

# Teoria Probabilităților

Teoria probabilităților este o disciplină a matematicii care se ocupă de studiul fenomenelor aleatoare.

- *aleator* = care depinde de o împrejurare viitoare și nesigură; supus întâmplării
- provine din latină: *aleatorius*; *alea* (lat.) = zar; joc cu zaruri; joc de noroc; șansă; risc

↪ se măsoară *șansele pentru succes* sau *riscul pentru insucces* al unor evenimente

Fenomene și procese aleatoare apar, de exemplu, în:

- jocuri de noroc, pariuri, loto (6 din 49)
- previziuni meteo
- previziuni economice / financiare (evaluarea stocurilor)
- sondaje de opinie, asigurări (evaluarea riscurilor, pierderilor)
- **în informatică**: sisteme de comunicare, prelucrarea informației, modelarea traficului în rețea, analiza probabilistică a unor algoritmi, fiabilitatea sistemelor, algoritmi de simulare, machine learning, data mining, recunoașterea formelor / a vocii, generarea de numere aleatoare, algoritmi aleatori: de tip Monte-Carlo, de tip Las Vegas etc.

Exemple de întrebări:

- Cum este concepută memoria cache pentru a maximiza viteza RAM a calculatoarelor?
- Rețelele de calculatoare: Care este probabilitatea ca un pachet de date să fie recepționat corect, atunci când canalul de transmisie este perturbat? În medie câte date sunt transmise corect?
- Structuri de date: Se caută un tabel hash eficient cu cea mai rapidă căutare (lookup)

## Algoritmi aleatori

**Def. 1.** *Un algoritm pe cursul executării căruia se iau anumite decizii aleatoare este numit **algoritm aleator (randomizat)**.*

- ▷ durata de execuție, spațiul de stocare, rezultatul obținut sunt variabile aleatoare (chiar dacă se folosesc aceleași valori input)
- ▷ la anumite tipuri de algoritmi corectitudinea e garantată doar cu o anumită probabilitate
- ▷ în mod paradoxal, incertitudinea ne poate oferi mai multă eficiență

Exemplu: Random QuickSort, în care elementul pivot este selectat aleator

- Algoritm de tip **Las Vegas** este un algoritm aleator, care returnează la fiecare execuție rezultatul corect (independent de alegerile aleatoare făcute); durata de execuție este o variabilă aleatoare.

Exemplu: Random QuickSort

- Un algoritm aleatoriu pentru care rezultatele obținute sunt corecte *doar* cu o anumită probabilitate se numește algoritm **Monte Carlo**.

↪ se examinează probabilitatea cu care rezultatul este corect; probabilitatea de eroare poate fi scăzută semnificativ prin execuții repetate, independente;

Exemplu:

1) testul Miller-Rabin, care verifică dacă un număr natural este prim sau este număr compus; rezultatul testului returnează fie răspunsul: “numărul este sigur un număr compus” sau răspunsul: “numărul este probabil un număr prim” (acest test nu returnează valorile divizorilor numărului compus);

2) problema tăieturii minime într-un graf (algoritmul lui D. Karger: random min-cut)

→ De care tip este următorul algoritm (scris în pseudocod)?

Input: Fie  $A(1), \dots, A(200)$  un vector cu 200 de elemente, din care 100 sunt egale cu 0 și restul egale cu 1 (ordinea lor este necunoscută).

Output: Să se găsească un 0 în vector.

```
algorithm(array A, size n)
begin
repeat
randomly select one element from A
until 0 is found
end
```

Răspuns: Algoritm de tip Las Vegas.

Versiunea Monte Carlo a problemei formulate anterior: se dă  $k$  numărul maxim de iterații

```
find_MC(array A, size n, k)
begin
  i=0
  repeat
    randomly select one element from A
    i = i + 1
  until i=k or 0 is found
end
```

▷ dacă 0 este găsit, atunci algoritmul se încheie cu rezultatul corect, altfel algoritmul nu găsește niciun 0; probabilitatea de a găsi pe 0 după  $k$  iterații este

$$P(\text{“0 este găsit după } k \text{ iterații”}) = 1 - (1/2)^k; \quad 1 - (1/2)^k \rightarrow 1, \text{ când } k \rightarrow \infty.$$

### Noțiuni introductive:

- **Experiența aleatoare** este acea experiență al cărei rezultat nu poate fi cunoscut decât după încheierea ei.
- **Evenimentul** este rezultatul unui experiment.

**Exemple:**

- ▷ Experiment: aruncarea a două zaruri, eveniment: ambele zaruri indică 1
- ▷ experiment: aruncarea unei monede, eveniment: moneda indică pajură
- ▷ experiment: extragerea unei cărți de joc, eveniment: s-a extras as
- ▷ experiment: extragerea unui număr la loto, eveniment: s-a extras numărul 27
- **evenimentul imposibil**, notat cu  $\emptyset$ , este evenimentul care nu se realizează niciodată la efectuarea experienței aleatoare
- **evenimentul sigur** este un eveniment care se realizează cu certitudine la fiecare efectuare a experienței aleatoare
- **spațiul de selecție**, notat cu  $\Omega$ , este mulțimea tuturor rezultatelor posibile ale experimentului considerat

◇ spațiul de selecție poate fi finit sau infinit

- dacă  $A$  este o submulțime a lui  $\Omega$  atunci  $A$  se numește **eveniment aleator**, iar dacă  $A$  are un singur element atunci  $A$  este un **eveniment elementar**.

▷ O analogie între evenimente și mulțimi permite o scriere și în general o exprimare mai comodă ale unor idei și rezultate legate de conceptul de eveniment aleator.

