

Învățare automată

— Licență, anul III, 2018-2019, examenul parțial I —

Nume student:

Grupa:

1.

(Formula lui Bayes; inferențe statistice;
exemplificarea noțiunii de ipoteză / ipoteze MAP
("Maximum A posteriori Probability"))

A. Mickey dă cu zarul de mai multe ori, sperând să obțină un 6. Secvența celor 10 rezultate obținute de el în urma acestor aruncări este următoarea: 1, 2, 4, 3, 2, 2, 3, 5, 1, 6. Mickey se întreabă dacă nu cumva zarul este măsluit (având tendința să producă de mai multe ori fața 2 decât ar fi normal dacă zarul ar fi perfect).

Se presupune că în general fiecare set de 100 de zaruri conține 5 zaruri măsluite (engl., unfair) în așa fel încât este favorizată apariția feței 2, rezultând următoarea distribuție de probabilitate a celor șase fețe, (1, 2, 3, 4, 5, 6): $P = [0.1, 0.5, 0.1, 0.1, 0.1, 0.1]$.

Folosind teorema lui Bayes, furnizați-i lui Mickey informația care-l interesează: în ce măsură putem spune că zarul este măsluit [sau nu]?

Observație: Puteți folosi următoarele aproximări: $\log_2 3 = 1.585$; $\log_2 5 = 2.322$; $\log_2 19 = 4.248$.

(Algoritmul Bayes Naiv și algoritmul Bayes Corelat: aplicare)

B. Se dă setul de date alăturat, cu A și B variabile de intrare, iar C variabilă de ieșire.

a. Care este numărul minim de probabilități ce trebuie estimate pentru a putea construi după aceea (pe acest set de date) un clasificator de tip Bayes Naiv? Justificați.

b. Similar, pentru clasificatorul Bayes Corelat. Justificați.

A	B	C	nr. apariții
0	0	1	3
0	1	0	1
0	1	1	4
1	0	0	5
1	1	0	2
1	1	1	1

c. Care este *decizia* clasificatorului Bayes Naiv pentru $A = 0, B = 1$? Precizați cu ce *probabilitate* este luată această decizie.

d. Care este *decizia* clasificatorului Bayes Corelat pentru $A = 0, B = 1$? Precizați cu ce *probabilitate* este luată această decizie.

e. Dacă rezultatele obținute la punctele c și d diferă (fie și numai în privința *probabilităților* cu care sunt luate deciziile), care este explicația? Justificați în mod riguros.

2.

(Algoritmul ID3: aplicare;
calculul erorii la antrenare, respectiv la validare)

A. Dați definiția noțiunii de arbore de decizie.

B. Folosiți următorul set de date pentru a învăța cu ajutorul unui arbore de decizie dacă o ciupercă este sau nu comestibilă, utilizând atributele discrete *Formă*, *Culoare* și *Miros*.

<i>Formă</i>	<i>Culoare</i>	<i>Miros</i>	<i>Comestibilă</i>
C	B	1	1
D	B	1	1
D	W	1	1
D	W	2	1
C	B	2	1
D	B	2	0
D	G	2	0
C	U	2	0
C	B	3	0
C	W	3	0
D	W	3	0

Observație (1): Dacă la calculul entropiilor și / sau al câștigului de informație folosiți *alte formule decât* cele date în definiția acestor noțiuni, *enunțați-le* la cazul general și apoi *demonstrați-le* în mod succint.

Observație (2): Dacă veți folosi corect *definițiile*, veți constata că veți avea de făcut foarte puține calcule!

a. Calculați entropia condițională specifică $H(\text{Comestibilă} \mid \text{Miros} = 1 \text{ sau } \text{Miros} = 3)$.

