

Laboratorul 4

1. Salvați fișierele `keywords_spam.txt` și `keywords_ham.txt`, care conțin liste de cuvinte cheie extrase din email-uri *spam*, respectiv *ham*.

Implementați în Octave/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 studierea clasificării naive Bayes din curs.

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:

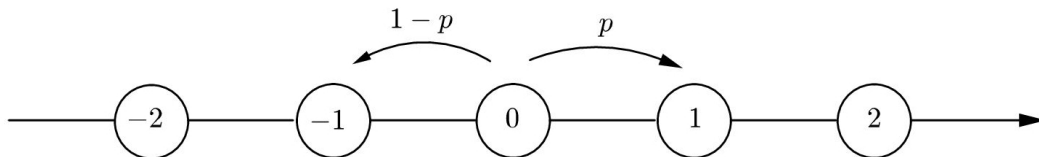
$$\begin{aligned} P(C = \text{spam} | E_1) &= \frac{P(E_1 | C = \text{spam})P(C = \text{spam})}{P(E_1)} \\ &= \frac{P(W_1 = \text{false}, W_2 = \text{true}, \dots, W_{14} = \text{false} | C = \text{spam})P(C = \text{spam})}{P(E_1)} \\ &= \frac{P(W_1 = \text{false} | C = \text{spam})P(W_2 = \text{true}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

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 punct material se deplasează pe axa reală dintr-un nod spre un nod vecin, la fiecare pas, cu probabilitatea $p \in (0, 1)$ la dreapta și cu probabilitatea $1 - p$ la stânga. Nodurile sunt centrate în numerele întregi, iar nodul inițial este 0:



a) Simulați o astfel de deplasare cu $k \in \mathbb{N}^*$ pași, cu probabilitatea $p \in (0, 1)$, și returnați pozițiile curente la fiecare pas.

b) Simulați de $m \in \mathbb{N}^*$ ori o astfel de deplasare cu $k \in \mathbb{N}^*$ pași, cu probabilitatea $p \in (0, 1)$, și afișați histograma pozițiilor finale. Care este poziția finală cel mai des întâlnită (sau pozițiile finale cel mai des întâlnite)?

3. Un jucător de “Loto 6/49” își cumpără câte un bilet pentru fiecare extrage efectuată de loteria română până când reușește să nimerească un bilet cu cel puțin 3 numere câștigătoare.

i) Folosind funcțiile `hygepdf` și `geornd`, generați un vector x care conține, pentru fiecare simulare, numărul de bilete necâștigătoare (care au cel mult 2 numere câștigătoare) până la primul bilet câștigător (care are cel puțin 3 numere câștigătoare).

ii) Estimați probabilitatea evenimentului:

“Cel puțin 10 bilete succesive sunt necâștigătoare până când jucătorul nimerește un bilet câștigător.”
Comparați probabilitatea estimată cu valoarea teoretică corespunzătoare, folosind funcția `geopdf`.