**Exemplu:** Experimentul: aruncarea unui zar, spațiul de selecție:  $\Omega = \{e_1, e_2, e_3, e_4, e_5, e_6\}$ ,  $e_i$ : s-a obținut numărul  $i$  ( $i = 1, \dots, 6$ );  $e_1, e_2, e_3, e_4, e_5, e_6$  sunt evenimente elementare

$A$ : s-a obținut un număr par  $\Rightarrow A = \{e_2, e_4, e_6\}$

$\bar{A}$ : s-a obținut un număr impar  $\Rightarrow \bar{A} = \{e_1, e_3, e_5\}$



### Operații cu evenimente

- dacă  $A, B \subseteq \Omega$ , atunci **evenimentul reuniune**  $A \cup B$  este un eveniment care se produce dacă cel puțin unul din evenimentele  $A$  sau  $B$  se produce
- dacă  $A, B \subseteq \Omega$ , atunci **evenimentul intersecție**  $A \cap B$  este un eveniment care se produce dacă cele două evenimente  $A$  și  $B$  se produc în același timp
- dacă  $A \subseteq \Omega$  atunci **evenimentul contrar** sau **complementar**  $\bar{A}$  este un eveniment care se realizează atunci când evenimentul  $A$  nu se realizează
- $A, B \subseteq \Omega$  sunt **evenimente incompatibile (disjuncte)**, dacă  $A \cap B = \emptyset$
- dacă  $A, B \subseteq \Omega$ , atunci **evenimentul diferență**  $A \setminus B$  este un eveniment care se produce dacă  $A$  are loc și  $B$  nu are loc, adică

$$A \setminus B = A \cap \bar{B}$$

### Relații între evenimente

- dacă  $A, B \subseteq \Omega$ , atunci  $A$  **implică**  $B$ , dacă producerea evenimentului  $A$  conduce la producerea evenimentului  $B$ :  $A \subseteq B$
- dacă  $A$  implică  $B$  și  $B$  implică  $A$ , atunci evenimentele  $A$  și  $B$  sunt **egale**:  $A = B$

### Proprietăți ale operațiilor între evenimente $A, B, C \subseteq \Omega$

Operațiile de reuniune și intersecție sunt operații **comutative**:

$$A \cup B = B \cup A, \quad A \cap B = B \cap A,$$

**asociative**

$$(A \cup B) \cup C = A \cup (B \cup C), \quad (A \cap B) \cap C = A \cap (B \cap C),$$

**și distributive**

$$(A \cup B) \cap C = (A \cap C) \cup (B \cap C), \quad (A \cap B) \cup C = (A \cup C) \cap (B \cup C);$$

satisfac legile lui De Morgan

$$\overline{A \cup B} = \bar{A} \cap \bar{B}, \quad \overline{A \cap B} = \bar{A} \cup \bar{B}.$$

Are loc  $\bar{\bar{A}} = A$ .

**Frecvența relativă și frecvența absolută**

**Def. 2.** Fie  $A$  un eveniment asociat unei experiențe, repetăm experiența de  $n$  ori (în aceleași condiții date) și notăm cu  $r_n(A)$  numărul de realizări ale evenimentului  $A$ ; **frecvența relativă** a evenimentului  $A$  este numărul

$$f_n(A) = \frac{r_n(A)}{n}$$

$r_n(A)$  este **frecvența absolută** a evenimentului  $A$ .

**Definiția clasică a probabilității**

**Def. 3.** Într-un experiment în care cazurile posibile sunt finite la număr și au aceleași șanse de a se realiza, **probabilitatea** unui eveniment  $A$  este numărul

$$P(A) = \frac{\text{numărul de cazuri favorabile apariției lui } A}{\text{numărul total de cazuri posibile}}.$$

▷ Prin repetarea de multe ori a unui experiment, în condiții practic identice, frecvența relativă  $f_n(A)$  de apariție a evenimentului  $A$  este aproximativ egală cu  $P(A)$

$$f_n(A) \approx P(A), \text{ dacă } n \rightarrow \infty.$$

**Exemplu:** Experiment: Se aruncă 4 monede. Evenimentul  $A$ : cele 4 monede indică pajură exact de 3 ori ; experimentul s-a repetat de  $n = 100$  de ori și evenimentul  $A$  a apărut de 22 de ori.

$$f_n(A) = ?, \quad P(A) = ?$$

$$\text{Răspuns: } f_n(A) = \frac{22}{100} = 0.22$$

$$\Omega = \{(c, c, c, c), (c, p, p, p), \dots, (p, p, p, c), (p, p, p, p)\}$$

$$A = \{(c, p, p, p), (p, c, p, p), (p, p, c, p), (p, p, p, c)\} \Rightarrow P(A) = \frac{4}{2^4} = 0.25$$



### Definiția axiomatică a probabilității

Definiția clasică a probabilității poate fi utilizată numai în cazul în care numărul cazurilor posibile este finit. Dacă numărul evenimentelor elementare este infinit, atunci există evenimente pentru care probabilitatea în sensul clasic nu are nici un înțeles.

**Probabilitatea geometrică:** Măsura unei mulțimi corespunde “lungimii” în  $\mathbb{R}$ , “ariei” în  $\mathbb{R}^2$ , “volumului” în  $\mathbb{R}^3$ . Fie  $M \subset D \subset \mathbb{R}^n$ ,  $n \in \{1, 2, 3\}$ , mulțimi cu măsură finită.

Alegem aleator un punct  $A \in D$  (în acest caz spațiul de selecție este  $D$ ). Probabilitatea geometrică a evenimentului “ $A \in M$ ” este

$$P(A \in M) := \frac{\text{măsura}(M)}{\text{măsura}(D)}.$$

O teorie formală a probabilității a fost creată în anii '30 ai secolului XX de către matematicianul rus **Andrei Nikolaevici Kolmogorov**, care, în anul **1933**, a dezvoltat teoria axiomatică a probabilității în lucrarea sa *Conceptele de bază ale Calculului Probabilității*.

$\Rightarrow P : \mathcal{K} \rightarrow \mathbb{R}$  este o funcție astfel încât oricărui eveniment aleator  $A \in \mathcal{K}$  i se asociază valoarea  $P(A)$ , **probabilitatea de apariție a evenimentului  $A$**

$\hookrightarrow \mathcal{K}$  este o mulțime de evenimente și are structura unei  $\sigma$ -algebre (vezi Def. 4)

$\hookrightarrow P$  satisface anumite axiome (vezi Def. 5)

**Def. 4.** O familie  $\mathcal{K}$  de evenimente din spațiul de selecție  $\Omega$  se numește  **$\sigma$ -algebră** dacă sunt satisfăcute condițiile:

- (i)  $\mathcal{K}$  este nevidă;
- (ii) dacă  $A \in \mathcal{K}$ , atunci  $\bar{A} \in \mathcal{K}$ ;
- (iii) dacă  $A_n \in \mathcal{K}$ ,  $n \in \mathbb{N}^*$ , atunci  $\bigcup_{n=1}^{\infty} A_n \in \mathcal{K}$ .

