

Disciplina: Învățare automată Timp: 1 oră și 10 minute	TEST	Probabilități și statistică ID3 Clasificare bayesiană AdaBoost Clusterizare ierarhică
---	------	---

1. (0.05p) Verificați dacă următoarea definiție pentru funcția masă de probabilitate pentru variabila aleatoare X [care ia valori în mulțimea {1,2,3}] este validă:

$$p(1) = 0.2; p(2) = 0.3; p(3) = 0.4.$$

2. Fie următorul experiment aleator: aruncarea unei monede. Fie X variabila aleatoare asociată [cu valorile posibile 0 și 1: 0 pentru *tails*, 1 pentru *heads*].

- a. (0.05p) **Presupunând** că rezultatele sunt echiprobabile, asigurați probabilități ca X să ia o anumită valoare.

- b. (0.1p) Experimentul a fost repetat de 10 ori și s-au înregistrat următoarele date: 1,1,1,1,1,1,1,0,0,1. **Estimați** în sensul verosimilității maxime (MLE) probabilitățile ca X să ia o anumită valoare.

3. (0.15p) Presupunem ca funcția densitate de probabilitate a unei variabile aleatoare continue

$$p(x) = \begin{cases} \frac{4}{3}(1-x^3) & \text{pentru } 0 \leq x \leq 1 \\ 0 & \text{în caz contrar.} \end{cases}$$

X este definită astfel: . Cat este $P(X < 0.5)$?

4. (0.2p) Fie X și Y variabile aleatoare continue având funcția densitate de probabilitate comună

$$p(x, y) = \begin{cases} 1 & \text{pentru } 0 < x < 1, |y| < x \\ 0 & \text{în caz contrar.} \end{cases}$$

definită astfel:

- a) Cat este $p(y|x=0.5)$?

- b) Este variabila X independentă de variabila Y?

5. (0.2p) Fie o variabilă aleatoare cu distribuția de probabilitate uniformă, definită pe intervalul $[0, \frac{1}{2}]$, și anume $f(x)=2$ pentru x din intervalul $[0, \frac{1}{2}]$. Intrucat suma probabilitatilor care corespund unei distributii aleatoare trebuie sa fie 1, Mike este nedumerit de ce valoarea lui $f(x)$ este mai mare decat suma totala. Explicati acest paradox. Altfel spus, aratati ce anume nu stie Mike.

6. (0.3p) Formula de inlanturie a entropiilor pentru cazul general este:

$$H(X, Y) = H(X) + H(Y|X) = H(Y) + H(X|Y).$$

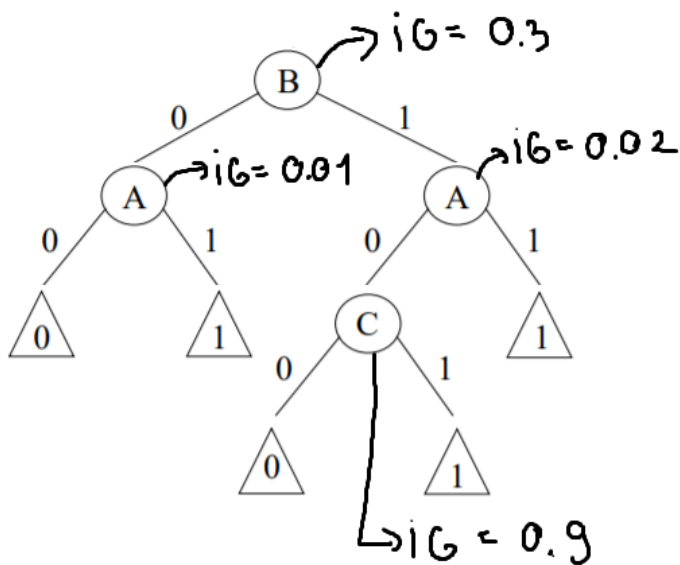
- a) (0.2p) Demonstrati ca daca X și Y sunt variabile aleatoare independente discrete, atunci $H(X, y) = H(X) + H(Y)$

- b) (0.1p) Este adevărata și reciprocă acestei afirmații? Adică, atunci când are loc egalitatea demonstrată la punctul a) rezulta că variabilele X și Y sunt independente?

