

Seminarul 4

Exemplu de clasificare naivă Bayes

Se dorește *clasificarea* traficului pe un anumit bulevard, în *clasele*: aglomerat sau relaxat, în funcție de următoarele *attribute* cu valorile lor posibile:

- vreme: ploaie, zăpadă, senin;
- temperatură: scăzută, medie, ridicată;
- timp: dimineață, amiază, seară, noapte.

Fie \mathbf{T} , V , Te , Ti variabilele aleatoare corespunzătoare și \mathbf{a} , \mathbf{r} , p , z , s , sc , m , ri , d , am , se și n valorile de mai sus, în ordinea în care sunt menționate.

În continuare, facem următoarea presupunere *naivă*: attributele sunt **condițional independente**, dacă se dă clasificarea.

De exemplu, avem:

$$P(V = p, Te = ri, Ti = d | \mathbf{T} = \mathbf{a}) = P(V = p | \mathbf{T} = \mathbf{a})P(Te = ri | \mathbf{T} = \mathbf{a})P(Ti = d | \mathbf{T} = \mathbf{a}),$$

$$P(Te = ri, Ti = d | \mathbf{T} = \mathbf{a}) = P(Te = ri | \mathbf{T} = \mathbf{a})P(Ti = d | \mathbf{T} = \mathbf{a}),$$

unde presupunem că $P(\mathbf{T} = \mathbf{a}) \neq 0$. Din acest exemplu se poate deduce că:

$$P(V = p | Te = ri, Ti = d, \mathbf{T} = \mathbf{a}) = P(V = p | \mathbf{T} = \mathbf{a}),$$

dacă $P(Te = ri, Ti = d, \mathbf{T} = \mathbf{a}) \neq 0$, altfel spus: probabilitatea să plouă, știind că traficul este aglomerat, nu depinde de temperatură sau timp. Într-adevăr, avem:

$$\begin{aligned} P(V = p | Te = ri, Ti = d, \mathbf{T} = \mathbf{a}) &= \frac{P(V = p, Te = ri, Ti = d, \mathbf{T} = \mathbf{a})}{P(Te = ri, Ti = d, \mathbf{T} = \mathbf{a})} \\ &= \frac{P(V = p, Te = ri, Ti = d | \mathbf{T} = \mathbf{a})P(\mathbf{T} = \mathbf{a})}{P(Te = ri, Ti = d | \mathbf{T} = \mathbf{a})P(\mathbf{T} = \mathbf{a})} = \frac{P(V = p | \mathbf{T} = \mathbf{a})P(Te = ri | \mathbf{T} = \mathbf{a})P(Ti = d | \mathbf{T} = \mathbf{a})}{P(Te = ri | \mathbf{T} = \mathbf{a})P(Ti = d | \mathbf{T} = \mathbf{a})}. \end{aligned}$$

Considerăm următorul *tabel de date* obținute în urma unor observații pe bulevard:

	<i>Vreme</i>	<i>Temperatură</i>	<i>Timp</i>	Trafic
1	ploaie	ridicăță	amiază	aglomerat
2	zăpadă	scăzută	seară	aglomerat
3	senin	ridicăță	noapte	relaxat
4	ploaie	scăzută	seară	aglomerat
5	senin	medie	amiază	relaxat
6	senin	scăzută	dimineață	aglomerat
7	ploaie	scăzută	amiază	aglomerat
8	ploaie	medie	noapte	relaxat
9	zăpadă	scăzută	noapte	relaxat
10	senin	ridicăță	seară	relaxat
11	zăpadă	scăzută	amiază	aglomerat
12	ploaie	medie	seară	aglomerat
13	ploaie	ridicăță	dimineață	aglomerat
14	zăpadă	medie	noapte	relaxat
15	senin	medie	dimineață	relaxat
16	ploaie	scăzută	noapte	aglomerat

i) Folosind datele din tabel, determinați probabilitățile claselor și probabilitățile condiționate ale atributelor, știind clasa.

ii) Considerăm evenimentul următor, denumit *vector de attribute*: $E = (V = p) \cap (Te = m) \cap (Ti = am)$. Alegeți o clasă pentru E , stabilind care din următoarele probabilități este mai mare: $P(\mathbf{T} = \mathbf{a}|E)$ sau $P(\mathbf{T} = \mathbf{r}|E)$.

