

# Învățare automată

— Licență, anul III, 2021-2022, examenul parțial II —

Nume student:

Grupa:

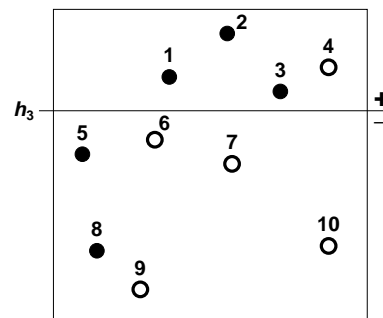
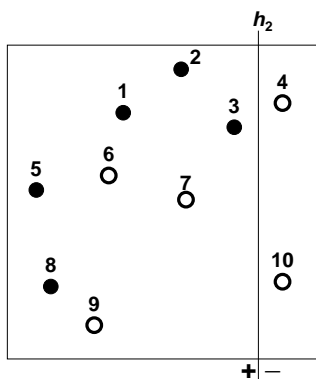
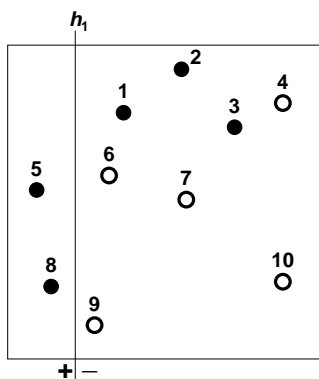
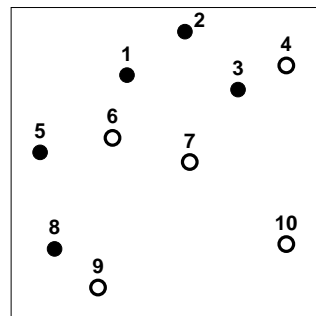
1. [3.5p]

(Algoritmul AdaBoost: aplicare pe un set de date din  $\mathbb{R}^2$ )

Se consideră că aplicăm algoritmul AdaBoost pe dataset-ul din figura alăturată. (Pentru ușurința exprimării la calcule, am notat pe figură indicii instanțelor de antrenament, în imediata apropiere a acestora.)

Folosim convenția noastră obișnuită de notare: simbolul  $\bullet$  desemnează instanțe pozitive, iar simbolul  $\circ$  instanțe negative.

La primele trei iterații ale algoritmului au fost selectați compașii de decizie  $h_1$ ,  $h_2$  și  $h_3$  (în această ordine), așa cum se indică în figurile de mai jos.



*Obiectivul* acestei probleme este să determinăm dacă la sfârșitul celor trei iterații algoritmul AdaBoost reușește să clasifice perfect toate instanțele de antrenament.

a. Calculați

- distribuțiile probabile corespunzătoare celor 3 iterații, adică  $D_1(x_i)$ ,  $D_2(x_i)$  și  $D_3(x_i)$ , pentru  $i = 1, \dots, 10$ ;
- pentru fiecare dintre cele 3 iterații ( $t = 1, 2, 3$ ): eroarea ponderată la antrenare ( $\varepsilon_t$ ) produsă de compasul de decizie  $h_t$ , precum și ponderea ( $\alpha_t$ ) asociată ipotezei / compasului de decizie  $h_t$ .

Veți completa tabelele următoare și veți indica succint(!) modul în care ați procedat pentru a ajunge la rezultatele respective.

$i$	1	2	3	4	5	6	7	8	9	10
$D_1(x_i)$										
$D_2(x_i)$	1/6			1/14						
$D_3(x_i)$	7/66			1/22		1/6				

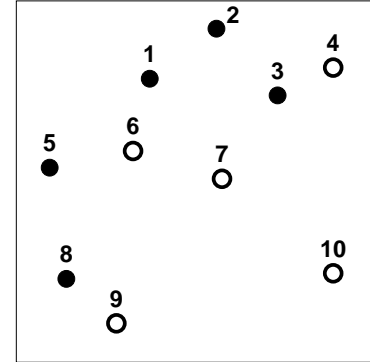
$t$	1	2	3
$\varepsilon_t$			
$\alpha_t$			

**Atenție!** Pentru a vă ușura munca, am completat noi câteva dintre elementele primului tabel. Vă puteți baza pe valorile indicate de noi, ca să vă simplificați raționamentele / calculele!

b. Folosind ipoteza combinată obținută de algoritmul AdaBoost la finalul celei de-a treia iterații, stabiliți

i. eroarea la antrenare produsă (pentru calcularea ei, puteți folosi tabelul de mai jos),

ii. zonele de decizie corespunzătoare acestui clasificator (veți justifica modul în care ați procedat!). Veți indica aceste zone de decizie, precum și granițele de decizie, pe desenul de alăturat.



