

Învățare automată

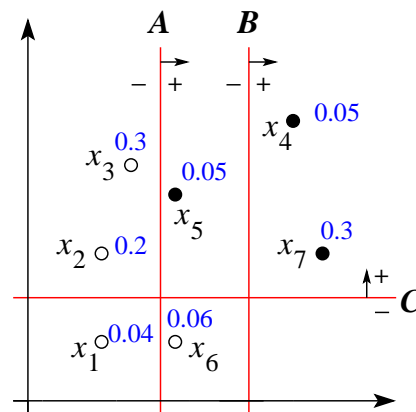
— Licență, anul III, 2020-2021, examenul parțial II —

Nume student:

Grupa:

1. (Algoritmul AdaBoost: exemplu de aplicare pe date din \mathbb{R}^2 ; întrebări calitative)

Fie setul de date de antrenament din figura alăturată. Am marcat cu semnul \circ exemplele / instanțele negative ($y_i = -1$) și cu semnul \bullet instanțele pozitive ($y_i = +1$). Figura conține de asemenea ponderile normalizate (adică, probabilitățile) asociate exemplilor de antrenament, așa cum au rezultat în urma executării unui anumit număr de iterații ale algoritmului AdaBoost. În figură sunt trasați și trei compași de decizie, $h(x; \theta_A)$, $h(x; \theta_B)$, and $h(x; \theta_C)$ sau, pe scurt, A , B și C .



- a. Care dintre acești trei compași de decizie considerați că a fost folosit la *precedenta iterație* a algoritmului AdaBoost, în așa fel încât să rezulte ponderile [asociate exemplilor] prezentate în figură? Veți răspunde indicând A , B sau C și veți justifica în mod riguros alegerea pe care ați făcut-o.
- b. Pe care dintre acești trei compași de decizie considerați că-l va selecta algoritmul AdaBoost la *iterația următoare*? Veți răspunde indicând A , B sau C și veți justifica riguros, prin calcule, alegerea pe care ați făcut-o.
- c. În figura dată, încercuiți instanțele de antrenament (este posibil să nu fie niciuna!) pe care ansamblul (i.e., combinația liniară de compași de decizie)

$$H_2(x) = \text{sign}(\alpha_A h(x; \theta_A) + \alpha_C h(x; \theta_C)) \text{ cu } \alpha_A = 0.3 \text{ și } \alpha_C = 0.5,$$

nu le poate clasifica în mod corect.

2.

(Concepte din \mathbb{R} *reprezentabile* cu ajutorul combinațiilor liniare de compași de decizie)

Presupunem că avem o problemă de clasificare a unor instanțe pe axa reală: fiecare instanță x_i este un număr real, iar etichetele pe care urmează să le prezicem sunt binare, $y_i \in \{-1, +1\}$.

Pentru această problemă de clasificare, veți folosi *ansambluri*, adică niște combinații liniare de separatori / ipoteze „slabe”. (Atenție! NU trebuie să folosiți algoritmul AdaBoost!) Vă readucem aminte că un astfel de *clasificator* are forma următoare:

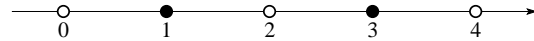
$$\hat{y} = \text{sign}\left(\sum_{t=1}^T \alpha_t h_t(x)\right), \quad (1)$$

unde \hat{y} este eticheta prezisă, $\text{sign}(x)$ este $+1$ dacă $x > 0$ și respectiv -1 în cazul contrar, α_t este o *pondere* (număr real strict pozitiv), iar $h_t(x)$ este predicția făcută de către ipoteza „slabă” h_t . Fiecare h_t ia una dintre următoarele forme:

$$h_t(x; s, +) = \begin{cases} -1 & \text{dacă } x \leq s \\ +1 & \text{dacă } x > s \end{cases} \quad h_t(x; s, -) = \begin{cases} +1 & \text{dacă } x \leq s \\ -1 & \text{dacă } x > s, \end{cases}$$

pentru un anumit *prag de separare* (engl., split threshold) $s \in \mathbb{R}$.¹

Considerăm setul de date din figura alăturată, format din 5 instanțe situate pe axa reală.



a. Arătați că acest set de date (LC: sau, acest *concept*) este *reprezentabil* cu ajutorul unei combinații liniare formată din 4 ipoteze „slabe”. Așadar, vă cerem să identificați în mod explicit 4 ipoteze „slabe” h_1, \dots, h_4 , precum și ponderile lor $\alpha_1, \dots, \alpha_4$, astfel încât ipoteza combinată $\text{sign}\left(\sum_{t=1}^T \alpha_t h_t(x)\right)$ să fie *consistentă* cu datele de mai sus.

