

3	6	2	1	2	6	5	1	6	1
9	6	0	1	2	6	3	1	2	6
1	1	1	3	1	2	3	1	2	6

Învățare automată

— Licență, anul III, 2016-2017, examenul parțial I —

Nume student: **Suciu Ana Maria** Grupa: **B2**

1. Variabile aleatoare discrete: calcul mediei; formula lui Bayes; măsuri statistice folosite în clasificare; ipoteze MAP)

O anumită boală afectează una din 500 de persoane în medie. Identificarea persoanelor care au această boală se poate face cu ajutorul unei analize a sângelui, care costă 100 de dolari de persoană. Această analiză indică în cazul unui rezultat pozitiv faptul că se poate ca persoana respectivă să sufere de acea boală.

Testul/analiza are o *sensibilitate* (engl., *sensitivity* sau *recall*) perfectă — adică raportul dintre numărul instanțelor pozitive identificate ca atare de acel test și numărul total de instanțe pozitive este 1 —, ceea ce înseamnă că pentru orice persoană care are boala respectivă, rezultatul testului este pozitiv cu probabilitate de 100%. Pe de altă parte, testul are o *specificitate* — raportul dintre numărul instanțelor negative identificate ca atare de acel test și numărul total de instanțe negative — de 99%, adică o persoană care nu suferă de acea boală va avea cu probabilitate de 1% rezultatul testului pozitiv.

a. Se testează o persoană selectată în mod aleatoriu, iar rezultatul este pozitiv. Care este probabilitatea ca persoana respectivă să sufere de acea boală? Interpretați (în manieră calitativă) rezultatul obținut.

b. Există și un al doilea test, care costă 10.000 de dolari și are atât sensibilitatea cât și specificitatea de 100%. Dacă am cere ca toate persoanele detectate pozitiv la testul precedent să fie supuse acestui test mult mai scump, care ar fi costul mediu pentru testarea/analiza unui individ?

Sugestie: Definiți următoarele variabile aleatoare:

B : ia valoarea 1/adevărat pentru persoanele care suferă de această boală și 0/fals în caz contrar
 T_1 : rezultatul primului test, care poate fi + (în caz de boală) sau -
 T_2 : rezultatul celui de-al doilea test, care poate fi tot + sau -

Folosind aceste variabile aleatoare, în vederea rezolvării problemei este bine ca mai întâi să formalizați enunțul sub forma următoare:

$$P(B) = \dots$$

$$P(T_1 = + | B) = \dots$$

$$P(T_1 = + | \bar{B}) = \dots$$

$$P(B) = \frac{1}{500} \Rightarrow P(\bar{B}) = 1 - P(B) = 1 - \frac{1}{500} = \frac{499}{500}$$

$$P(T_1 = + | B) = 1$$

$$P(T_1 = + | \bar{B}) = \frac{1}{100}$$

$$C_1 = 100 \$$$

$$C_2 = 10.000 \$$$

$$P(T_2 = + | B) = 1$$

$$P(T_2 = + | \bar{B}) = 0$$

a). $P(B | T_1 = +) = ?$ aplicăm formula lui Bayes

$$P(B | T_1 = +) = \frac{P(T_1 = + | B) \cdot P(B)}{P(T_1 = + | B) \cdot P(B) + P(T_1 = + | \bar{B}) \cdot P(\bar{B})}$$

$$= \frac{1 \cdot \frac{1}{500}}{1 \cdot \frac{1}{500} + \frac{1}{100} \cdot \frac{499}{500}} = \frac{1}{1 + 499} = \frac{1}{500}$$

b) Luăm o var aleatoare C , care reprezintă costul mediu

$$C = \begin{cases} C_1 & \text{dacă persoana nu suferă de boală} \\ C_1 + C_2 & \text{dacă persoana suferă de boală} \end{cases}$$

$$P(C = C_1) = P(T_1 = -) \quad P(C = C_1 + C_2) = P(T_1 = +)$$

$$E[C] = \sum_{c \in C} c \cdot p(c) = C_1 \cdot P(C = C_1) + (C_1 + C_2) \cdot P(C = C_1 + C_2)$$

$$= C_1 \cdot P(T_1 = -) + (C_1 + C_2) \cdot P(T_1 = +)$$

$$= C_1 \cdot P(T_1 = -) + (C_1 + C_2) \cdot (1 - P(T_1 = -))$$

$$= P(T_1 = -) (C_1 - C_1 - C_2) + C_1 + C_2$$

$$= C_1 + C_2 - C_2 \cdot P(T_1 = -)$$

$$P(T_1 = +) = P(T_1 = + | B) \cdot P(B) + P(T_1 = + | \bar{B}) \cdot P(\bar{B})$$

$$= 1 \cdot \frac{1}{500} + \frac{1}{100} \cdot \frac{499}{500} = \frac{599}{5000}$$

$$E[C] = C_1 \cdot P(T_1 = -) + (C_1 + C_2) \cdot P(T_1 = +)$$

$$= C_1 (1 - P(T_1 = +)) + (C_1 + C_2) \cdot P(T_1 = +)$$

$$= C_1 + P(T_1 = +) \cdot (C_1 + C_2 - C_1)$$

$$= C_1 + P(T_1 = +) \cdot C_2$$

$$= 100 + \frac{599}{5000} \cdot 5000 = 100 + \frac{599}{5} = 100 + 119.8 = 219.8 \approx 220 \$$$

\Rightarrow costul mediu = 220 \$