

Laboratorul 4

1. Descărcați următoarele liste de cuvinte cheie extrase din email-uri:

- *spam*: http://math.ubbcluj.ro/~miancu/teaching/PS/keywords_spam.txt
- *ham*: http://math.ubbcluj.ro/~miancu/teaching/PS/keywords_ham.txt

Implementați în Matlab o clasificare naivă Bayes care stabilește dacă un email este *spam* sau *ham*, calculând probabilitățile corespunzătoare date de frecvențele relative de apariții, respectiv neapariții, ale cuvintelor cheie în listele date (neglijați cuvintele din email care nu apar în liste). Testați programul cu următoarele email-uri:

- 1) “invite your friend today to click here”
- 2) “call your friend today it’s urgent thank you”.

Pentru rezolvarea acestei probleme este necesară înțelegerea exemplului de clasificare naivă Bayes din curs: http://math.ubbcluj.ro/~hanne/teaching/proba_stat_2019.pdf

Se pot urma pașii:

- *atributele* sunt date de cuvintele distincte din cele două fișiere: $W_1 : \text{call}, W_2 : \text{click}, \dots, W_{14} : \text{urgent}$; valorile atributelor sunt *true* sau *false*; exemplu: $W_2 = \text{true}$ reprezintă evenimentul că email-ul are cuvântul *click*; clasele sunt date de *ham* și *spam*; exemplu: $C = \text{spam}$ reprezintă evenimentul că email-ul este *spam*.
- se calculează probabilitățile claselor. Exemplu: probabilitatea ca un email să fie *spam* este:

$$P(C = \text{spam}) = \frac{\text{numărul de cuvinte din fișierul keywords_spam.txt}}{\text{numărul de cuvinte din ambele fișiere}}.$$

- se calculează probabilitățile atributelor, știind clasa. Exemplu: probabilitatea de apariție a lui W_1 într-un email, știind că email-ul este *spam*, este:

$$P(W_1 = \text{true} | C = \text{spam}) = \frac{\text{numărul de apariții ale cuvântului call în keywords_spam.txt}}{\text{numărul de cuvinte din keywords_spam.txt}},$$

iar probabilitatea de neapariție a lui W_1 într-un email, știind că email-ul este *spam*, este:

$$P(W_1 = \text{false} | C = \text{spam}) = 1 - P(W_1 = \text{true} | C = \text{spam}).$$

- pentru vectorul de attribute E_1 dat de cuvintele din primul email, se calculează produsele probabilităților de mai sus, iar pe baza formulei lui Bayes și a condițional independenței avem:

$$\begin{aligned} P(C = \text{spam} | E_1) &= \frac{P(E_1 | C = \text{spam})P(C = \text{spam})}{P(E_1)} \\ &= \frac{P(W_1 = \text{false}, W_2 = \text{true}, \dots, W_{14} = \text{false} | C = \text{spam})P(C = \text{spam})}{P(E_1)} \\ &= \frac{P(W_1 = \text{false} | C = \text{spam})P(W_2 = \text{true} | C = \text{spam}) \dots P(W_{14} = \text{false} | C = \text{spam})P(C = \text{spam})}{P(E_1)} \end{aligned}$$

și

$$\begin{aligned} &P(C = \text{ham} | E_1) \\ &= \frac{P(W_1 = \text{false} | C = \text{ham})P(W_2 = \text{true} | C = \text{ham}) \dots P(W_{14} = \text{false} | C = \text{ham})P(C = \text{ham})}{P(E_1)}. \end{aligned}$$

- se compară cele două probabilități (adică cei doi numărători ai fracțiilor de mai sus) și se decide clasificarea dată de probabilitatea mai mare.

2. Un număr natural aleatoriu N (mai mic decât 64) este generat conform unei rețele Bayes în care nodurile sunt variabile aleatoare B_i care indică valoarea bitului de pe poziția i , $i = \overline{1, 6}$, în reprezentarea binară a numărului N (numerotarea pozițiilor biților se face de la dreapta la stânga), cu următoarele probabilități:

B_1	$P(B_1 = \dots)$
1	0,8
0	0,2

B_2	$P(B_2 = \dots B_1 = 1)$	$P(B_2 = \dots B_1 = 0)$
1	0,9	0,6
0	0,1	0,4

B_3	$P(B_3 = \dots B_2 = 1, B_1 = 1)$	$P(B_3 = \dots B_2 = 1, B_1 = 0)$	$P(B_3 = \dots B_2 = 0, B_1 = 1)$	$P(B_3 = \dots B_2 = 0, B_1 = 0)$
1	0,6	0,2	0,9	0,4
0	0,4	0,8	0,1	0,6

B_4	$P(B_4 = \dots B_3 = 1)$	$P(B_4 = \dots B_3 = 0)$
1	0,3	0,5
0	0,7	0,5

B_5	$P(B_5 = \dots B_3 = 1)$	$P(B_5 = \dots B_3 = 0)$
1	0,5	0,8
0	0,5	0,2

B_6	$P(B_6 = \dots B_5 = 1, B_4 = 1)$	$P(B_6 = \dots B_5 = 1, B_4 = 0)$	$P(B_6 = \dots B_5 = 0, B_4 = 1)$	$P(B_6 = \dots B_5 = 0, B_4 = 0)$
1	0,5	0,3	0,8	0,5
0	0,5	0,7	0,2	0,5

a) Simulați de $n \in \{500, 1000, 1500\}$ ori valoarea numărului N . Afișați o histogramă a valorilor obținute.

b) Comparați rezultatele obținute cu cele teoretice pentru o valoare particulară a lui N .

Exemplu:

$$\begin{aligned}
 P(N = 23) &= P(B_6 = 0, B_5 = 1, B_4 = 0, B_3 = 1, B_2 = 1, B_1 = 1) \\
 &= P(B_6 = 0 | B_5 = 1, B_4 = 0) \cdot P(B_5 = 1 | B_3 = 1) \cdot P(B_4 = 0 | B_3 = 1) \\
 &\quad \cdot P(B_3 = 1 | B_2 = 1, B_1 = 1) \cdot P(B_2 = 1 | B_1 = 1) \cdot P(B_1 = 1) \\
 &= 0,7 \cdot 0,5 \cdot 0,7 \cdot 0,6 \cdot 0,9 \cdot 0,8 = 0,10584.
 \end{aligned}$$

Pentru rezolvarea acestei probleme este necesară înțelegerea exemplelor de rețea Bayes din curs: http://math.ubbcluj.ro/~hanne/teaching/proba_stat_2019.pdf