Puneți în evidență cele 4 ipoteze „slabe” h_1, \dots, h_4 pe desenul de mai sus, folosind praguri de separare s_1, \dots, s_4 precum și linii verticale, cărora le veți asocia o etichetare adecvată. Concret, veți desemna cu semnele $+$ și $-$ la stânga și la dreapta fiecărei linii verticale — asociată unui anumit prag s_j — zonele de decizie determinate de această ipoteză „slabă”.

Observație: Ținând cont de modul în care a fost definită funcția *sign* mai sus, *regula de predicție* (1) va trata cazurile de „paritate” etichetând cu -1 instanțele pentru care suma ponderată (engl., weighted sum) a predicțiilor făcute de ipotezele „slabe” este 0. Acest mod de tratare a cazurilor de „paritate” este [foarte] util pentru rezolvarea punctului a.

Indicație: Pentru a justifica ușor consistența *ansamblului* ales de dumneavoastră cu datele din figura de mai sus, vă cerem să completați un tabel similar cu cel pe care-l folosim la algoritmul AdaBoost atunci când facem calculul erorii la antrenare produse de ipoteza combinată $\text{sign}\left(\sum_{t=1}^T \alpha_t h_t(x)\right)$.

α_t	h_t	0	1	2	3	4
$\alpha_1 =$	$h_1(x_i)$					
$\alpha_2 =$	$h_2(x_i)$					
$\alpha_3 =$	$h_3(x_i)$					
$\alpha_4 =$	$h_4(x_i)$					
$\sum_{t=1}^4 \alpha_t h_t(x_i)$						

¹Remarcați faptul că definiția dată aici pentru noțiunea de *ipoteză „slabă”* corespunde bine-cunoscutului *compas de decizie* (engl., decision stump), pentru cazul particular când datele sunt din \mathbb{R} .

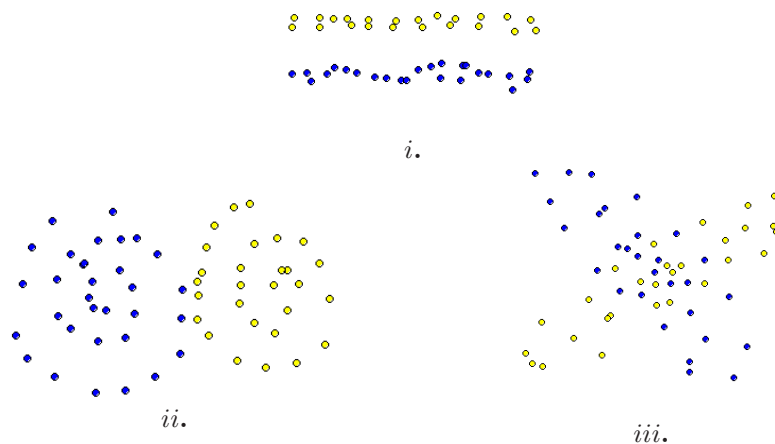
b. Demonstrați că setul de date de mai sus NU este reprezentabil cu mai puțin de 4 ipoteze „slabe“.

c. *Generalizați* rezultatul obținut la punctul a, referindu-vă la posibilitatea (sau, dimpotrivă, imposibilitatea) de a *reprezenta* — folosind combinații liniare de ipoteze „slabe“ așa cum au fost definite mai sus — dataset-uri arbitrare [formate din instanțe situate] pe axa reală.

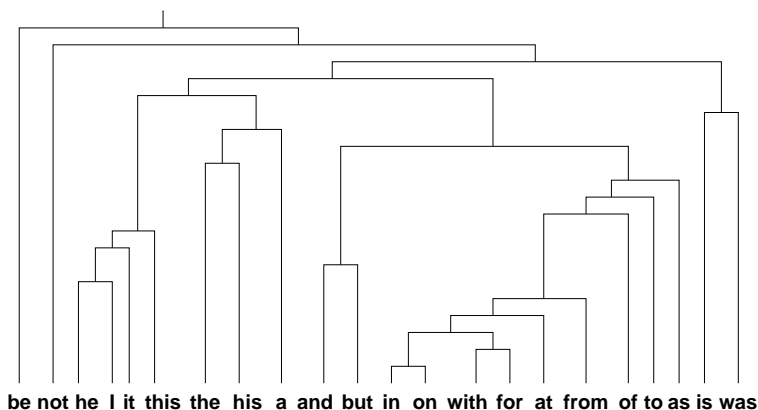
3. (Care dintre cei trei algoritmi de clusterizare — clusterizare ierarhică, K -means și EM/GMM — este adecvat pentru a clusteriza un anumit set de date?)