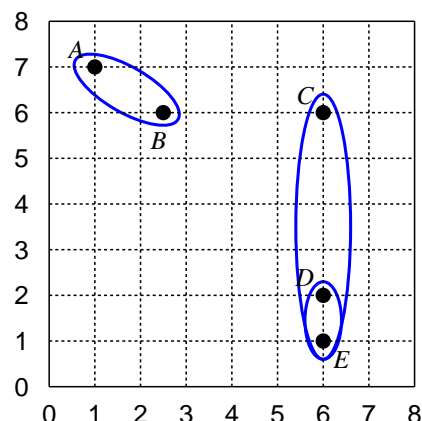
$t$	$\alpha_t$	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$x_7$	$x_8$	$x_9$	$x_{10}$
1	$\alpha_1 = \dots$										
2	$\alpha_2 = \dots$										
3	$\alpha_3 = \dots$										
	$H_3(x_i)$										

**Atenție!** Pentru a vă ușura munca, vă furnizăm noi următoarele valori numerice:  $\ln \sqrt{7/3} \approx 0.423$ ,  $\ln \sqrt{11/3} \approx 0.649$ ,  $\ln \sqrt{15/7} \approx 0.923$ .

2. [1.75p] (Clusterizare ierarhică aglomerativă: un exemplu simplu de aplicare, cu single-, complete- și average-linkage, pe date din  $\mathbb{R}^2$ )

Considerăm punctele  $A(1, 7)$ ,  $B(2.5, 6)$ ,  $C(6, 6)$ ,  $D(6, 2)$  și  $E(6, 1)$  din planul euclidian. Pe acest dataset veți aplica algoritmul de clusterizare ierarhică aglomerativă (i.e., bottom-up), folosind pe rând (separat) funcțiile de similaritate single-linkage, complete-linkage și average-linkage.

Care dintre aceste funcții de similaritate va conduce după executarea a trei iterații consecutive la ierarhia aplatizată prezentată în figura alăturată?



*Indicație:*

Veți folosi grid-urile de mai jos.

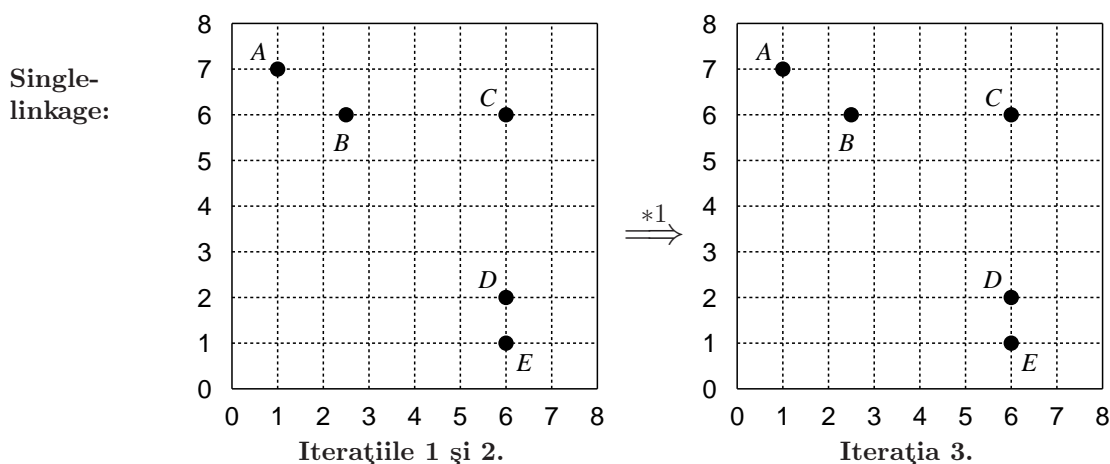
Pentru fiecare dintre cele 3 tipuri de măsuri de similaritate, *veți justifica riguros* trecerea de la iterația 2 la iterația 3 determinând (*numeric!*) minimul *distanțelor* dintre perechile de *cluster* care au fost formate la iterația 2.