Sugestie: Aplicați (deci, scrieți mai întâi) definiția noțiunii de entropie condițională specifică.

b. Ce atribut va alege algoritmul ID3 ca rădăcină a arborelui de decizie?

c. Elaborați întregul arbore de decizie care va fi învățat din datele de mai sus (fără pruning).

d. Exprimați cu ajutorul unui set de reguli din calculul propozițional clasificarea produsă de arborele de decizie obținut. (IF ... THEN *Comestibilă*; IF ... THEN \neg *Comestibilă*.)

e. Să presupunem că avem un set de date de validare:

<i>Formă</i>	<i>Culoare</i>	<i>Miros</i>	<i>Comestibilă</i>
C	B	2	0
D	B	2	0
C	W	2	1

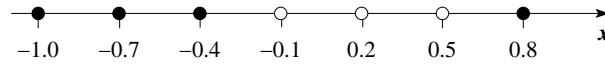
Care va fi eroarea produsă de arborele de decizie pe mulțimea de date de antrenare respectiv pe datele de validare? (Exprimați răspunsul ca număr de exemple clasificate greșit.)

3.

(Algoritmul AdaBoost: întrebări în legătură cu aplicarea algoritmului pe un set de date din \mathbb{R})

A. Demonstrați că minimul funcției $Z(\alpha_t) \stackrel{\text{def.}}{=} (1 - \varepsilon_t) \exp(-\alpha_t) + \varepsilon_t \exp(\alpha_t)$ se atinge pentru valoarea argumentului $\bar{\alpha}_t = \ln \sqrt{\frac{1 - \varepsilon_t}{\varepsilon_t}}$. Calculați $Z(\bar{\alpha}_t)$.

B. Considerăm setul de date de antrenament din figura următoare. Vom folosi algoritmul AdaBoost cu compași de decizie în rolul de ipoteze „slabe“.



a. Determinați separatorul decizional corespunzător primei ipoteze „slabe“, h_1 . Desenați-l pe figura de mai sus și indicați [eventual printr-o mică săgeată perpendiculară pe acest separator] care este zona clasificată cu +.

b. Calculați ε_1 și α_1 . Cât este acuratețea obținută de AdaBoost la antrenare dacă oprim acum algoritmul?

c. Cât va fi valoarea noilor probabilități / ponderi $D_2(i)$ pentru fiecare dintre cele șapte exemple de antrenament? (Atenție! Nu uitați să faceți normalizarea cu ajutorul factorului Z_1 .)

d. Desenați separatorul decizional corespunzător celei de-a doua ipoteze „slabe“, h_2 . Indicați iarăși zona de decizie corespunzătoare clasei +.

e. Care sunt exemplele de antrenament cărora le va fi asignată cea mai mică pondere / probabilitate după ce algoritmul AdaBoost va fi terminat cea de-a doua sa iterație?

f. Se îmbunătățește oare acuratețea la antrenare obținută de AdaBoost la a doua iterație în raport cu cea obținută la prima iterație?

g. Va putea oare AdaBoost să obțină eroare la antrenare 0 pe acest set de date? Justificați riguros, citând un rezultat teoretic prezentat la curs.

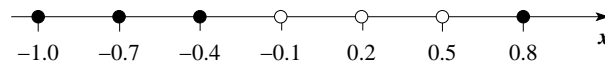
Observație importantă: Dacă veți folosi rezultate teoretice importante prezentate la curs, calculele pe care va trebui să le faceți la acest exercițiu vor fi foarte simple!

Răspuns:

[a.-c.] Iterația $t = 1$:

i	1	2	3	4	5	6	7
$D_1(i)$							

s	-1.15	-0.25	+0.65
$err_{D_1}(X < s)$			
$err_{D_1}(X \geq s)$			



$\varepsilon_1 =$

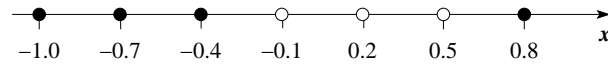
$\alpha_1 =$

i	1	2	3	4	5	6	7
$D_2(i)$							

$$Z_1 =$$

[d.-f.] Iterația $t = 2$:

s	-1.15	-0.25	+0.65
$err_{D_2}(X < s)$			
$err_{D_2}(X \geq s)$			



$$\varepsilon_2 =$$

$$\alpha_2 =$$

i	1	2	3	4	5	6	7
$D_3(i)$							

$$Z_2 =$$

t	α_t	x_1	x_2	x_3	x_4	x_5	x_6	x_7
1	$\alpha_1 =$							
2	$\alpha_2 =$							
	H							

Observație: În acest tabel, veți trece clasificările făcute de ipotezele „slabe“ h_t (pentru $t = 1$ și $t = 2$) și respectiv ipoteza „combinată“ H livrată de AdaBoost după efectuarea primelor două iterații. Veți încercui (în tabelul acesta) clasificările eronate, sub forma \oplus sau \ominus .

