

Învățare automată

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

Nume student:

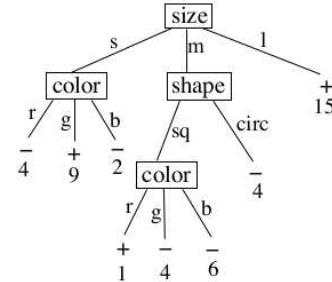
Grupa:

1. (Algoritmul ID3: 2 exemple simple de aplicare (antrenare + testare))

a. Considerăm arbore de decizie alăturat.

Numerele situate lângă etichetele $+/ -$ reprezintă cât de multe exemple de antrenament au fost asignate la nodul de decizie respectiv în cursul antrenării.

Folosind acest arbore de decizie, clasificați următoarele instanțe de test:



Size	Color	Shape	Smell	Classification
small	green	square	pine	
large	blue	square	pine	
medium	green	circle	rotten egg	
medium	red	square	lemon	

Observație: Pentru conveniență, pe arborele din imaginea de mai sus, valorile atributelor au fost scrise în mod prescurtat.

b. Acum vom presupune că aplicați algoritmul ID3 pe un set de date de antrenament (neprecizat) și că folosiți criteriul câștigului de informație maxim (sau, un criteriu echivalent!) cu acesta. Presupunem că

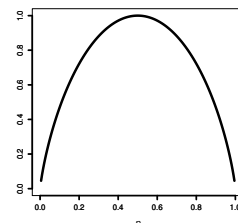
- setul de date de antrenament este constituit din 128 de exemple, dintre care 64 sunt etichetate pozitiv, iar 64 sunt etichetate negativ;
- instanțele de antrenament sunt caracterizate de trei atribute booleene: *Home-Owner*, *In-Debt* și *Rich*;
- pentru 64 de exemple de antrenament, atributul *Home-Owner* este *true*, iar dintre aceste exemple $1/4$ sunt negative și restul sunt pozitive;
- pentru 96 de exemple de antrenament, atributul *In-Debt* este *true*, iar dintre aceste exemple $1/2$ sunt pozitive și restul sunt negative;
- pentru 32 de exemple de antrenament, atributul *Rich* este *true*, iar dintre aceste exemple $3/4$ sunt pozitive și restul sunt negative.

i. Ce atribut trebuie pus în nodul rădăcină al arborelui ID3?

(În prealabil, veți desena compașii de decizie corespunzători celor trei atribute de intrare care au fost indicate mai sus.)

Observație: Atunci când veți calcula entropii și / sau câștiguri de informație, le veți exprima în biți.

ii. Puteți determina atributul care trebuie pus în nodul rădăcină al arborelui ID3 fără să faceți calcule efectiv [pentru entropii], ci doar bazându-vă pe monotonia funcției entropie ($H(p)$) pentru distribuția aleatoare Bernoulli de parametru p (vedeți graficul alăturat)?

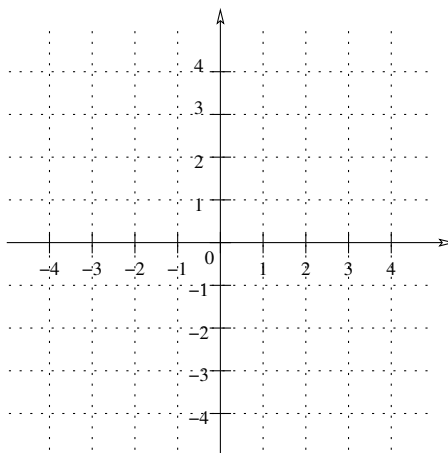


2. (Algoritmul k -NN, cu aplicare pe date din \mathbb{R}^2 , respectiv \mathbb{R} :
desenarea de diagrame Voronoi / identificarea zonelor de decizie;
calcularea erorii la CVLOO)

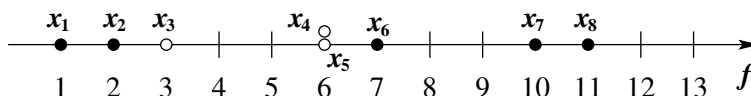
A. Se dă un set de date de antrenament din planul euclidian (\mathbb{R}^2), format din instanțele negative $(-1, 0)$, $(2, 1)$ și $(2, -2)$ și instanțele pozitive $(0, 0)$ și $(1, 0)$.

a. Pe reperul de axe de coordonate ortogonale alăturat marcați aceste exemple de antrenament. Veți folosi *convenția* noastră de notare: simbolul \bullet desemnează instanțe pozitive (+), iar simbolul \circ instanțe negative (-). Desenați apoi granițele de decizie determinate de algoritmul 1-NN. Veți hașura zona / zonele de decizie corespunzătoare instanțelor pozitive.

b. Cum va clasifica instanța de test $(1, -1.01)$ de către algoritmul 1-NN? Dar de către algoritmul 3-NN? Justificați, în ambele cazuri.

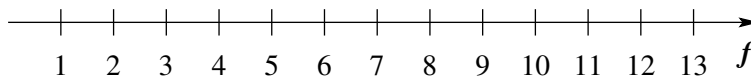


B. În desenul de mai jos este reprezentat un set de antrenament care conține 8 instanțe, fiecare dintre ele având doar o trăsătură (engl., feature), notată cu f .