Perechea  $(\Omega, \mathcal{K})$  se numește **spațiu măsurabil**.

**Exemple:** 1) Dacă  $\emptyset \neq A \subset \Omega$  atunci  $\mathcal{K} = \{\emptyset, A, \bar{A}, \Omega\}$  este o  $\sigma$ -algebră.

2)  $\mathcal{P}(\Omega) :=$  mulțimea tuturor submulțimilor ale lui  $\Omega$  este o  $\sigma$ -algebră.

3) Dacă  $(\Omega, \mathcal{K})$  este un spațiu măsurabil și  $\emptyset \neq B \subseteq \Omega$ , atunci

$$B \cap \mathcal{K} = \{B \cap A : A \in \mathcal{K}\}$$

este o  $\sigma$ -algebră pe mulțimea  $B$ , iar  $(B, B \cap \mathcal{K})$  este un spațiu măsurabil.

**P. 1.** *Proprietăți ale unei  $\sigma$ -algebre:* Dacă  $\mathcal{K}$  este o  $\sigma$ -algebră în  $\Omega$ , atunci au loc proprietățile:

- (1)  $\emptyset, \Omega \in \mathcal{K}$ ;
- (2)  $A, B \in \mathcal{K} \implies A \cap B, A \setminus B \in \mathcal{K}$ ;

$$(3) A_n \in \mathcal{K}, n \in \mathbb{N}^* \implies \bigcap_{n=1}^{\infty} A_n \in \mathcal{K}.$$

**Def. 5.** Fie  $\mathcal{K}$  o  $\sigma$ -algebră în  $\Omega$ . O funcție  $P : \mathcal{K} \rightarrow \mathbb{R}$  se numește **probabilitate** dacă satisface axiomele:

- (i)  $P(\Omega) = 1$ ;
- (ii)  $P(A) \geq 0$  pentru orice  $A \in \mathcal{K}$ ;
- (iii) pentru orice șir  $(A_n)_{n \in \mathbb{N}^*}$  de evenimente două câte două disjuncte (adică  $A_i \cap A_j = \emptyset$  pentru orice  $i \neq j$ ) din  $\mathcal{K}$  are loc

$$P\left(\bigcup_{n=1}^{\infty} A_n\right) = \sum_{n=1}^{\infty} P(A_n).$$

Tripletul  $(\Omega, \mathcal{K}, P)$  format din spațiul măsurabil  $(\Omega, \mathcal{K})$  și probabilitatea  $P : \mathcal{K} \rightarrow \mathbb{R}$  se numește **spațiu de probabilitate**.

**P. 2.** Fie  $(\Omega, \mathcal{K}, P)$  un spațiu de probabilitate. Au loc proprietățile:

- (1)  $P(\bar{A}) = 1 - P(A)$  și  $0 \leq P(A) \leq 1$ ;
- (2)  $P(\emptyset) = 0$ ;
- (3)  $P(A \setminus B) = P(A) - P(A \cap B)$ ;
- (4)  $A \subseteq B \implies P(A) \leq P(B)$ , adică  $P$  este monotonă;
- (5)  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$ .

**Exercițiu:** a) Să se arate că pentru  $\forall A, B, C \in \mathcal{K}$  are loc:

$$P(A \cup B \cup C) = P(A) + P(B) + P(C) - P(A \cap B) - P(A \cap C) - P(B \cap C) + P(A \cap B \cap C).$$

b) Pentru  $A_1, \dots, A_n \in \mathcal{K}$  care e formula similară de calcul pentru  $P(A_1 \cup A_2 \cup \dots \cup A_n)$ ?

**Exemplu:** Dintr-un pachet de 52 de cărți de joc se extrage o carte aleator. Care este probabilitatea  $p$  de a extrage a) un as sau o damă de pică? b) o inimă sau un as?

R.: a)  $A$ : s-a extras un as;  $D$ : s-a extras damă de pică;  $A$  și  $D$  sunt două evenimente incompatibile (disjuncte)

$$p = P(A \cup D) = P(A) + P(D) = \frac{4 + 1}{52};$$

b)  $I$ : s-a extras inimă;  $I$  și  $A$  nu sunt evenimente incompatibile

$$p = P(I \cup A) = P(I) + P(A) - P(I \cap A) = \frac{13 + 4 - 1}{52} = \frac{4}{13}.$$



## Probabilitate condiționată

**Def. 6.** Fie  $(\Omega, \mathcal{K}, P)$  un spațiu de probabilitate și fie  $A, B \in \mathcal{K}$ . **Probabilitatea condiționată a evenimentului  $A$  de evenimentul  $B$**  este  $P(\cdot|B) : \mathcal{K} \rightarrow [0, 1]$  definită prin

$$P(A|B) = \frac{P(A \cap B)}{P(B)},$$

dacă  $P(B) > 0$ .  $P(A|B)$  este probabilitatea apariției evenimentului  $A$ , știind că evenimentul  $B$  s-a produs.

**Exemplu:** Se extrag succesiv fără reținere două bile dintr-o urnă cu 4 bile albe și 5 bile roșii.

a) Știind că prima bilă este roșie, care este probabilitatea ca a doua bilă să fie albă?

b) Care este probabilitatea ca ambele bile să fie roșii?

R.: pentru  $i \in \{1, 2\}$  fie evenimentele

$R_i$ : la a  $i$ -a extragere s-a obținut o bilă roșie;

$A_i = \bar{R}_i$ : la a  $i$ -a extragere s-a obținut o bilă albă;

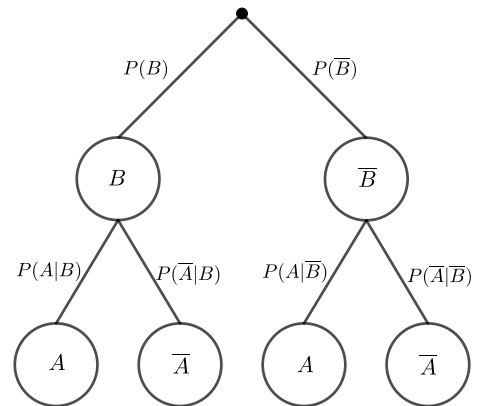
a)  $P(A_2|R_1) = \frac{4}{8}$ . b)  $P(R_1 \cap R_2) = P(R_2|R_1)P(R_1) = \frac{4}{8} \cdot \frac{5}{9}$ .



