

## Seminarul 4

• **Modelul urnei cu bile de 2 culori și bilă nereturnată:** fie  $n_1, n_2, n \in \mathbb{N}$  cu  $n \leq n_1 + n_2$  și fie  $k \in \mathbb{N}$  astfel încât  $k \leq n_1$  și  $n - k \leq n_2$ ; considerând o urnă, care are inițial  $n_1$  bile albe și  $n_2$  bile negre, avem

$$\begin{aligned} p(k; n) &= \text{probabilitatea de a obține } k \text{ bile albe din } n \text{ extrageri fără returnarea bilei extrase,} \\ &\quad \text{în care ordinea de extragere a bilelor nu contează} \\ &= \frac{C_{n_1}^k \cdot C_{n_2}^{n-k}}{C_{n_1+n_2}^n}. \end{aligned}$$

▷ Acest model corespunde **distribuției hipergeometrice**.

1. Dintr-un set de 52 de cărți de joc se extrag aleator, pe rând, fără returnare, 13 cărți (*bridge hand*). Calculați probabilitățile următoarelor evenimente:

- a)  $A$ : nu s-a extras nicio treflă;
- b)  $B$ : s-au obținut 5 inimi;
- c)  $C$ : s-a obținut cel mult un as.

• **Modelul urnei cu  $r$  culori și bilă nereturnată:** fie  $n_i$  = numărul inițial de bile cu culoarea  $i$  din urnă,  $i = \overline{1, r}$ ;

$$\begin{aligned} p(k_1, \dots, k_r; n) &= \text{probabilitatea de a obține } k_i \text{ bile cu culoarea } i, i = \overline{1, r}, \\ &\quad \text{din } n = k_1 + \dots + k_r \text{ extrageri fără returnarea bilei extrase,} \\ &\quad \text{în care ordinea de extragere a bilelor de diverse culori nu contează} \\ &= \frac{C_{n_1}^{k_1} \cdot \dots \cdot C_{n_r}^{k_r}}{C_{n_1+\dots+n_r}^n}. \end{aligned}$$

▷ Cazul  $r = 2$  corespunde **distribuției hipergeometrice**.

**Observație:** Extragerea fără returnare (engl. *sampling without replacement*) este folosită în **metoda validării încrucișate** (engl. *k-fold cross validation*): În cazul validării încrucișate eșantionul original de date este împărțit aleatoriu în  $k$  sub-eșantioane de dimensiuni egale. Din cele  $k$  sub-eșantioane, un singur sub-eșantion este folosit ca date de validare pentru testarea modelului, iar celelalte  $k - 1$  sub-eșantioane sunt utilizate ca date de antrenament. Procesul de validare încrucișată se repetă apoi de  $k$  ori, fiecare dintre cele  $k$  sub-eșantioane fiind utilizat exact o dată ca date de validare.

2. O echipă formată din 4 cercetători este aleasă aleator dintr-un grup de 4 matematicieni, 3 informaticieni și 5 fizicieni. Care este probabilitatea ca echipa să fie formată din 2 matematicieni, 1 informatician și 1 fizician?

### 3. Clasificarea naivă Bayes

Clasificatorii bayesieni naivi sunt o familie de clasificatori probabilistici simpli, bazați pe aplicarea formulei lui Bayes cu ipoteze “naive” de **independență condiționată între atribute** (engl. *features*), cunoscând **clasificarea**. În aplicații practice pentru modelele bayesiene naive se folosește **metoda probabilității maxime**. Noțiunea folosită în acest context este condițional independența între v.a.

**Def. 1** Fie  $U, X, Y, Z$  v.a. discrete, care iau valori în mulțimile  $\mathcal{U}, \mathcal{X}, \mathcal{Y}, \mathcal{Z}$ . V.a.  $U, X, Y$  sunt **condițional independente**, cunoscând (știind) v.a.  $Z$ , dacă pentru fiecare  $u \in \mathcal{U}, x \in \mathcal{X}, y \in \mathcal{Y}, z \in \mathcal{Z}$  are loc

$$P(U = u, X = x, Y = y | Z = z) = P(U = u | Z = z)P(X = x | Z = z)P(Y = y | Z = z).$$

Considerăm următoarea problemă de *clasificare naivă Bayes* a unor restaurante ( $\mathbf{R}$ ), în

- *clasele*: recomandat sau nerecomandat,
- în funcție de următoarele *atribute* cu valorile lor posibile:
- cost ( $C$ ): ieftin, mediu, scump;
- timp de așteptare ( $T$ ): puțin, mediu, îndelungat;
- mâncare ( $M$ ): fadă, acceptabilă, bună, delicioasă.

$\mathbf{R}$ ,  $C$ ,  $T$ ,  $M$  sunt variabilele aleatoare (catoriale) și  $\mathbf{r}$ ,  $\mathbf{n}$ ,  $i$ ,  $m$ ,  $s$ ,  $p$ ,  $m$ ,  $\hat{i}$ ,  $f$ ,  $a$ ,  $b$ ,  $d$  valorile de mai sus, în ordinea în care sunt menționate.