A. *Comentariu:*²

La seminar am analizat problema 61 de la capitolul *Clusterizare*, în care sunt prezentate mai multe seturi de date (vedeți figura de mai jos) și se cere să se precizeze care dintre cele trei tipuri de metode de clusterizare predate la curs — și anume, clusterizare ierarhică (cu similaritate *single*-, *complete*-, și respectiv *average-linkage*), algoritmul K -means și algoritmul EM pentru modele de tip mixtură de distribuții gaussiene³ (GMM) — este [mai] adecvat pentru a produce clusterurile care au fost indicate în această figură, folosind cele două culori, galben și albastru.



B. La curs am arătat că putem face clusterizarea celor mai frecvente cuvinte care apar într-un corpus lingvistic,⁴ folosind clusterizarea ierarhică. În acest scop, am arătat că se poate utiliza ca *măsură de similaritate* dintre două cuvinte oarecare numărul de „contexte“ (identice) în care apar aceste două cuvinte.⁵ Iată exemplul de *dendrogramă* care a fost prezentat la curs:



²Nu veți avea nimic de rezolvat la secțiunea aceasta, dar este indicat să o lecturați, în particular pentru a putea să răspundeți la punctul b de la secțiunea B.

³Aici, ne referim în mod implicit la cazul distribuțiilor gaussiene bivariate / bidimensionale. (Însă acest *specific* nu afectează linia raționamentului din exercițiul prezent.)

⁴Un *corpus lingvistic* este o colecție de documente. Exemplul dat la curs se referă la corpusul [Universității] Brown, foarte cunoscut și studiat în domeniul lingvisticii computaționale.

⁵Pentru un cuvânt oarecare dintr-un text dat, putem considera drept „context“ o „fereastră“ constituită dintr-un anumit număr de cuvinte la stânga și la dreapta cuvântului respectiv.

a. *Întrebare:* Putem oare să aplicăm algoritmul K -means cu distanța euclidiană și / sau algoritmul EM/GMM pe același set de date ca mai sus (i.e., corpusul Brown), în așa fel încât să obținem rezultate bune?

Precizări:

1. Întrucât algoritmi K -means și EM/GMM lucrează pe date numerice, pentru a putea răspunde la *întrebarea* de mai sus ne propunem să asociem fiecărui cuvânt câte o valoare numerică, și anume indicele poziției lui într-un dicționar / lexicon. (Practic, acesta este o listă formată din toate cuvintele din corpus.)

2. Pentru acești doi algoritmi, K — numărul de clustere în cazul lui K -means și, respectiv, numărul de componente ale mixturii în cazul lui EM/GMM — poate fi ales, de exemplu, ca fiind numărul „natural” de clustere stabilit în urma clusterizării ierarhice, așa cum am arătat de pildă la problema 1.b de la capitolul de *Clusterizare*.

Explicați succint, pentru fiecare dintre cei doi algoritmi de clusterizare menționați, motivele care au dus la acceptarea sau, dimpotrivă, la respingerea lui de către dumneavoastră ca răspuns la *întrebarea* pusă mai sus.

- K -means (cu distanța euclidiană):

- EM/GMM:

b. Cele două task-uri de clustering (cel de la secțiunea A și respectiv cel de la secțiunea B) din această problemă pun în evidență o anumită *caracteristică* a clusterizării ierarhice — în comparație cu clustering-ul partițional și cu clustering-ul bazat pe modelare probabilistă —, alta decât (1) capabilitatea de a produce ierarhii de clustere și (2) faptul că nu necesită fixarea în avans a numărului de clustere care trebuie formate.

Care credeți că este această *caracteristică*?

4.