**P. 3.** Pentru  $A, B \in \mathcal{K}$ ,  $P(A) > 0$ ,  $P(B) > 0$  au loc:

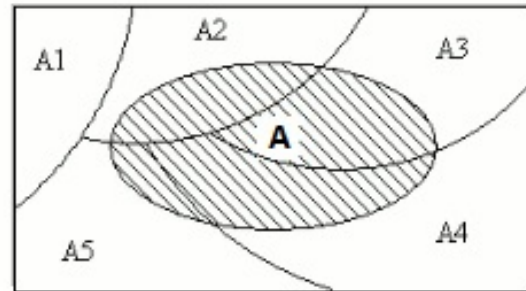
$$P(A \cap B) = P(B)P(A|B) = P(A)P(B|A),$$

$$P(\bar{A}|B) = 1 - P(A|B).$$



**Fig.1. Probabilități condiționate**

**Def. 7.** O familie  $A_1, \dots, A_n$  de evenimente din  $\Omega$  se numește **partiție** sau **sistem complet de evenimente** a lui  $\Omega$ , dacă  $\bigcup_{i=1}^n A_i = \Omega$  și pentru fiecare  $i, j \in \{1, \dots, n\}$ ,  $i \neq j$ , evenimentele  $A_i$  și  $A_j$  sunt disjuncte, adică  $A_i \cap A_j = \emptyset$ .



Partiție  $A_1 \cup A_2 \cup A_3 \cup A_4 \cup A_5 = \Omega$

**Exemplu:** Dacă  $B \subset \Omega$  atunci  $\{B, \bar{B}\}$  formează o partiție a lui  $\Omega$ .



**P. 4. (Formula probabilității totale)** Într-un spațiu de probabilitate  $(\Omega, \mathcal{K}, P)$  considerăm partiția  $H_1, \dots, H_n$  a lui  $\Omega$  cu  $P(H_i) > 0$  și  $H_i \in \mathcal{K} \forall i \in \{1, \dots, n\}$ , și fie  $A \in \mathcal{K}$ . Atunci are loc

$$P(A) = \sum_{i=1}^n P(A|H_i)P(H_i).$$

**Exemplu:** Într-o urnă sunt 7 bile albe, notate cu 1,2,3,4,5,6,7, și 6 bile roșii notate cu 8,9,10,11,12,13. Se extrage o bilă. **a)** Știind că bila extrasă este roșie, care este probabilitatea  $p_1$ , ca numărul înscris să fie divizibil cu 4? **b)** Știind că prima bilă este roșie, care este probabilitatea  $p_2$ , ca o a doua bilă extrasă să fie un număr impar? (Prima bilă nu s-a returnat în urnă!)

R.: Se consideră evenimentele:

$A_1$ : prima bilă extrasă are înscris un număr divizibil cu 4;

$B_1$ : prima bilă extrasă este roșie;

$C_1$ : prima bilă extrasă are înscris un număr impar;

$C_2$ : a doua bilă extrasă are înscris un număr impar.

**a)**  $p_1 = P(A_1|B_1) = \frac{2}{6}$ .

**b)**  $p_2 = P(C_2|B_1) = ?$  Folosim Def.6 și P.3, scriem succesiv

$$\begin{aligned} p_2 = P(C_2|B_1) &= \frac{P(C_2 \cap B_1)}{P(B_1)} = \frac{P(C_2 \cap B_1 \cap C_1) + P(C_2 \cap B_1 \cap \bar{C}_1)}{P(B_1)} \\ &= \frac{P(C_2|B_1 \cap C_1)P(B_1 \cap C_1) + P(C_2|B_1 \cap \bar{C}_1)P(B_1 \cap \bar{C}_1)}{P(B_1)} = \frac{\frac{6}{12} \cdot \frac{3}{13} + \frac{7}{12} \cdot \frac{3}{13}}{\frac{6}{13}} = \frac{13}{24}. \end{aligned}$$



## Formula lui Bayes

Formula lui Bayes este o metodă de a "corecta" (a revizui, a îmbunătăți) pe baza unor noi informații (date) disponibile o probabilitate determinată apriori. Se pornește cu o estimare pentru probabilitatea unei anumite ipoteze  $I$ . Dacă avem noi informații (date)  $D$ , ce privesc ipoteza  $I$ , se poate calcula o probabilitate "corectată" pentru ipoteza  $I$ , numită probabilitate posterioară (a-posteriori).

$\hookrightarrow P(I)$  probabilitatea ca ipoteza  $I$  să fie adevărată, numită și *probabilitatea apriori*;

$\hookrightarrow$  probabilitatea condiționată  $P(I|D)$  este *probabilitatea posterioară* (corectată de noile informații);

$\hookrightarrow P(D|I)$  probabilitatea ca să apară informațiile, dacă ipoteza  $I$  este adevărată;

$\hookrightarrow P(D|\bar{I})$  probabilitatea ca să apară informațiile, dacă ipoteza  $I$  este falsă.

Folosind P.5 are loc:

$$P(D) = P(D|I) \cdot P(I) + P(D|\bar{I}) \cdot P(\bar{I}) = P(D|I) \cdot P(I) + P(D|\bar{I}) \cdot (1 - P(I)).$$

Formula lui Bayes este în acest caz

$$P(I|D) = \frac{P(D|I) \cdot P(I)}{P(D)} = \frac{P(D|I) \cdot P(I)}{P(D|I) \cdot P(I) + P(D|\bar{I}) \cdot P(\bar{I})}.$$



**Exemplu:** Considerăm evenimentele (în teste clinice):

$I$ : o persoană aleasă aleator dintr-o populație are o anumită alergie  $\mathcal{A}$

$D_+$ : testul clinic returnează pozitiv privind alergia  $\mathcal{A}$

$\bar{D}_+$ : testul clinic returnează negativ privind alergia  $\mathcal{A}$

▷ din statistici anterioare sunt cunoscute:

$p = P(I)$ , probabilitatea ca o persoană selectată aleator din populație să sufere de alergia  $\mathcal{A}$ ;

sensibilitatea testului  $s_1 = P(D_+|I)$ ;

specificitatea testului  $s_2 = P(\bar{D}_+|\bar{I})$ ;

⇒ probabilitatea de a obține răspuns fals pozitiv este  $P(D_+|\bar{I}) = 1 - s_2$ ;

▷ un test clinic bun implică valori apropiate de 1 pentru  $s_1$  și  $s_2$ ;

► cunoscând  $p, s_1, s_2$  se dorește a se determina valoarea predictivă  $P(I|D_+)$

$$P(I|D_+) = \frac{P(D_+|I) \cdot P(I)}{P(D_+)} = \frac{P(D_+|I) \cdot P(I)}{P(D_+|I) \cdot P(I) + P(D_+|\bar{I}) \cdot P(\bar{I})} = \frac{s_1 \cdot p}{s_1 \cdot p + (1 - s_2) \cdot (1 - p)}.$$



**P. 5. (Formula lui Bayes)** Într-un spațiu de probabilitate  $(\Omega, \mathcal{K}, P)$  considerăm partiția  $H_1, \dots, H_n$  a lui  $\Omega$  cu  $P(H_i) > 0$  și  $H_i \in \mathcal{K} \forall i \in \{1, \dots, n\}$ , și fie  $E \in \mathcal{K}$  astfel încât  $P(E) > 0$ . Atunci,

$$P(H_j|E) = \frac{P(H_j)P(E|H_j)}{P(E)} = \frac{P(H_j)P(E|H_j)}{\sum_{i=1}^n P(H_i)P(E|H_i)} \quad \forall j \in \{1, 2, \dots, n\}.$$

▷ pentru  $i \in \{1, 2, \dots, n\}$   $P(H_i)$  sunt **probabilități apriori** pentru  $H_i$ , numite și ipoteze (asertiuni),  $E$  se numește **evidență** (dovadă, premisă, informație); cu formula lui Bayes se calculează probabilitățile pentru ipoteze, cunoscând evidența:  $P(H_i|E)$ , acestea se numesc **probabilități posterioare** (ulterioare);  $P(E|H_i)$  reprezintă verosimilitatea datelor observate.

▷ Se pot calcula probabilitățile *cauzelor*, date fiind *efectele*; formula lui Bayes ne ajută să diagnosticăm o anumită situație sau să testăm o ipoteză.

**Exemplu:** Ce probabilități calculează programul de mai jos? Ce tip de algoritm aleator este?

► `randi(imax, n, m)` generează o  $n \times m$  matrice cu valori întregi aleatoare (pseudoaleatoare) între 1 și imax.

```
clear all
ci=0;
cp=0;
c=0;
a=0;
b=0;
N=1000;
for i=1:N
```

```

A=[randi(5,1,5),5+randi(8,1,5),13+randi(7,1,10)];
r= randi(length(A));
v=A(r);
ci=ci+mod(v,2);
cp=cp+(mod(v,2)==0);
c=c+ mod(v,2)*(mod(v,3)==0);
a=a+ mod(v,2)*(6<=r && r<=10);
b=b+ (mod(v,2)==0)*(r>=10);
end
p1=c/ci
p2=a/ci
p3=b/cp

```

R.: Se dă un șir  $A$  format din 20 de elemente, în care

- 25% provin din generarea aleatoare și cu aceeași probabilitate (care e  $1/5$ ) a unui număr din  $\{1, 2, 3, 4, 5\}$
- 25% provin din generarea aleatoare și cu aceeași probabilitate (care e  $1/8$ ) a unui număr din  $\{6, 7, 8, 9, 10, 11, 12, 13\}$
- 50% provin din generarea aleatoare și cu aceeași probabilitate (care e  $1/7$ ) a unui număr din  $\{14, 15, 16, 17, 18, 19, 20\}$ .

Se extrage aleator un număr din șir.

►  $p1$  estimează probabilitatea condiționată ca numărul ales aleator să fie divizibil cu 3, *știind* că s-a extras un număr impar;

►  $p2$  estimează probabilitatea condiționată ca numărul ales aleator să provină din mulțimea  $\{6, 7, 8, 9, 10, 11, 12, 13\}$ , *știind* că s-a extras un număr impar;

►  $p3$  estimează probabilitatea condiționată ca numărul ales aleator să provină din mulțimea  $\{14, 15, 16, 17, 18, 19, 20\}$ , *știind* că s-a extras un număr par.

Algoritmul este de tip Monte-Carlo!



**Exercițiu:** Să se calculeze valorile teoretice pentru probabilitățile  $p1$ ,  $p2$ ,  $p3$  din exemplul anterior!



**P. 6. (Regula de înmulțire)** Fie  $(\Omega, \mathcal{K}, P)$  un spațiu de probabilitate și fie  $A_1, \dots, A_n \in \mathcal{K}$  astfel încât  $P(A_1 \cap \dots \cap A_{n-1}) > 0$ . Atunci,

$$P(A_1 \cap \dots \cap A_n) = P(A_1)P(A_2|A_1) \dots P(A_n|A_1 \cap \dots \cap A_{n-1}).$$

**Exemplu:** Într-o urnă sunt 2 bile verzi și 3 bile albastre. Se extrag 2 bile succesiv, fără returnare. Care este probabilitatea ca

- prima bilă să fie verde, iar cea de-a doua albastră?
- cele 2 bile să aibă aceeași culoare?
- a doua bilă să fie albastră?

d) prima bilă să fie verde, *știind* că a doua este albastră?

e) se mai extrage o a treia bilă; se cere probabilitatea ca prima bilă să fie verde, cea de-a doua albastră și a treia tot albastră.