Rezolvare:

i)

$\mathbf{T} = \mathbf{a}$	$\mathbf{T} = \mathbf{r}$	$P(\mathbf{T} = \mathbf{a})$	$P(\mathbf{T} = \mathbf{r})$
9	7	$\frac{9}{16}$	$\frac{7}{16}$

V	$\mathbf{T} = \mathbf{a}$	$\mathbf{T} = \mathbf{r}$	$P(V = \dots \mathbf{T} = \mathbf{a})$	$P(V = \dots \mathbf{T} = \mathbf{r})$
p	6	1	$\frac{6}{9}$	$\frac{1}{7}$
z	2	2	$\frac{2}{9}$	$\frac{2}{7}$
s	1	4	$\frac{1}{9}$	$\frac{4}{7}$

Te	$\mathbf{T} = \mathbf{a}$	$\mathbf{T} = \mathbf{r}$	$P(Te = \dots \mathbf{T} = \mathbf{a})$	$P(Te = \dots \mathbf{T} = \mathbf{r})$
sc	6	1	$\frac{6}{9}$	$\frac{1}{7}$
m	1	4	$\frac{1}{9}$	$\frac{4}{7}$
ri	2	2	$\frac{2}{9}$	$\frac{2}{7}$

Ti	$\mathbf{T} = \mathbf{a}$	$\mathbf{T} = \mathbf{r}$	$P(Ti = \dots \mathbf{T} = \mathbf{a})$	$P(Ti = \dots \mathbf{T} = \mathbf{r})$
d	2	1	$\frac{2}{9}$	$\frac{1}{7}$
am	3	1	$\frac{3}{9}$	$\frac{1}{7}$
se	3	1	$\frac{3}{9}$	$\frac{1}{7}$
n	1	4	$\frac{1}{9}$	$\frac{4}{7}$

ii) Pe baza formulei lui Bayes și a ipotezei de independență condiționată, deducem că:

$$\begin{aligned}
 P(\mathbf{T} = \mathbf{a}|E) &= \frac{P(E|\mathbf{T} = \mathbf{a})P(\mathbf{T} = \mathbf{a})}{P(E)} = \frac{P(V = p, Te = m, Ti = am|\mathbf{T} = \mathbf{a})P(\mathbf{T} = \mathbf{a})}{P(E)} \\
 &= \frac{P(V = p|\mathbf{T} = \mathbf{a})P(Te = m|\mathbf{T} = \mathbf{a})P(Ti = am|\mathbf{T} = \mathbf{a})P(\mathbf{T} = \mathbf{a})}{P(E)} = \frac{\frac{6}{9} \cdot \frac{1}{9} \cdot \frac{3}{9} \cdot \frac{9}{16}}{P(E)} = \frac{1}{P(E)} \cdot \frac{1}{72}
 \end{aligned}$$

și

$$\begin{aligned}
 P(\mathbf{T} = \mathbf{r}|E) &= \frac{P(E|\mathbf{T} = \mathbf{r})P(\mathbf{T} = \mathbf{r})}{P(E)} = \frac{P(V = p, Te = m, Ti = am|\mathbf{T} = \mathbf{r})P(\mathbf{T} = \mathbf{r})}{P(E)} \\
 &= \frac{P(V = p|\mathbf{T} = \mathbf{r})P(Te = m|\mathbf{T} = \mathbf{r})P(Ti = am|\mathbf{T} = \mathbf{r})P(\mathbf{T} = \mathbf{r})}{P(E)} = \frac{\frac{1}{7} \cdot \frac{4}{7} \cdot \frac{1}{7} \cdot \frac{7}{16}}{P(E)} = \frac{1}{P(E)} \cdot \frac{1}{196}.
 \end{aligned}$$

Deoarece $P(\mathbf{T} = \mathbf{a}|E) > P(\mathbf{T} = \mathbf{r}|E)$, asociem vectorului de attribute E clasa $\mathbf{T} = \mathbf{a}$.

Probleme cu variabile aleatoare discrete și caracteristici numerice

1. Cinci cartonașe sunt numerotate cu 2, 4, 6, 8, respectiv 10. Se aleg aleator trei cartonașe, repunând de fiecare dată înapoi în teanc cartonașul ales. Știind că suma numerelor extrase este 12 și notând cu X variabila aleatoare care indică de câte ori a fost extras numărul 2, determinați:

- a) funcția de repartiție a lui X ;
- b) valoarea medie a lui X ;
- c) varianța lui X .

