

Învățare automată

— Licență, anul III, 2018-2019, re-examinare, parțial II —

Nume student:

Grupa:

1.

(Clusterizare ierarhică aglomerativă:
aplicare pe date din \mathbb{R}^2 , folosind măsurile de similaritate
single-linkage, complete-linkage și metrica lui Ward)

Pe setul de date din \mathbb{R}^2

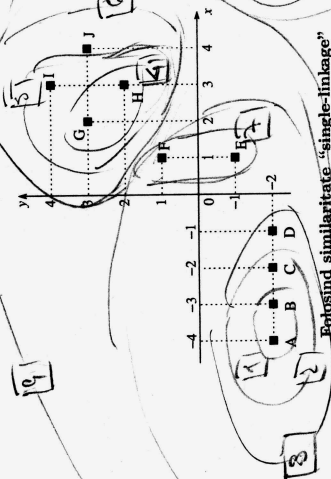
$$A: (-4, -2), B: (-3, -2), C: (-2, -2), D: (-1, -2), E: (+1, -1) \\ F: (+1, +1), G: (+2, +3), H: (+3, +2), I: (+3, +4), J: (+4, +3)$$

veți aplica algoritmul de clusterizare ierarhică aglomerativă conform specificațiilor de la fiecare din punctele următoare.

Prezervare: Dacă la o iterație a algoritmului de clusterizare distanțele (adică similaritățile) dintre două perechi de clustere au aceeași valoare, prioritatea la alcătuirea noului cluster este dictată de ordinea alfabetică.

a. Folosind măsurile de similaritate single-linkage și complete-linkage, reprezentați pe desenele de mai jos rezultatele aplicării algoritmului de clusterizare ierarhică aglomerativă, sub forma unor *dendrograme* (adică, ierarhii) *aplatizate* (engl., flat hierarchies).

Pentru fiecare cluster non-singleton veți folosi câte o curbă închisă de formă elipsoidală pentru a încadra punctele din clusterul respectiv. Fiecărei elipse îi veți asocia câte un index, scris sub forma $[1], [2], [3] \dots$ (chiar pe conturul elipsei respective) pentru a indica ordinea în care sunt formate clusterele.



b. La curs am prezentat încă o funcție de similaritate, numită *metrica lui Ward*. Potrivit acestei metrice, distanța dintre două clustere disjuncte X și Y se definește astfel:

$$\Delta(X, Y) = \sum_{x_i \in X, y_j \in Y} \|x_i - \mu_{X \cup Y}\|^2 - \sum_{x_i \in X} \|x_i - \mu_X\|^2 - \sum_{y_j \in Y} \|y_j - \mu_Y\|^2 \quad (1)$$

unde, spre exemplu, μ_X este centrulidul [sau „centrul de greutate” al clusterului X , iar x_i este o instanță generică dintr-un cluster [oarecare, fixat]. Prin definiție, aici vom considera $\mu_X = \frac{1}{n_X} \sum_{x_i \in X} x_i$, unde n_X este numărul de elemente din X . (Similar sunt definiți centrulidii μ_Y și $\mu_{X \cup Y}$.)

Se poate arăta că

$$\Delta(X, Y) = \frac{n_X n_Y}{n_X + n_Y} \|\mu_X - \mu_Y\|^2. \quad (2)$$

Observații:

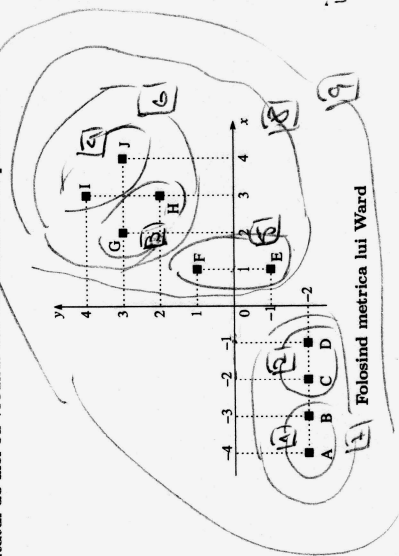
1. Pentru perechi de clustere (X, Y) și (X', Y') astfel încât $n_Y = n_{Y'}$ și $n_{Y'} = n_{Y'}$, formula (2) arată că la clusterizare ierarhică este „favorizată” acea pereche pentru care centrulidii $(\mu_X$ și $\mu_{Y'}$, respectiv $\mu_{X'}$ și $\mu_{Y'})$ sunt mai apropiați.

2. Invers, dacă $\|\mu_X - \mu_{Y'}\| = \|\mu_{X'} - \mu_{Y'}\|$, atunci este favorizată perechea pentru care ponderea¹ (adică $\frac{n_X n_Y}{n_X + n_Y}$, respectiv $\frac{n_{X'} n_{Y'}}{n_{X'} + n_{Y'}}$) este mai mică.

Aceste două observații vă vor ajuta să simplificați / reduceți foarte mult calculele pe care ar trebui să le faceți la punctul b).

Aplicați algoritmul de clusterizare ierarhică aglomerativă pe același set de date ca mai sus, însă folosind de această dată metrica lui Ward. Ca și la punctul a, veți folosi elipse (și indici) pentru a reprezenta noua *dendrogramă aplatizată*.

Coincide rezultatul de aici cu vreunul din rezultatele de la punctul a?



Observație:

La punctul b, la fiecare iterație a algoritmului veți justifica riguros alegerea făcută, scriind [doar] calculele care sunt determinante!

Răspuns: (pentru punctul b)

Iterația 1:

$$\dots \dots \dots \Delta(\{X\}, \{Y\}) = \min_{X' \in \{A, \dots, J\}} \|X - Y'\|^2 \\ \dots \dots \dots X \neq Y$$

Iterația 2:

$$\dots \dots \dots \Delta(\{A, B\}, \{C\}) = \frac{3}{5} \cdot 15^2 = \frac{3}{5} \cdot \left(\frac{3}{2}\right)^2 = \frac{3}{2} = \Delta(\{C\}, \{D\}) \\ \dots \dots \dots \Delta(\{A, B\}, \{X\}) = \Delta(\{A, B\}, \{C\}) \quad \forall X \in \{D, \dots, J\}$$

¹ Această pondere este jumătate din *media armonică* a cardinalilor n_X și n_Y .