4.

(BONUS!

Cât de multe date de antrenament necesită algoritmul Bayes Naiv?

[LC: complexitatea la eşantionare])

Fie variabilele aleatoare $X_1, X_2, \dots, X_d \stackrel{\text{not.}}{=} Y$, fiecare dintre ele urmând o distribuție Bernoulli, de parametru $1/2$.

Comentariu: Variabilele X_1, X_2, \dots, X_d fiind de tip *Bernoulli*($1/2$), rezultă că a „observa“ (sau, a genera) o valoare (x_i) pentru o variabilă oarecare, X_i , fixată este echivalent cu aruncarea unei monede perfecte. Similar, a „observa“ (sau, a genera) o instanță etichetată $\bar{x} \stackrel{\text{not.}}{=} (X_1 = x_1, \dots, X_d = x_d)$ este echivalent cu aruncarea a d monede perfecte în mod independent (sau, echivalent, cu d aruncări independente ale unei monede perfecte).

Scopul acestui exercițiu este să calculăm N_{NB} , o margine inferioară pentru numărul de exemple de antrenament de care este nevoie [în prezentul context] astfel încât, dat fiind un prag $\varepsilon \in (0, 1)$, fixat, să putem spune că în setul de date reprezentat de cele N_{NB} exemple de antrenament probabilitatea să nu fi întâlnit toate valorile [tuturor] variabilelor X_1, X_2, \dots, X_d este mai mică decât ε . (Așadar, în acest context, cu probabilitate de $1 - \varepsilon$ putem estima toți parametrii algoritmului Bayes Naiv.)

Sugestie de lucru:

i. Considerând (din nou) X_i fixat, care este *probabilitatea* ca în N exemple de antrenament (unde $N \in \mathbb{N}$) să nu fi întâlnit ambele valori ale lui X_i ?¹

ii. Folosind inegalitatea $P(E_1 \cup \dots \cup E_n) \leq \sum_{i=1}^n P(E_i)$ care este valabilă pentru orice evenimente aleatoare E_1, \dots, E_n , și considerând din nou N fixat, calculați o *margine superioară* pentru *probabilitatea* ca în N exemple de antrenament să nu fi întâlnit toate valorile [tuturor] celor d variabile X_1, X_2, \dots, X_d .

Observație importantă: Pentru fixarea ideilor, la punctele *i.* și *ii.* puteți considera $d = 5$ și $N = 10$. Veți da răspunsurile atât în această variantă (adică, folosind aceste date concrete) cât și în cazul general.

iii. Dacă impuneți condiția ca marginea superioară a probabilității de la punctul *ii.* să fie egală cu un $\varepsilon > 0$ fixat, puteți obține — raționând analitic, pas cu pas, prin relații de echivalență — cât trebuie să fie N_{NB} , în funcție de d și ε (deci vom scrie $N_{NB}(d, \varepsilon)$) astfel încât, dat fiind un set de N_{NB} exemple de antrenament (generate cu ajutorul variabilelor aleatoare X_1, X_2, \dots, X_d), *probabilitatea* de a nu întâlni în acest set toate valorile [tuturor] variabilelor X_1, X_2, \dots, X_d să fie mai mică decât ε .

iv. Calculați efectiv marginile superioare $N_{NB}(d, \varepsilon)$ pentru $\varepsilon = 0.01$ și $d \in \{2, 5, 10\}$. Puteți folosi aproximarea $\log_2 5 = 2.3219$.

¹Echivalent spus, care este *probabilitatea* ca pentru X_i , în toate cele N exemple de antrenament, să apară în mod exclusiv(!) fie una fie cealaltă dintre valorile lui X_i ?