De *exemplu*, dacă în figura de mai sus notăm  $C_1 = \{D, E\}$ ,  $C_2 = \{A, B\}$ , atunci *justificarea* acestui rezultat ar putea fi scrisă sub forma

$$\underbrace{d(C_1, C)}_{= \dots} < \underbrace{d(C_2, C)}_{= \dots} \text{ și } d(C_1, C) < \underbrace{d(C_1, C_2)}_{= \dots}, \quad (*)$$

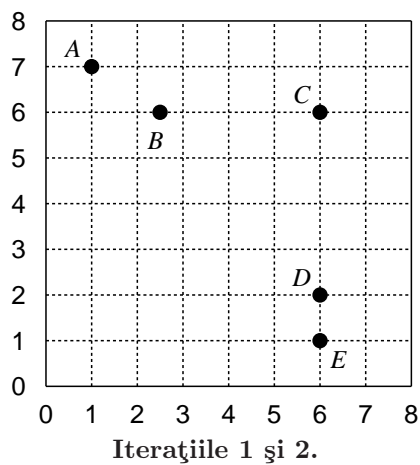
unde  $d$  va fi una dintre „măsurile de distanță“  $d_{SL}$ ,  $d_{CL}$  sau  $d_{AL}$ , corespunzătoare măsurilor de similaritate single-linkage, complete-linkage și respectiv average-linkage.

Răspuns:

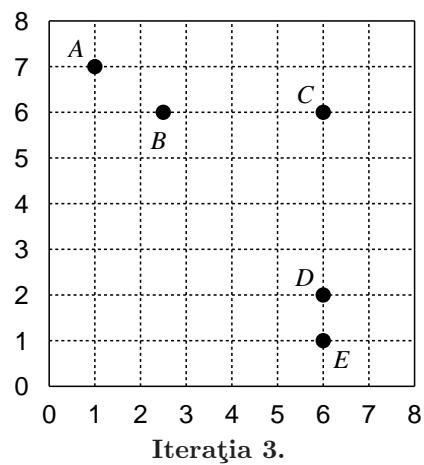


\*1:

Complete-linkage:

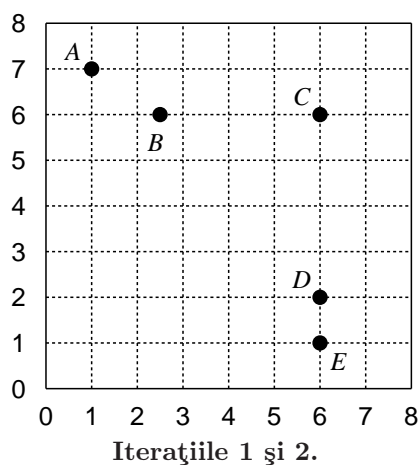


$\Rightarrow$  \*2

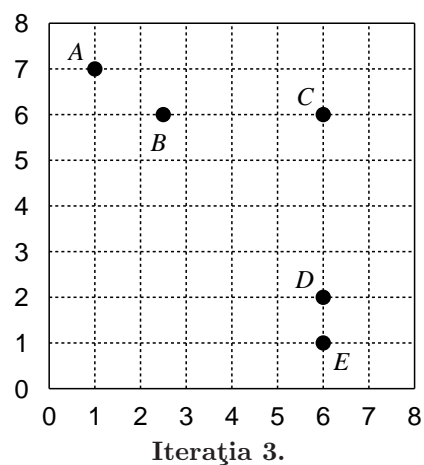


\*2:

Average-linkage:



$\Rightarrow$  \*3



\*3:

3. [1.75p] (Algoritmul  $K$ -means: chestiuni legate de valoarea criteriului  $J$ )

Vă readucem aminte că la clusterizare cu algoritmul  $K$ -means obiectivul este să găsim pozițiile celor  $K$  centroizi ai clusterelor, notați cu  $c_j \in \mathbb{R}^d, j \in \{1, \dots, K\}$ , astfel încât suma pătratelor distanțelor dintre fiecare instanță  $x_i$  și cel mai apropiat centroid să fie minimizată. Așadar, funcția obiectiv este

$$\sum_{i=1}^n \min_{j \in \{1, \dots, K\}} \|x_i - c_j\|^2, \quad (1)$$

unde  $n$  este numărul de instanțe de clusterizat. Altfel spus, încercăm să găsim  $c_1, \dots, c_k$  astfel încât să fie minimizată expresia (1). Pentru aceasta, efectuăm mai multe iterații în cadrul cărora asignăm fiecare instanță  $x_i$  la cel mai apropiat centroid și apoi actualizăm poziția fiecărui centroid  $c_j$  la media instanțelor asignate la clusterul  $j$ .

Însă prietenul tău Ionuț, în loc să mențină numărul de clustere  $K$  fixat, încearcă să minimizeze valoarea expresiei (1) variindu-l pe  $K$ . Tu ești de părere că această idee nu este bună.