(Algoritmul K -means: aplicare pe date din \mathbb{R}^2
monotonia criteriului B)

Considerăm setul de date din tabelul alăturat. Fiecare linie din tabel reprezintă un punct din \mathbb{R}^2 ; pentru conveniență, identificăm fiecare punct / instanță cu câte o literă.

a. Executați manual două iterații complete ale algoritmului K -means pe acest set de date, folosind $K = 3$ și distanța euclidiană. Clusterele vor fi inițializate după cum urmează:

$$C_1^{(0)} = \{A, B, F\}, C_2^{(0)} = \{C, H, I\}, C_3^{(0)} = \{D, E, G\}.$$

A	1	1
B	3	3
C	6	6
D	6	12
E	9	9
F	11	11
G	0	3
H	3	0
I	9	3

Veți face reprezentarea datelor pe gridurile de mai jos, și anume câte unul pentru fiecare iterație.

Indicație: În acele cazuri în care *metoda geometrică* [bazată pe reprezentarea mediatoarelor determinate de perechi de centroizi] nu este într-un totu concludentă în privința asignării unei instanțe oarecare la un anumit cluster / centroid, va trebui să faceți verificarea riguroasă, folosind *metoda analitică* [bazată pe calculul distanțelor].

b. La problema 42 am arătat că are loc următoarea *proprietate*: Dacă definim „criteriul variației inter-cluster” prin expresia

$$B(\mu^{(t)}) \stackrel{\text{def.}}{=} \sum_{j=1}^K \frac{n_j^{(t)}}{n} \|\mu_j^{(t)} - \bar{x}\|^2,$$

unde

n este numărul total de instanțe de clusterizat,

$n_j^{(t)}$ este numărul de instanțe din clusterul C_j la iterația t în momentul (re)calculării mediilor,

$\mu_j^{(t)}$ este media instanțelor din clusterul C_j , calculată la iterația t (înainte de reasignarea instanțelor la noii centroizi), iar

\bar{x} este media tuturor instanțelor de clusterizat (adică, centrul de greutate al întregii mulțimi de instanțe),

atunci acest „criteriu” $B(\mu^{(t)})$ este crescător (dar nu neapărat strict crescător) de la o iterație la alta.⁶

Folosind rezultatele obținute la punctul a, demonstrați (numeric!) că

$$B(\mu^{(1)}) \leq B(\mu^{(2)}).$$

Indicație: Pentru a rezolva cu mai multă ușurință această cerință, vă sugerăm să

– arătați că pe aceste date avem $\sum_{j=1}^3 \mu_j^{(1)} = \sum_{j=1}^3 \mu_j^{(2)}$, iar apoi

– dovediți că această egalitate implică (pe aceste date!) echivalența

$$B(\mu^{(1)}) \leq B(\mu^{(2)}) \Leftrightarrow \sum_{j=1}^3 \|\mu_j^{(1)}\|^2 \leq \sum_{j=1}^3 \|\mu_j^{(2)}\|^2.$$

Ultima inegalitate este ușor de verificat numeric (și chiar trebuie să o verificați!).

⁶De fapt, ceea ce am enunțat aici este o versiune simplificată a proprietății care a fost enunțată (și demonstrată) la problema 42.

Răspuns:

a.

Inițializare:

$$C_1^{(0)} = \{A, B, F\}$$

$$C_2^{(0)} = \{C, H, I\}$$

$$C_3^{(0)} = \{D, E, G\}.$$

Iterația 1:

$$\mu_1^{(1)} =$$

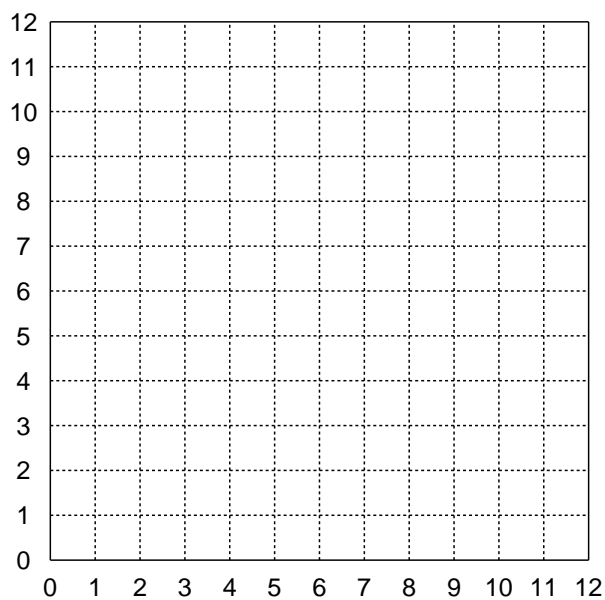
$$\mu_2^{(1)} =$$

$$\mu_3^{(1)} =$$

$$C_1^{(1)} =$$

$$C_2^{(1)} =$$

$$C_3^{(1)} =$$



Iterația 2:

$$\mu_1^{(2)} =$$

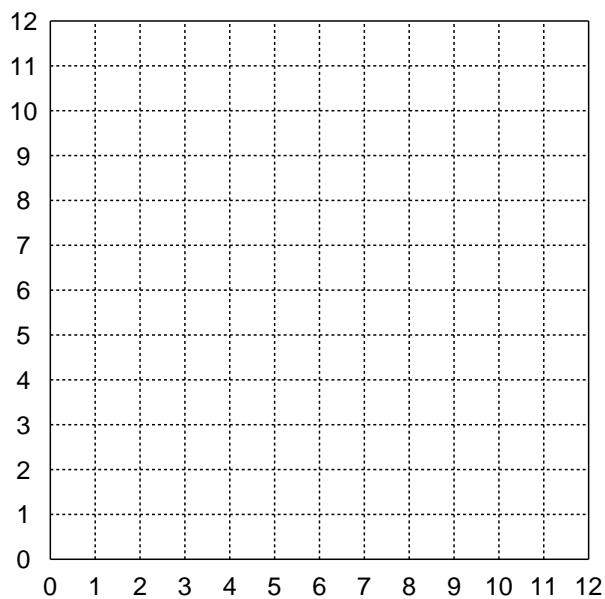
$$\mu_2^{(2)} =$$

$$\mu_3^{(2)} =$$

$$C_1^{(2)} =$$

$$C_2^{(2)} =$$

$$C_3^{(2)} =$$



b. ...

5. (EM/GMM, cazul unidimensional:
executarea manuală a pasului M (precum și a următorului pas E),
pentru o mixtură de un tip particular,
parametrii liberi fiind π_1, μ_1, μ_2)

Să zicem că dorim să antrenăm (engl., to fit) un model de tip mixtură de distribuții gaussiene unidimensionale, folosind $K = 2$ componente. Avem $n = 5$ instanțe, ale căror valori (notate cu x_i , cu $i = 1, \dots, 5$) sunt:

5, 15, 25, 30, 40.

Folosim algoritmul EM pentru a găsi estimările de verosimilitate maximă (engl., maximum likelihood estimates) pentru parametrii modelului, și anume: probabilitățile de selecție (sau, mixare) ale celor două componente (vom nota aceste probabilități cu π_1 și π_2) și mediile acestor componente (notate cu μ_1 și respectiv μ_2). Deviațiile standard ale celor două componente sunt considerate fixate și au valoarea 10.

Presupunem că în timpul execuției algoritmului EM, la o anumită iterație t , la pasul E au fost determinate mediile variabilelor neobservabile Z_{ij} — aceste medii se mai numesc *responsabilități*; engl., responsibilities — corespunzătoare celor două componente, pentru fiecare dintre cele cinci instanțe, după cum urmează:

$p_{i1}^{(t)}$	0.2	0.2	0.8	0.9	0.9
$p_{i2}^{(t)}$	0.8	0.8	0.2	0.1	0.1

a. Deduceți expresia funcției „auxiliare“ $Q(\pi_1, \mu_1, \mu_2 | \pi_1^{(t-1)}, \mu_1^{(t-1)}, \mu_2^{(t-1)})$ care urmează să fie optimizată la execuția pasului M al respectivei iterații (t) a algoritmului EM.

b. Ce valori vor fi atribuite parametrilor π_1, π_2, μ_1 și μ_2 la execuția respectivului pas M?

Veți optimiza funcția „auxiliară“ Q pe care ați determinat-o la punctul a. [În cazul în care nu ați rezolvat punctul a sau nu știți cum se optimizează respectiva funcție Q — dar doar în acest caz! —, puteți folosiți *formule* de genul celor pe care le-am predat la curs. Rezultatele trebuie să fie aceleași, în ambele variante!]

c. i. Deduceți expresiile necesare pentru a calcula $p_{i1}^{(t+1)}$ și $p_{i2}^{(t+1)}$, mediile variabilelor neobservabile Z_{ij} corespunzătoare următorului pas E al algoritmului EM.

ii. Calculați aceste probabilități pentru $i = 1$.

Indicație: Puteți folosi următoarele aproximări ale funcției de densitate de probabilitate (p.d.f.) pentru distribuția gaussiană standard: $\mathcal{N}(0.9 | 0, 1) = 0.2660$ și $\mathcal{N}(2.4 | 0, 1) = 0.0223$.