$$\text{R: a) } X \sim \begin{pmatrix} 0 & 1 & 2 \\ \frac{1}{10} & \frac{6}{10} & \frac{3}{10} \end{pmatrix}, \text{ deci } F(x) = P(X \leq x) = \begin{cases} 0, & x < 0, \\ \frac{1}{10}, & 0 \leq x < 1, \\ \frac{7}{10}, & 1 \leq x < 2, \\ 1, & x \geq 2. \end{cases} \quad \text{b) } E(X) = \frac{6}{5}. \quad \text{c) } V(X) =$$

$$E(X^2) - E^2(X) = \frac{18}{10} - \frac{36}{25} = \frac{9}{25}.$$

2. Un sistem electronic are $n \in \mathbb{N}^*$ componente care funcționează independent unele de altele. Fiecare componentă funcționează cu probabilitatea $p \in (0, 1)$. Fie X variabila aleatoare care indică numărul de componente funcționale ale sistemului. Determinați distribuția lui X și apoi calculați valoarea sa medie și varianța sa.

Consecință: Dacă $X \sim \text{Bino}(n, p)$, atunci $E(X) = np$ și $V(X) = np(1 - p)$.

R: $E(X) = E(X_1) + \dots + E(X_n) = np$ și $V(X) = V(X_1) + \dots + V(X_n) = np(1 - p)$, unde $X_i \sim \text{Bernoulli}(p)$ indică funcționarea componentei i , $i = \overline{1, n}$.

3. Un mesaj este transmis printr-un canal de comunicație cu perturbări. Probabilitatea ca mesajul să fie recepționat complet este $p \in (0, 1)$. Dacă mesajul nu este recepționat complet, atunci se reia transmisia mesajului, independent de transmisiile anterioare. Fie X variabila aleatoare care indică numărul de transmisi până la prima transmisie în care mesajul este recepționat complet. Determinați valoarea medie și varianța lui X .

Consecință: Dacă $X \sim \text{Geo}(p)$, atunci $E(X) = \frac{1-p}{p}$ și $V(X) = \frac{1-p}{p^2}$.

R: $E(X) = \sum_{k=0}^{\infty} kp(1-p)^k = (1-p) \sum_{k=1}^{\infty} (k-1)p(1-p)^{k-1} + p \sum_{k=1}^{\infty} (1-p)^k = (1-p)E(X) + 1-p$, deci $E(X) = \frac{1-p}{p}$.

$V(X) = E(X^2) - E^2(X) = E(X^2) - \frac{(1-p)^2}{p^2}$; $E(X^2) = \sum_{k=0}^{\infty} k^2 p(1-p)^k = (1-p) \sum_{k=1}^{\infty} (k^2 - 2k + 1)p(1-p)^{k-1} + 2 \sum_{k=1}^{\infty} kp(1-p)^k - p \sum_{k=1}^{\infty} (1-p)^k = (1-p)E(X^2) + 2E(X) - (1-p) = (1-p)E(X^2) + \frac{(1-p)(2-p)}{p}$, deci $E(X^2) = \frac{(1-p)(2-p)}{p^2}$ și $V(X) = \frac{(1-p)(2-p)}{p^2} - \frac{(1-p)^2}{p^2} = \frac{1-p}{p^2}$.

4. O monedă este aruncată de n ori ($n \in \mathbb{N}^*$). Fie X variabila aleatoare care indică diferența dintre numărul de capete și numărul de pajuri obținute. Determinați:

- i) distribuția lui X ;
- ii) valoarea medie a lui X ;
- iii) varianța lui X .

R: Dacă C indică numărul de capete, atunci $C \sim \text{Bino}(n, \frac{1}{2})$ și $X = C - (n - C) = 2C - n$. $E(X) = E(2C - n) = 2E(C) - n = 2 \cdot \frac{n}{2} - n = 0$. $V(X) = 4V(C) = 4 \cdot \frac{n}{4} = n$.

5. Trei prieteni decid cine va plăti nota la restaurant astfel: fiecare aruncă pe rând o monedă; plătește cel care obține un simbol diferit de al celorlalți doi; dacă toți au obținut același simbol, atunci se reia seria de trei aruncări ale monedei, ș.a.m.d. până se decide cine plătește nota. Determinați valoarea medie a numărului de serii de aruncări până la seria de aruncări care va decide cine plătește nota.

R: $X \sim \text{Geom}(\frac{3}{4})$, $E(X) = \frac{1}{3}$.