În mod concret, tu îl convingi pe Ionuț dându-i două valori:  $\alpha$ , care reprezintă minimul valorilor posibile pentru expresia (1), și  $\beta$ , care este valoarea lui  $K$  atunci când expresia (1) își atinge valoarea minimă.

a. Cât este valoarea expresiei  $\alpha + \beta$  în cazul în care  $n = 100$ ?

b. Presupunem că datele de clusterizat este format din 4 instanțe situate pe axa reală, și anume  $x_1 = 1, x_2 = 2, x_3 = 5$  și  $x_4 = 7$ , iar  $K$  are valoarea 3. Cât este în acest caz valoarea optimă a funcției obiectiv (1)?

4. [1.25p] (Distribuția Bernoulli: estimare în sensul verosimilității maxime)

La curs am spus că la estimarea de verosimilitate maximă (engl., Maximum Likelihood Estimation, MLE), definim

$$\hat{\theta}_{MLE} = \arg \max_{\theta} P(\mathcal{D}|\theta) = \arg \max_{\theta} \ln P(\mathcal{D}|\theta),$$

unde  $\mathcal{D}$  este setul de date pe care se estimează parametrul  $\theta$  al distribuției de probabilitate  $P$ .

Imaginează-ți că ești un *data scientist* care lucrează pentru o firmă de publicitate (engl., advertisement). Această firmă a făcut recent o campanie de publicitate [care constă în publicarea unui anunț publicitar pe ecranul calculatorului / telefonului unor anumite persoane], iar ție ți se cere să evaluezi succesul acestei campanii de publicitate.

În cadrul acestei campanii, au fost targetate / vizate  $N$  persoane. Variabila aleatoare  $Y_i$  primește valoarea  $y_i = 1$  dacă persoana  $i$  a făcut “click” pe anunțul publicitar și  $y_i = 0$  în cazul contrar. În total,  $\sum_{i=1}^N y_i = k$  persoane au decis să facă “click” pe anunțul publicitar. Vom presupune că probabilitatea ca persoana  $i$  să facă “click” pe anunțul publicitar respectiv este  $\theta$ , iar probabilitatea ca persoana  $i$  să nu facă “click” pe anunțul publicitar este  $1 - \theta$ .

Desigur,  $P(\mathcal{D}|\theta) = P(Y_1, \dots, Y_N|\theta)$ .

a. Dedu expresia analitică pentru calcularea lui  $\hat{\theta}_{MLE}$ .

b. Presupunem că  $N = 100$  și  $k = 10$ . Calculează  $\hat{\theta}_{MLE}$ .

5. [2.85p]

(O problemă à la C. Do și S. Batzoglou:  
rezolvarea unei mixturi de doi vectori de distribuții Bernoulli i.i.d.  
identificarea parametrilor și a variabilelor neobservabile;  
scrierea funcției de log-verosimilitate a datelor complete)

Să zicem că avem două monede,  $A$  și  $B$ . Sarcina ta este aceea de a afla (a „estima“)  $\theta_A$  și  $\theta_B$ , probabilitățile de apariție a feței *stemă* (engl., head) pentru fiecare dintre cele două monede. Însă eu sunt cam răutacios și nu-ți dau voie să arunci tu monedele. În schimb, decid să procedez astfel: voi arunca eu însumi monedele și după aceea îți voi comunica ție rezultatele aruncărilor. În mod concret, îți voi spune ceva de forma următoare: am ales una dintre cele două monede (nu-ți spun care anume),<sup>1</sup> am aruncat-o de 10 ori și am obținut în total de 7 ori stema și de 3 ori banul. Apoi am ales din nou una dintre cele două monede (poate aceeași cu cea dinainte, poate nu), am aruncat-o de 10 ori și am obținut în total de 5 ori stema și de 5 ori banul. În total, îți comunic de  $N$  ori câte o astfel de informație. (Așadar, la final vei dispune de rezultatele a  $10N$  aruncări ale monedelor.)

Facem *presupunerea* că la fiecare dintre aceste  $N$  serii cele 10 aruncări sunt independente unele de altele și că am reținut ordinea / succesiunea rezultatelor obținute în urma acestor aruncări și ți-o comunic.

a. Formulează aceasta ca pe o problemă de tip EM (Expectation-Maximization).

i. Care sunt datele observabile?

ii. Care sunt variabilele neobservabile / ascunse / latente (engl., hidden variables)?

iii. Care sunt parametrii modelului (adică parametrii distribuțiilor probabiliste folosite)?

b. Calculează [adică, stabilește pas cu pas care este] expresia funcției de *log-verosimilitate a datelor complete* la iterația  $t$ , pentru această problemă.

(Atenție! Nu ți se cere să rezolvi efectiv problema EM.)

---

<sup>1</sup>Vom considera că moneda  $A$  este aleasă întotdeauna cu probabilitatea  $\pi$ , iar moneda  $B$  cu probabilitatea  $1 - \pi$ .