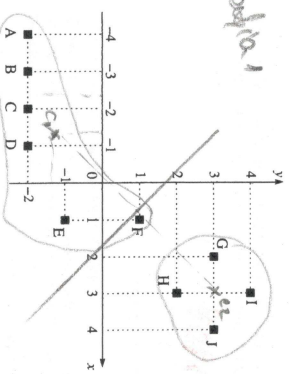
$$d(A, B, C, D, E, F, G, H, I, J) = \dots = 6.3908$$

c. La curs am afirmat că una dintre metodele de inițializare pentru algoritmul K-means este următoarea: clusterelor inițiale sunt cele K cluster de la vârful dendrogramei (engl., top clusters) obținute cu ajutorul unei metode de clusterizare ierarhică.

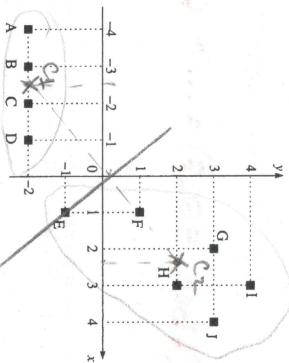
Executați algoritmul K-means pe datele din enunț, pentru fiecare dintre cele două variante de inițializare, folosind rezultatele de la punctele precedente. (Vezi lucră cu $K = 2$.) Folosiți pentru aceasta reprezentările de mai jos. Indicați la fiecare iterație coordonatele centrozilor și componența fiecărui cluster.

Setul de date din acest exercițiu este oarecum particular. Judecând pe un caz [mai] general, care dintre cele două tipuri de funcții de similaritate folosite la punctele a și b vi se pare mai adecvată pentru pasul de inițializare a algoritmului K-means?

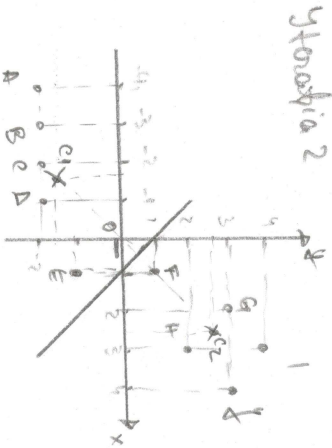
d. Considerând cele două cluster de "top" obținute de către algoritmul de clusterizare ierarhică care folosește similaritate de tip "single-linkage" ca fiind niște clase, trasați pe desenul din stânga de mai jos, separatorul decizional produs de algoritmul 1-NN.



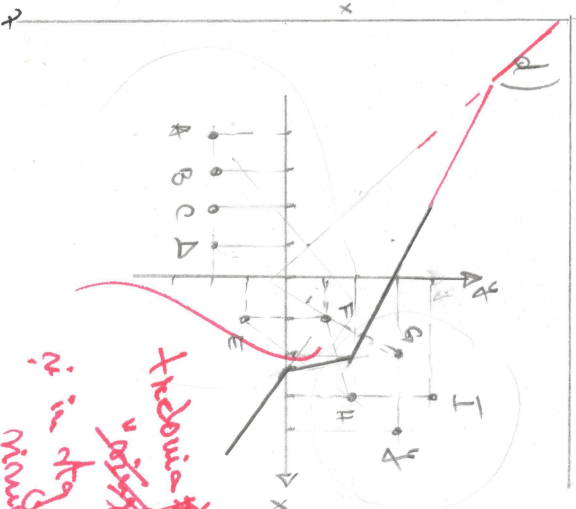
K-means cu clusterelor inițiale cf. clusterizării ierarhice single-linkage



K-means cu clusterelor inițiale cf. clusterizării ierarhice average-linkage



Ytrodofia 2



Ytrodofia 2
ni in dta mas!
naua glauz!

c) Cozud single linkage!

$$x_{c1} = \frac{-4-3-2-1+1+1}{6} = -2.16$$

$$x_{c2} = \frac{2+2.3+4}{4} = 3$$

$$y_{c1} = \frac{-2.4-1+1}{6} = -1.33$$

$$y_{c2} = \frac{2+2.3+4}{4} = 3$$

$$c_1(-1,3); -1,3)$$

$$c_2(3,3)$$

$$d(F, c_1) = \sqrt{(-1,3+1)^2 + (-1,3+1)^2} = 2.25$$

$$d(F, c_2) = \sqrt{(3-1)^2 + (3-1)^2} = \sqrt{8} = 2.82 \quad \rightarrow F \in \text{Cluster 2}$$

Ytrodofia 2:

$$x_{c1} = \frac{-4-3-2-1+1}{5} = -1.8$$

$$x_{c2} = \frac{1+2+2.3+4}{5} = 2.6$$

$$y_{c1} = \frac{-2.4-1}{5} = -1.8$$

$$y_{c2} = \frac{1+2+2.3+4}{5} = 2.6$$

$$c_1(-1,8); -1,8)$$

$$c_2(2,6); 2,6)$$

Clusterela nu re mai selamuta, centrul nu nău
la fel, algoritmul s-a stabilizat!

$$c_1: 4, 1, 3, 1, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13$$

Cozud Average-linkage!

$$x_{c1} = \frac{-4-3-2-1}{4} = -2.5$$

$$x_{c2} = \frac{2.1+2+2.3+4}{6} = 2.13$$

$$y_{c1} = -2$$

$$y_{c2} = \frac{-1+1+2+2.3+4}{6} = 2$$

$$c_1(-2,5); -2)$$

$$c_2(2,13); 2)$$

$$d(F, c_1) = \sqrt{(1+2.5)^2 + (1+1+2)^2} = \sqrt{17.25+1} = 4.34 \rightarrow F \in \text{Cluster 2}$$

$$d(F, c_2) = \sqrt{(1-2.13)^2 + (1-2)^2} = \sqrt{1.63+1} = 1.63$$

\Rightarrow Clusterela si centrul nu nău selamuta, algoritmul s-a stabilizat

$$c_1: 4, 1, 3, 1, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13$$

8