6. O persoană are două cutii de Tic Tac în buzunar. De fiecare dată când dorește un drajeu, scoate aleator din buzunar una dintre cutii, din care ia un drajeu. La un moment dat scoate din buzunar o cutie și constată că e goală. Fie X numărul de drajeuri din cealaltă cutie. Determinați distribuția lui X , știind că fiecare cutie conținea inițial câte n drajeuri ($n \in \mathbb{N}$).

$$\text{R: } X \sim \left(\binom{k}{C_{2n-k}^n \frac{1}{2^{2n-k}}} \right)_{k=\overline{0,n}}.$$

7. n bile sunt puse aleator în k cutii ($k, n \in \mathbb{N}^*$), numerotate de la 1 la k , astfel: pentru fiecare bilă se alege aleator o cutie din cele k astfel încât probabilitatea de a alege cutia i să fie p_i , $i = \overline{1,k}$, unde $p_i \in (0, 1)$, $i = \overline{1,k}$, sunt fixate astfel încât $p_1 + \dots + p_k = 1$. Determinați valoarea medie a numărului de cutii care conțin exact m bile, $m \in \mathbb{N}$, $m \leq n$.

$$\text{R: Fie } X_i = \begin{cases} 1, & \text{cutia } i \text{ conține exact } m \text{ bile, } i = \overline{1,k}. \\ 0, & \text{altfel,} \end{cases} \text{ Avem } P(X_i = 1) = C_n^m p_i^m (1 - p_i)^{n-m},$$

unde am folosit modelul urnei cu bila întoarsă cu 2 stări și n extrageri. Deci $E(X) = E(\sum_{i=1}^k X_i) = \sum_{i=1}^k C_n^m p_i^m (1 - p_i)^{n-m}$.

Aplicații ale valorii medii în teoria grafurilor

1. Fie X o variabilă aleatoare discretă care ia un număr finit de valori. Atunci $P(X \geq E(X)) \neq 0$ și $P(X \leq E(X)) \neq 0$.

R: Fie $\{x_1, \dots, x_n\}$ mulțimea de valori pe care le poate lua X cu probabilități nenule. Dacă am presupune că $P(X \geq E(X)) = 0$, atunci am avea $x_i < E(X)$, pentru fiecare $i \in \{1, \dots, n\}$, și astfel am obține următoarea contradicție:

$$E(X) = \sum_{i=1}^n x_i \cdot P(X = x_i) < \sum_{i=1}^n E(X) \cdot P(X = x_i) = E(X).$$

Similar se demonstrează cealaltă relație.

Definiții: În următoarele probleme vom folosi următoarele definiții:

- $G = (V, E)$ este un graf neorientat cu V mulțimea nodurilor/vârfurilor și E mulțimea muchiilor;
- $n = |V| \geq 2$ și $m = |E| \geq 1$;
- G este complet, dacă oricare două noduri sunt adiacente;
- o muchie este bicoloră, dacă nodurile ei sunt colorate diferit;
- un subgraf este monocolor, dacă toate muchiile sale sunt colorate la fel;
- $I \subseteq V$ este independentă, dacă nu există două noduri în I care să fie adiacente.

2. Demonstrați că există o colorare cu roșu și albastru a nodurilor lui G astfel încât cel puțin $\frac{m}{2}$ muchii sunt bicolore.

Idee de rezolvare: colorați aleator nodurile și calculați valoarea medie a numărului de muchii bicolore.

R: Colorăm aleator nodurile cu ajutorul unei monede: se aruncă moneda pentru fiecare muchie; dacă obținem *cap* colorăm muchia cu roșu, dacă obținem *pajură* colorăm muchia cu albastru.

Fie $X_i = \begin{cases} 1, & \text{muchia } i \text{ este bicoloră,} \\ 0, & \text{altfel.} \end{cases}$ Observăm că $X = \sum_{i=1}^m X_i$ indică numărul de muchii bicolore.

Avem $E(X) = \sum_{i=1}^m E(X_i) = \sum_{i=1}^m \frac{1}{2} = \frac{m}{2}$. Problema 1. implică $P(X \geq \frac{m}{2}) \neq 0$, deci există cel puțin o colorare cu cel puțin $\frac{m}{2}$ muchii bicolore.

3. Dacă G este complet și avem $2C_n^k < 2C_k^2$, unde $k \in \mathbb{N}^*$, atunci există o colorare cu roșu și albastru a muchiilor astfel încât niciun subgraf complet cu k noduri nu este monocolor.

Idee de rezolvare: colorați aleator muchiile și calculați valoarea medie a numărului de subgrafuri complete care sunt monocore.