R.: Notăm pentru  $i \in \{1, 2, 3\}$  evenimentele:

$A_i$ : la a  $i$ -a extragere s-a obținut bilă albastră;  $V_i$ : la a  $i$ -a extragere s-a obținut bilă verde;

a) folosim P.3:  $P(V_1 \cap A_2) = P(A_2|V_1)P(V_1) = \frac{3}{4} \cdot \frac{2}{5}$

b)  $P((V_1 \cap V_2) \cup (A_1 \cap A_2)) = P(V_1 \cap V_2) + P(A_1 \cap A_2) = P(V_2|V_1)P(V_1) + P(A_2|A_1)P(A_1) = \frac{1}{4} \cdot \frac{2}{5} + \frac{2}{4} \cdot \frac{3}{5}$

c) folosim formula probabilității totale P.5:

$P(A_2) = P(A_2|V_1)P(V_1) + P(A_2|A_1)P(A_1) = \frac{3}{4} \cdot \frac{2}{5} + \frac{2}{4} \cdot \frac{3}{5}$

d) folosim P.3:  $P(V_1|A_2) = \frac{P(V_1 \cap A_2)}{P(A_2)} = \frac{P(A_2|V_1)P(V_1)}{P(A_2)} = \frac{\frac{3}{4} \cdot \frac{2}{5}}{\frac{3}{4} \cdot \frac{2}{5} + \frac{2}{4} \cdot \frac{3}{5}}$

e) formula de înmulțire a probabilităților P.6:

$P(V_1 \cap A_2 \cap A_3) = P(V_1) \cdot P(A_2|V_1) \cdot P(A_3|V_1 \cap A_2) = \frac{2}{5} \cdot \frac{3}{4} \cdot \frac{2}{3}$ .

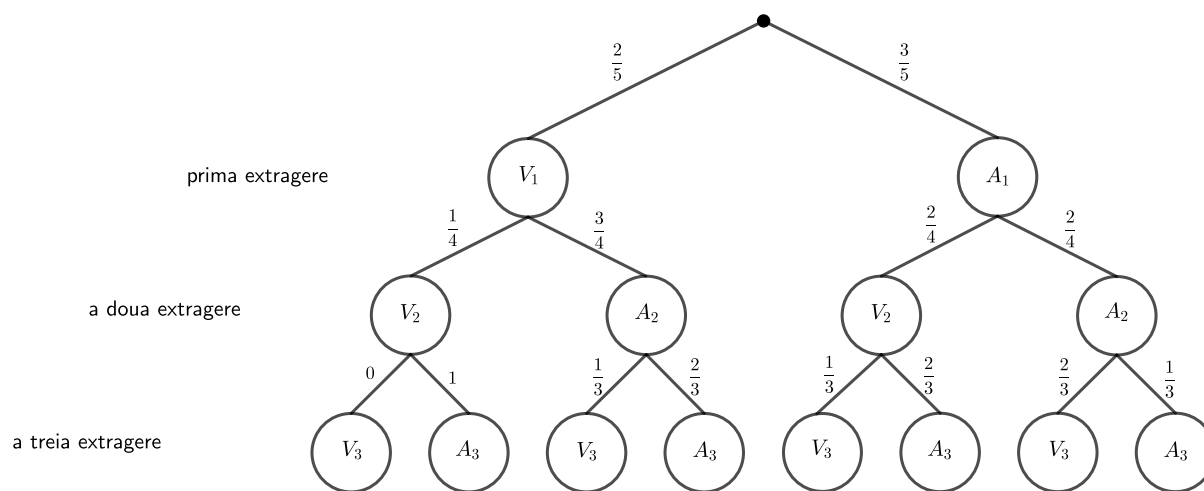


Fig. 3. Extragere fără returnare



## Evenimente independente

**Def. 8.** Fie  $(\Omega, \mathcal{K}, P)$  un spațiu de probabilitate. Evenimentele  $A, B \in \mathcal{K}$  sunt **evenimente independente** dacă

$$P(A \cap B) = P(A)P(B).$$

**Observație:** Fie evenimentele  $A, B \in \mathcal{K}$  astfel încât  $P(A) > 0$  și  $P(B) > 0$ . Evenimentele  $A$  și  $B$  sunt **independente**, dacă apariția evenimentului  $A$ , nu influențează apariția evenimentului  $B$  și invers, adică

$$P(A|B) = P(A) \text{ și } P(B|A) = P(B).$$

Două evenimente se numesc **dependente** dacă probabilitatea realizării unuia dintre ele depinde de faptul că celălalt eveniment s-a produs sau nu.

**Exercițiu: 1)** Se extrag succesiv fără returnare două cărți de joc dintr-un pachet de cărți;

A: prima carte extrasă este as;

B: prima carte extrasă este damă

C: a doua carte extrasă este 1;

Sunt A și B evenimente independente?

Sunt A și C evenimente independente?

Sunt B și C evenimente dependente?

**2)** Se aruncă un zar de două ori.

A: primul număr este 6;

B: al doilea număr este 5

C: primul număr este 1;

Sunt A și B evenimente independente?

Sunt A și C evenimente independente?

Sunt B și C evenimente dependente?



**P. 7.** Fie  $(\Omega, \mathcal{K}, P)$  un spațiu de probabilitate și fie  $A, B \in \mathcal{K}$ . Sunt echivalente afirmațiile:

(1)  $A$  și  $B$  sunt independente.

(2)  $\bar{A}$  și  $B$  sunt independente.

(3)  $A$  și  $\bar{B}$  sunt independente.

(4)  $\bar{A}$  și  $\bar{B}$  sunt independente.

**Def. 9.** Fie  $(\Omega, \mathcal{K}, P)$  un spațiu de probabilitate.  $B_1, \dots, B_n$  sunt  $n$  **evenimente independente** (în totalitate) din  $\mathcal{K}$  dacă

$$P(B_{i_1} \cap \dots \cap B_{i_m}) = P(B_{i_1}) \cdot \dots \cdot P(B_{i_m})$$

pentru orice submulțime finită  $\{i_1, \dots, i_m\} \subseteq \{1, 2, \dots, n\}$ .

**Exemplu: 1)**  $A, B, C \in \mathcal{K}$  sunt trei evenimente independente (în totalitate), dacă

$$P(A \cap B) = P(A)P(B), \quad P(A \cap C) = P(A)P(C), \quad P(B \cap C) = P(B)P(C),$$

$$P(A \cap B \cap C) = P(A)P(B)P(C).$$

**2)** Cele 4 fețe ale unui tetraedru regulat sunt vopsite astfel: una este roșie, una este albastră, una este verde și una este colorată având cele trei culori. Se aruncă tetraedrul și se consideră evenimentele:  $R$ : tetraedrul cade pe partea roșie;  $A$ : tetraedrul cade pe partea albastră;  $V$ : tetraedrul cade pe partea verde. Se poate afirma că  $P(R \cap A \cap V) = P(R)P(A)P(V)$ ? Sunt

cele 3 evenimente independente în totalitate?