7. (0.5p) În cadrul algoritmului **ID3**, fie următorul tabel în care doar atributele C și D sunt continue:

A	B	C (continuu)	D (continuu)	Y (output)
0	0	2	1	0
0	1	2	2	1
1	0	1	3	1
1	1	5	4	1
1	0	6	4	0
0	1	5	6	0

- a. (0.15p) Câte praguri distincte trebuie să considerăm pentru C atunci când căutăm atributul (optim) care trebuie pus în rădăcină? Dar pentru D? **Justificați.**
 - b. (0.05p) Care sunt nodurile ce trebuie luate în calcul atunci când se caută nodul rădăcină? **NU mai scrieți partițiile [a+,b-], ci doar numele din cercuri.**
 - c. Presupunem că valorile câștigurilor de informație dintre atributul de ieșire și fiecare atribut candidat pentru nodul rădăcină sunt: 0.9 pentru A, 0.8 pentru B, 0.7 pentru restul.
 - i. (0.05p) Ce atribut/nod va fi ales ca rădăcină? **Justificați.**
 - ii. (0.1p) Care sunt nodurile ce trebuie luate în calcul atunci când se caută nodul corespunzător ramurii *Rădăcină* == 0? **NU mai scrieți partițiile [a+,b-], ci doar numele din cercuri.**
 - d. (0.15p) Fie tabelul în care apar doar coloanele C și Y. Aplicăm ID3 **doar** pe aceste date. Cum vor fi clasificate instanțele (C=0) și (C=7)? **Justificați fără a face arborele!.**
8. (0.5p) Fie următorul arbore produs de ID3 pe un set de date. Pentru fiecare nod este trecut și IG-ul corespunzător calculat în cadrul algoritmului.



- a. (0.2p) Pentru un astfel de arbore de decizie, o strategie simplă de trunchiere (engl., pruning) în vederea contracarării fenomenului de “overfitting” constă în a parcurge arborele de sus în jos, începând deci cu nodul-rădăcină și identificând fiecare nod de test pentru care câștigul de informație are o valoare mai mică decât o valoare pozitivă, mică, fixată de la început, ϵ . Orice astfel de nod de test este imediat înlocuit — împreună cu subarborele corespunzător lui — cu un nod de decizie, conform etichetei majoritare a instanțelor asignate nodului de test (**în acest exercițiu veți lăsa nodul de decizie, adică triunghiul, gol, pentru că nu aveți informații suficiente**). Această strategie se numește “top-down pruning”.
Care este arborele de decizie obținut aplicând această strategie pe arborele de mai sus, dacă se consideră $\epsilon = 0.03$?
 - b. (0.2p) O altă posibilitate de a face pruning este să parcurgem arborele de decizie începând cu părinții nodurilor-frunză și să eliminăm în mod recursiv acele noduri de test pentru care câștigul de informație (sau un alt criteriu ales) este mai mic decât ϵ . Aceasta este strategia de “bottom-up pruning”.
Observație: Spre deosebire de strategia top-down, în varianta de pruning de tip bottom-up nu vor fi eliminate noduri (cu $IG < \epsilon$) pentru care există descendenți al căror câștig de informație este mai mare sau egal cu ϵ .
Ce arbore se obține făcând “bottom-up pruning” pe arborele dat mai sus, dacă se consideră $\epsilon = 0.03$?
 - c. (0.1p) În acest exercițiu am folosit următoarea metodă pentru a determina dacă un atribut de intrare este independent de atributul de ieșire: $IG < \epsilon$, unde ϵ e valoare mică, pozitivă apropiată de zero. Există o altă metodă de a verifica o astfel de independență în contextul trunchierii arborilor de decizie. Menționați numele ei.
9. (0.2p) Se consideră variabilele aleatoare X_1, X_2, X_3 și X_4 . Aceste variabile sunt independente condițional două câte două în raport cu variabila Y , cu excepția perechii X_3, X_4 . (Așadar, dacă am aplica algoritmul Bayes Naiv, acesta ar produce

erori de clasificare.) Cum am putea modifica regula de decizie a algoritmului Bayes Naiv pentru a ține cont de această particularitate a datelor?