R: Colorăm aleator muchiile cu ajutorul unei monede (similar procedului din problema anterioară). Fie G_1, \dots, G_N subgrafurile complete cu k noduri ale grafului G , unde $N = C_n^k$.

Fie $X_i = \begin{cases} 1, & \text{subgraful } G_i \text{ este monocolor,} \\ 0, & \text{altfel.} \end{cases}$ Observăm că $X = \sum_{i=1}^N X_i$ indică numărul de subgrafuri

monocore. Avem $E(X) = \sum_{i=1}^N E(X_i) = \sum_{i=1}^N 2 \cdot \left(\frac{1}{2}\right)^{C_k^2} = \frac{2C_n^k}{2^{C_k^2}} < 1$. Problema 1. implică $P(X \leq E(X)) \neq 0$, deci $P(X < 1) = P(X = 0) \neq 0$, așadar există cel puțin o colorare astfel încât niciun subgraf complet cu k noduri nu este monocolor.

4. Demonstrați că există o mulțime $I \subseteq V$ independentă astfel încât $|I| \geq \frac{n^2}{4m}$.

Idee de rezolvare: considerați algoritmul aleator (Turan):

1. Delete each vertex of G (together with its incident edges) independently with probability $1-1/d$, where $d=2m/n$.
2. For each remaining edge, remove it and one of its adjacent vertices.

Fie X și Y numărul de noduri, respectiv de muchii, care supraviețuiesc după prima etapă din algoritm. Calculați $E(X - Y)$.

R: Observăm că *output-ul* algoritmului este o mulțime de noduri independentă I .

$$\text{Fie } X_i = \begin{cases} 1, & \text{nodul } i \text{ supraviețuiește după prima etapă,} \\ 0, & \text{altfel,} \end{cases} \quad i = \overline{1, n}, \text{ și}$$

$$Y_j = \begin{cases} 1, & \text{muchia } j \text{ supraviețuiește după prima etapă,} \\ 0, & \text{altfel,} \end{cases} \quad j = \overline{1, m}. \text{ Observăm că } E(X_i) = \frac{1}{d} \text{ și } E(Y_j) = \frac{1}{d^2},$$

deoarece o muchie supraviețuiește dacă și numai dacă supraviețuiesc nodurile ei. Deci $E(X - Y) = E(X) - E(Y) = \sum_{i=1}^n E(X_i) - \sum_{j=1}^m E(Y_j) = \frac{n}{d} - \frac{m}{d^2}$. Numărul de noduri Z care supraviețuiesc după etapa a doua este cel puțin egal cu $X - Y$, deci $E(Z) \geq \frac{nd-m}{d^2} = \frac{4m}{n^2}$. Problema 1. implică $P(Z \geq \frac{4m}{n^2}) \neq 0$, deci există o mulțime I , returnată de algoritm, care are cel puțin $\frac{4m}{n^2}$ noduri.

Observații:

- Pentru $k, l \in \mathbb{N}^*$, numărul lui Ramsey $R(k, l)$ este definit astfel: $R(k, l)$ este cel mai mic număr natural n cu proprietatea că orice colorare cu roșu și albastru a muchiilor grafului complet G cu n vârfuri conține un subgraf complet cu k noduri colorat cu roșu sau un subgraf complet cu l noduri colorat cu albastru.

- Numărul lui Ramsey $R(k, l)$ se poate defini și astfel: este cel mai mic număr de persoane pe care trebuie să le invitați la o petrecere astfel încât să aveți cel puțin un grup de k persoane în care oricare două se cunosc sau cel puțin un grup de l persoane în care oricare două nu se cunosc.

- Problema 3. ne spune: dacă $2C_n^k < 2C_k^2$, atunci $R(k, k) > n$.

- Sunt cunoscute doar câteva valori ale numerelor lui Ramsey. De exemplu: $R(1, 1) = 1$, $R(2, 2) = 2$, $R(3, 3) = 6$, $R(4, 4) = 18$. Numerele lui Ramsey sunt recunoscute pentru dificultatea aflării valorilor lor. Valorile $R(k, k)$ pentru $k \geq 5$ sunt încă necunoscute.

- “Erdős asks us to imagine an alien force, vastly more powerful than us, landing on Earth and demanding the value of $R(5, 5)$ or they will destroy our planet. In that case, he claims, we should marshal all our computers and all our mathematicians and attempt to find the value. But suppose, instead, that they ask for $R(6, 6)$. In that case, he believes, we should attempt to destroy the aliens” - din cartea Joel H. Spencer - *Ten Lectures on the Probabilistic Method* (1994).