3) Pentru a verifica dacă  $n$  evenimente distincte  $B_1, \dots, B_n$  sunt independente (în totalitate) câte relații trebuie verificate? R.:  $C_n^2 + C_n^3 + \dots + C_n^n = 2^n - n - 1$  ♦

### Variable aleatoare

**Exemplu:** Un jucător aruncă două monede  $\Rightarrow \Omega = \{(c, p), (c, c), (p, c), (p, p)\}$  ( $c$ =cap;  $p$ =pajură)

$X$  indică de câte ori a apărut pajură:  $\Rightarrow X : \Omega \rightarrow \{0, 1, 2\}$

$\Rightarrow P(X = 0) = P(X = 2) = \frac{1}{4}, P(X = 1) = \frac{1}{2}$  ■

**Notăție 1.** *variabilă/variabile aleatoare  $\rightarrow$  v.a.*

O variabilă aleatoare este:

► **discretă**, dacă ia un număr finit de valori  $(x_1, \dots, x_n)$  sau un număr infinit numărabil de valori  $(x_1, \dots, x_n, \dots)$

► **continuă**, dacă valorile sale posibile sunt nenumărabile și sunt într-un interval (sau reunine de intervale) sau în  $\mathbb{R}$

**V.a. discrete:** exemple de v.a. numerice discrete: suma numerelor obținute la aruncarea a 4 zaruri, numărul produselor defecte produse de o anumită firmă într-o săptămână; numărul apelurilor telefonice într-un call center în decursul unei ore; numărul de accesări ale unei anumite pagini web în decursul unei anumite zile (de ex. duminică); numărul de caractere transmise eronat într-un mesaj de o anumită lungime; exemple de v.a. categoriale ( $\rightarrow$  se clasifică în categorii): prognoza meteo: ploios, senin, înnorat, cețos; calitatea unor servicii: nesatisfăcătoare, satisfăcătoare, bune, foarte bune, excepționale ...)

**V.a. continue** sunt v.a. numerice: timpul de funcționare până la defectare a unei piese electronice, temperatura într-un oraș, viteza înregistrată de radar pentru mașini care parcurg o anumită zonă ...

### Variabile aleatoare numerice - definiție formală

**Def. 10.** Fie  $(\Omega, \mathcal{K}, P)$  spațiu de probabilitate  $X : \Omega \rightarrow \mathbb{R}$  este o variabilă aleatoare, dacă

$$\{\omega \in \Omega : X(\omega) \leq x\} \in \mathcal{K} \text{ pentru fiecare } x \in \mathbb{R}.$$

**Variabile aleatoare discrete**  $X : \Omega \rightarrow \{x_1, x_2, \dots, x_i, \dots\}$

**Def. 11.** Distribuția de probabilitate a v.a. discrete  $X$

$$X \sim \begin{pmatrix} x_1 & x_2 & \dots & x_i & \dots \\ p_1 & p_2 & \dots & p_i & \dots \end{pmatrix} = \begin{pmatrix} x_i \\ p_i \end{pmatrix}_{i \in I}$$

$I \subseteq \mathbb{N}$  (mulțime de indici nevidă);  $p_i = P(X = x_i) > 0, i \in I$ , cu  $\sum_{i \in I} p_i = 1$ .

▷ Variabilele aleatoare discrete sunt caracterizate de distribuția de probabilitate.

▷ Notăm  $\{X = x_i\} = \{\omega \in \Omega : X(\omega) = x_i\}$ ;  $\{X = x_i\}$  este un eveniment din  $\mathcal{K}$  pentru fiecare  $i \in I$ .

### Distribuții discrete clasice

**Distribuția discretă uniformă:**  $X \sim Unif(n)$

$$X \sim \begin{pmatrix} 1 & 2 & \dots & n \\ \frac{1}{n} & \frac{1}{n} & \dots & \frac{1}{n} \end{pmatrix}$$

**Exemplu:** Se aruncă un zar, fie  $X$  v.a. care indică numărul apărut

$$\Rightarrow X \sim \begin{pmatrix} 1 & 2 & \dots & 6 \\ \frac{1}{6} & \frac{1}{6} & \dots & \frac{1}{6} \end{pmatrix}$$

Matlab/Octave: `unidrnd(n,...)`, `randi(n,...)`

**Distribuția Bernoulli:**  $X \sim Bernoulli(p)$ ,  $p \in (0, 1)$

$$X \sim \begin{pmatrix} 0 & 1 \\ 1-p & p \end{pmatrix}$$

**Exemplu:** în cadrul unui experiment poate să apară evenimentul  $A$  (*succes*) sau  $\bar{A}$  (*insucces*)

$X = 0 \Leftrightarrow$  dacă  $\bar{A}$  apare;  $X = 1 \Leftrightarrow$  dacă  $A$  apare

$\Rightarrow X \sim Bernoulli(p)$  cu  $p := P(A)$

$$X \sim \begin{pmatrix} 0 & 1 \\ 1 - P(A) & P(A) \end{pmatrix}$$



generare în Matlab/Octave:

```
n=1000;
p=0.3;
nr=rand(1,n);
X=(nr<=p) % vector de date avand distributia Bernoulli(p)
%%%%%%%%%
Y=floor(rand(1,n)+p)% vector de date avand distributia Bernoulli(p)
%%%%%%%%%
Z=binornd(1,p,1,n)% vector de date avand distributia Bernoulli(p)
```

**Distribuția binomială:**  $X \sim Bino(n, p)$ ,  $n \in \mathbb{N}^*$ ,  $p \in (0, 1)$

în cadrul unui experiment poate să apară evenimentul  $A$  (*succes*) sau  $\bar{A}$  (*insucces*)

•  $A = \text{succes}$  cu  $P(A) = p$ ,  $\bar{A} = \text{insucces}$   $P(\bar{A}) = 1 - p$

- se repetă experimentul de  $n$  ori
- v.a.  $X$  = numărul de succese în  $n$  repetări independente ale experimentului  $\Rightarrow$  valori posibile:  $X \in \{0, 1, \dots, n\}$

$$P(X = k) = C_n^k p^k (1 - p)^{n-k}, \quad k = 0, \dots, n.$$

**Exemplu:** Un zar se aruncă de 10 ori, fie  $X$  v.a. care indică de câte ori a apărut numărul 6  $\Rightarrow Bino(10, \frac{1}{6})$ .