10. (0.3p) În cadrul unui experiment de ML un student a împărțit setul de date de care dispunea în 3: antrenare, validare și testare. Studentul trebuie să aleagă între 3 modele: ID3, ID3 cu index Gini, Bayes Naiv și calculează următoarele erori la antrenare și validare după ce a antrenat pe setul de antrenare.

Model	Eroare la antrenare	Eroare la validare
ID3	0.05 = 5%	0.06 = 6%
ID3 cu index Gini	0.05 = 5%	0.4 = 40%
Bayes Naiv	0.54 = 54%	0.55 = 55%

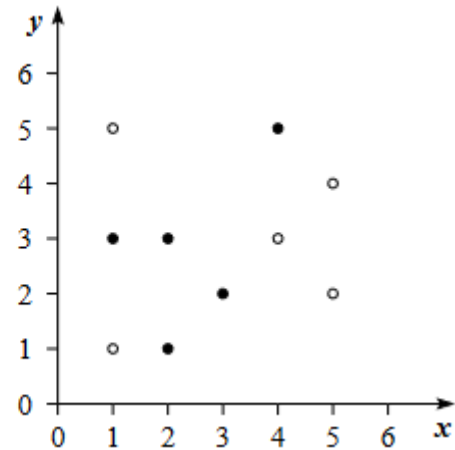
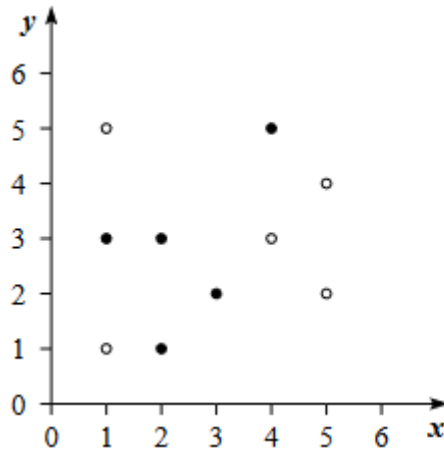
- (0.1p) Ce model va alege din cele 3? De ce?
- (0.1p) Ce probleme (underfitting/overfitting) prezintă celelalte două modele?
- (0.1p) Pentru modelul cu overfitting ce s-ar putea face pentru a scăpa de această problemă?

11. (0.25p) Fie următorul set de date, unde A,B,Y sunt discrete:

A	B	Y (output)
0	0	0
1	1	1
2	0	1
1	1	0

- (0.1p) Estimați $P(A = 0 | Y=1)$ în sensul verosimilității maxime (MLE).
- (0.15p) Estimați $P(A = 0 | Y=1)$ folosind regula *add-one* de netezire a lui Laplace.

12. (0.4p) Pe setul de date de mai jos desenați granițele de decizie și apoi hașurați suprafețele de decizie produse de:
- algoritmul 1-NN (veți obține deci diagrama Voronoi);
 - algoritmul ID3 extins cu capacitatea de a procesa attribute cu valori continue

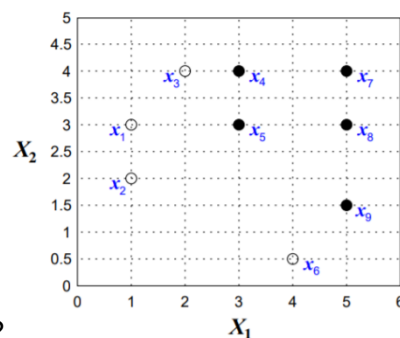


13. (0.4p) Folosind metoda 1-NN cu distanța euclidiană, învățăm un clasificator cu două valori pentru atributul de ieșire, $Y = 0$ și $Y = 1$, pornind de la datele de antrenament din tabelul de mai jos (X_1 și X_2 sunt atribute de intrare).