Considerăm următorul *tabel de date* furnizat de clienții unor restaurante:

	<i>Cost</i>	<i>Timp de așteptare</i>	<i>Mâncare</i>	<b>Restaurant</b>
1	mediu	îndelungat	acceptabilă	<b>nerecomandat</b>
2	scump	puțin	bună	<b>recomandat</b>
3	ieftin	îndelungat	delicioasă	<b>recomandat</b>
4	mediu	puțin	bună	<b>recomandat</b>
5	ieftin	mediu	acceptabilă	<b>nerecomandat</b>
6	ieftin	puțin	fadă	<b>nerecomandat</b>
7	mediu	puțin	acceptabilă	<b>nerecomandat</b>
8	mediu	mediu	delicioasă	<b>recomandat</b>
9	scump	puțin	delicioasă	<b>recomandat</b>
10	ieftin	îndelungat	bună	<b>nerecomandat</b>
11	scump	puțin	acceptabilă	<b>nerecomandat</b>
12	mediu	mediu	bună	<b>recomandat</b>
13	mediu	îndelungat	fadă	<b>nerecomandat</b>
14	scump	mediu	delicioasă	<b>recomandat</b>
15	ieftin	mediu	fadă	<b>nerecomandat</b>
16	mediu	puțin	delicioasă	<b>recomandat</b>
17	ieftin	puțin	acceptabilă	<b>recomandat</b>
18	scump	îndelungat	bună	<b>nerecomandat</b>
19	ieftin	puțin	fadă	<b>recomandat</b>
20	scump	îndelungat	delicioasă	<b>nerecomandat</b>

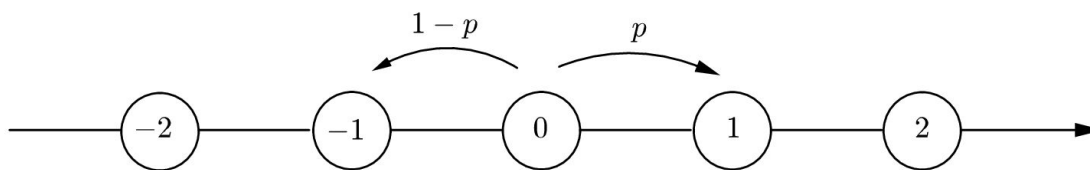
- Folosind datele din tabel, determinați probabilitățile claselor și probabilitățile condiționate ale atributelor, știind clasa.
- Considerăm evenimentul dat de *vectorul de atribut*:  $E = (C = s) \cap (T = m) \cap (M = b)$ . Alegeți o clasă pentru  $E$ , stabilind care din următoarele probabilități este mai mare:  $P(\mathbf{R} = \mathbf{r}|E)$  sau  $P(\mathbf{R} = \mathbf{n}|E)$ .
- Determinați  $P(E)$ .

4. Un mesaj este transmis printr-un canal de comunicare cu perturbări. Probabilitatea ca mesajul să fie recepționat este 10%. Dacă mesajul nu este recepționat, atunci se reia transmisia mesajului, independent de transmisiile anterioare. Fie  $X$  variabila aleatoare care indică numărul de transmisi până la prima transmisie în care este recepționat mesajul. Determinați valoarea medie a lui  $X$ .

5. 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:

Fie  $X$  variabila aleatoare care indică poziția finală a punctului material după  $n \in \mathbb{N}$  pași ai unei deplasări ce pornește din nodul 0. Determinați distribuția și valoarea medie lui  $X$ .

6. Considerăm vectorul aleator discret  $(X, Y)$  cu distribuția dată sub formă tabelară:



$\begin{array}{c} Y \\ \backslash X \end{array}$	-2	1	2
1	0,2	0,1	0,2
2	0,1	0,1	0,3

- Să se determine distribuțiile de probabilitate ale variabilelor aleatoare  $X$  și  $Y$ .
- Calculați probabilitatea ca  $|X - Y| = 1$ , știind că  $Y > 0$ .
- Sunt evenimentele  $X = 2$  și  $Y = 1$  independente?
- Sunt variabilele aleatoare  $X$  și  $Y$  independente?
- Sunt evenimentele  $X = 1$  și  $Y = 1$  condițional independente, cunoscând  $X + Y = 2$ ?
- Este variabila aleatoare  $X$  condițional independentă de  $Y$ , cunoscând  $X + Y$ ?
- Calculați valoarea medie a variabilei aleatoare  $2X + Y^2$ .

**7.** O monedă este aruncată de 10 ori. Fie  $X$  variabila aleatoare care indică diferența dintre numărul de capete și numărul de pajuri obținute. Determinați:

- distribuția de probabilitate a lui  $X$ ;
- valoarea medie a lui  $X$ .