Remarcați faptul că sunt două instanțe pentru care valoarea trăsăturii f este aceeași, și anume 6. Aceste două instanțe sunt reprezentate prin două simboluri \circ , situate unul deasupra celuilalt, dar de fapt ele ar fi trebuit să fie reprezentate ca două simboluri \circ suprapuse (unul peste celălalt), întrucât aceste instanțe au exact aceeași valoare pentru trăsătura f .

c. La acest punct al problemei veți folosi algoritmul 1-NN. Vă cerem ca pe linia de mai jos să hașurați — și, bineînțeles, să delimitați prin separatori decizionali — zonele în care algoritmul 1-NN va prezice semnul +, dat fiind setul de date de antrenament din figura de mai sus.



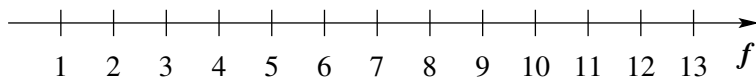
Cu ce etichetă vor fi clasificate instanțele de test $f = 2.5$ și $f = 6.5$? Justificați.

Atenție!

1. În cazul în care există două sau mai multe instanțe situate exact pe „marginea” [adică, pe conturul circular al] k -NN-vecinătății asociate instanței de clasificat, se va considera că toate aceste instanțe aparțin respectivei vecinătăți, iar fiecare dintre ele dispune de un vot întreg.
2. În caz de *paritate de voturi*, veți alege eticheta instanței de antrenament care este situată la *distanță minimă* către stânga față de instanța de test respectivă. (Distanță minimă este 0 atunci când instanța de test coincide cu o instanță de antrenament!)
3. Aceste reguli vor fi aplicate și la punctele următoare.

d. Dacă faceți cross-validare folosind metoda “leave-one-out” pe acest set de date, în conjuncție cu algoritmul 1-NN, cât va fi eroarea produsă? Veți arăta în mod clar — alcătuind un *tabel* care să conțină 1-NN-vecinătățile respective — cum anume ați ajuns la rezultatul pe care l-ați indicat.

e. Similar cu punctul c, însă aici veți folosi algoritmul 5-NN.¹



Indicație: Justificați în mod riguros rezultatul, referindu-vă la diverse intervale de valori (sau valori particulare) ale lui f :

Cazul 1: $f \in (-\infty, 4)$: ...

Cazul 2: $f = 4$: ...

...

Cu ce etichetă vor fi clasificate instanțele de test $f = 4$, $f = 6$ și $f = 7$? Justificați.

¹Fără să aplicați 5-NN pentru instanțele de test $f = 2.5$ și $f = 6.5$.

3. (Un exemplu de clasificator de tip bayesian care combină avantajele algoritmilor Bayes Naiv și Bayes Optimal)

Fie A, B și C variabile aleatoare binare independente, fiecare dintre ele având posibilitatea să ia valoarea 0 cu probabilitate de 50%. Considerăm funcția

$$Y = ((\neg B) \wedge (\neg C)) \vee (A \wedge B).$$

a. Scrieți *tabela de adevăr* a funcției Y .

b. i. Folosind *tabela de adevăr* a funcției Y — văzută acum ca un set de date de antrenament pentru un clasificator de tip bayesian — *estimați* valorile următoarelor probabilități condiționate:²

$P(A = 0 B = 0, Y = 0) = \dots$	$P(A = 0 C = 0, B = 0, Y = 0) = \dots$ $P(A = 0 C = 1, B = 0, Y = 0) = \dots$
$P(A = 0 B = 0, Y = 1) = \dots$	$P(A = 0 C = 0, B = 0, Y = 1) = \dots$ $P(A = 0 C = 1, B = 0, Y = 1) = \dots$
$P(A = 0 B = 1, Y = 0) = \dots$	$P(A = 0 C = 0, B = 1, Y = 0) = \dots$ $P(A = 0 C = 1, B = 1, Y = 0) = \dots$
$P(A = 0 B = 1, Y = 1) = \dots$	$P(A = 0 C = 0, B = 1, Y = 1) = \dots$ $P(A = 0 C = 1, B = 1, Y = 1) = \dots$

Atenție! Unele dintre aceste probabilități condiționate s-ar putea să nu fie definite. De exemplu, $P(A = 0|C = 0, B = 0, Y = 0)$ este nedefinită, fiindcă $P(C = 0, B = 0, Y = 0) = 0$, adică (veți vedea) nu există în tabelul de date nicio combinație de forma $C = 0, B = 0, Y = 0$. În dreptul unor astfel de probabilități condiționate nedefinite veți pune semnul $*$ (*nedefinit*).

ii. Analizând rezultatele estimărilor obținute mai sus, se poate deduce o relație de tip *independență condițională* între variabilele A, B, C și Y . Care este această relație?

c. i. Pe baza rezultatului de la punctul b.ii, scrieți *regula de decizie* pe care o poate folosi un clasificator de tip bayesian — diferit de Bayes Naiv și Bayes Optimal / Comun; să-i spunem *New-Bayes* — pentru învățarea conceptului / funcției Y , astfel încât

- să producă eroare 0 pe setul de date [de antrenament] de la punctul a, dar
- să necesite un număr de parametri „liberi” (engl., free) care trebuie estimați mai mic decât cel al algoritmului Bayes Optimal.

ii. *Justificați* în mod riguros regula de decizie pe care ați scris-o și că eroarea la antrenare produsă de la punctul a este într-adevăr 0.

iii. Care sunt acești parametri „liberi”?

iv. *Comparați* numărul acestor parametri „liberi” cu numărul minimal de parametri necesari de estimat de către clasificatorii Bayes Naiv și respectiv Bayes Optimal pe același set de date.

²Estimarea va fi făcută în mod clasic, adică în sensul verosimilității maxime (engl., Maximum Likelihood Estimation, MLE). Nu veți aplica nicio regulă de netezire a acestor probabilități, cum ar fi regula “add-one” a lui Laplace.