- (0.1p) Care este eroarea la antrenare (exprimată ca număr de exemple clasificate eronat)?
- (0.1p) Care este eroarea la cross-validare folosind metoda "Leave-One-Out"?
- (0.2p) Răspundeți la întrebările de mai sus, considerând acum arbori de decizie cu atribute numerice continue în locul metodei 1-NN.

X_1	X_2	Y
0	0	1
1	0	1
2	0	1
2.5	2	1
3	0	1
1	2	0
1.5	0	0
2	2	0
3	2	0
4	2	0

14. (0.4p) Aplicați 2 iterații de AdaBoost (deci, până la calculul distribuției D_3 inclusiv) pe următorul set de date. Care va fi eroarea la antrenare produsă de AdaBoost după



aceste 2 iterații?

15. (0.5p) Adevărat sau Fals? Justificați riguros.

- et, eroarea ponderată produsă la antrenare de către ipoteza h_t / clasificatorul „slab” A (măsurată relativ la ponderile de la începutul iterației t) tinde să crească în raport cu t .

- b) În decursul iterațiilor executate de algoritmul AdaBoost, erorile ponderate et (produse la antrenare, pe rând, de către ipotezele „slabe“ ht, în ocurență, compașii de decizie) pe de o parte, și erorile produse la antrenare de către clasificatorii combinați Ht pe de altă parte, variază aproximativ la fel.
- c) Ponderile / „voturile“ at asiguate de către algoritmul AdaBoost clasificatorilor „slabi“ ht asamblați sunt întotdeauna nenegative.
- d) robabilitățile / ponderile Dt(i) alocate de către algoritmul AdaBoost exemplilor de antrenament care au fost clasificate eronat [de către ipoteza
- e) Întotdeauna după ce algoritmul AdaBoost execută suficient de multe iterații, eroarea la antrenare produsă de ipoteza combinată Ht descrește la o valoare care este oricât [dorim să fie] de apropiată de zero, indiferent de tipul de clasificatori „slabi“ folosiți.

16. (1,2p) Tabelul de mai jos reprezintă matricea de distanțe pentru [o mulțime formată din] șase obiecte.

	A	B	C	D	E	F
A	0					
B	0.12	0				
C	0.51	0.25	0			
D	0.84	0.16	0.14	0		
E	0.28	0.77	0.70	0.45	0	
F	0.34	0.61	0.93	0.20	0.67	0

- a) Aplicați algoritmul de clusterizare ierarhică aglomerativă pe aceste date, folosind mai întâi similaritate single-linkage și apoi similaritate complete-linkage. La fiecare pas al algoritmului, rescrieți în mod corespunzător matricea de distanțe. (De la o iterație la alta, se micșorează cu 1 numărul liniilor precum și al coloanelor folosite.) La final, desenați dendrogramele rezultate [indicație: Înălțimea corespunzătoare fiecărui cluster non-singleton (adică, a fiecărui nod intern) din dendrogramă va fi considerată ca fiind egală cu distanța (i.e., conform măsurii de similaritate) dintre cele două sub-clustere constitutive.]
- b) Dacă ați lucrat corect, atunci cele două dendrograme obținute la punctul a nu coincid [nici măcar] ca structură. Modificați două valori din matricea de distanțe dată mai sus, în așa fel încât de data aceasta cele două dendrograme care obțin să fie identice ca structură
- c) Procedați similar pentru average-linkage. La actualizarea matricei de distanțe (sau, de „proximitate“) veți ține cont de formula:

$$\Delta(X \cup Y, Z) \stackrel{def.}{=} \frac{1}{(|X| + |Y|)|Z|} \sum_{x \in X \cup Y} \sum_{z \in Z} d(x, z)$$

$$\stackrel{calcul}{=} \frac{1}{|X| + |Y|} (|X| \Delta(X, Z) + |Y| \Delta(Y, Z)).$$

unde X, Y și Z

sunt clustere disjuncte două câte două, iar notația $|X|$ desemnează cardinalul lui X (adică, numărul de elemente din X).

d) Demonstrați formula enunțată la punctul c.