$\rightarrow$  are loc **formula binomială**

$$(a + b)^n = \sum_{k=0}^n C_n^k a^k b^{n-k}$$

pentru  $a = p$  și  $b = 1 - p$  se obține

$$1 = \sum_{k=0}^n C_n^k p^k (1 - p)^{n-k}.$$

Matlab/Octave: `binornd(n, p, ...)`

$\triangleright$  Distribuția binomială corespunde modelului cu extragerea bilelor dintr-o urnă cu returnarea bilelor după fiecare extragere.

**Exemplu:** Într-o urnă sunt  $n_1$  bile albe și  $n_2$  bile negre. Se extrag cu returnare  $n$  bile; fie v.a.  $X$  = numărul de bile albe extrase  $\Rightarrow X \sim Bino(n, p)$  cu  $p = \frac{n_1}{n_1 + n_2}$ .

2) Fie un canal de comunicare binară care transmite cuvinte codificate de  $N$  biți fiecare. Probabilitatea transmiterii cu succes a unui singur bit este  $p$ , iar probabilitatea unei erori este  $1 - p$ . Presupunem, de asemenea, că un astfel de cod este capabil să corecteze până la  $m$  erori, unde  $0 \leq m \leq N$ . Se știe că transmiterea biților succesivi este independentă, atunci probabilitatea transmiterii cu succes a cuvântului este  $p = P(A)$ , unde

$A$ : cel mult  $m$  erori apar în transmiterea celor  $N$  biți

$$p = P(A) = \sum_{k=0}^m C_N^k p^{N-k} (1 - p)^k.$$



**Exercițiu:** 1) Un client accesează o dată pe zi o anumită pagină web, care oferă produse bio, cu probabilitatea 0.4. Cu ce probabilitate clientul accesează această pagină în total de 3 ori în următoarele 6 zile?

2) O rețea de laborator este compusă din 15 calculatoare. Rețeaua a fost atacată de un virus nou, care atacă un calculator cu o probabilitatea 0.4, independent de alte calculatoare. Care este probabilitatea ca virusul a atacat a) cel mult 10 computere; b) cel puțin 10 calculatoare; c) exact 10 calculatoare?



**Distribuția hipergeometrică:**  $X \sim Hyge(n, n_1, n_2)$

Într-o urnă sunt  $n_1$  bile albe și  $n_2$  bile negre. Se extrag **fără returnare**  $n$  bile.

Fie v.a.  $X$  = numărul de bile albe extrase  $\Rightarrow$  valori posibile pentru  $X$  sunt  $\{0, 1, \dots, n^*\}$  cu

$$n^* = \min(n_1, n) = \begin{cases} n_1 & \text{dacă } n_1 < n \text{ (mai puține bile albe decât numărul de extrageri)} \\ n & \text{dacă } n_1 \geq n \text{ (mai multe bile albe decât numărul de extrageri)} \end{cases}$$

Fie  $n_1, n_2, n \in \mathbb{N}$  cu  $n \leq n_2$  și notăm  $n^* = \min(n_1, n)$ .

$$\Rightarrow P(X = k) = \frac{C_{n_1}^k C_{n_2}^{n-k}}{C_{n_1+n_2}^n}, \quad k = 0, \dots, n^*.$$

Matlab/Octave: `hygernd( $n_1 + n_2, n_1, n, \dots$ )`

**Exemplu:** Loto 6 din 49  $\rightarrow$  Care este probabilitatea de a nimeri exact 4 numere câștigătoare?  
R.: Între cele 49 de bile exact  $n_1 = 6$  sunt câștigătoare (“bilele albe”) și  $n_2 = 43$  necâștigătoare (“bilele negre”). Care este probabilitatea ca din  $n = 6$  extrageri fără returnare, exact  $k = 4$  numere să fie câștigătoare?

$$\Rightarrow P(X = 4) = \frac{C_6^4 C_{43}^2}{C_{49}^6}$$

◇

**Distribuția geometrică**  $X \sim Geo(p)$

În cadrul unui experiment poate să apară evenimentul  $A$  (*succes*) sau  $\bar{A}$  (*insucces*)

- $A$  = succes cu  $P(A) = p$ ,  $\bar{A}$  = insucces  $P(\bar{A}) = 1 - p$
- se repetă (independent) experimentul până apare prima dată  $A$  (“succes”)
- v.a.  $X$  arată de câte ori apare  $\bar{A}$  (numărul de “insuccese”) *până* la apariția primului  $A$  (“succes”)  $\Rightarrow$  valori posibile:  $X \in \{0, 1, \dots\}$

$$P(X = k) = p(1 - p)^k \quad \text{pentru } k \in \{0, 1, 2, \dots\}.$$

Matlab/Octave: `geornd( $p, \dots$ )`

**Exemplu:**  $X$  v.a. ce indică numărul de retransmisii printr-un canal cu zgomot (canal cu perturbări) până la prima recepționare corectă a mesajului;  $X$  are distribuție geometrică. ♠

### Variabile aleatoare independente

**Def. 12.** Variabilele aleatoare discrete  $X$  (care ia valorile  $\{x_i, i \in I\}$ ) și  $Y$  (care ia valorile  $\{y_j, j \in J\}$ ) sunt **independente**, dacă și numai dacă

$$P(X = x_i, Y = y_j) = P(X = x_i)P(Y = y_j) \quad \forall i \in I, j \in J.$$

**Notăție 2.**  $P(X = x_i, Y = y_j) = P(\{X = x_i\} \cap \{Y = y_j\}) \quad \forall i \in I, j \in J.$



**P. 8.** Fie variabilele aleatoare discrete  $X$  (care ia valorile  $\{x_i, i \in I\}$ ) și  $Y$  (care ia valorile  $\{y_j, j \in J\}$ ). Sunt echivalente afirmațiile:

- (1)  $X$  și  $Y$  sunt v.a. sunt independente;
- (2)  $P(X = x_i | Y = y_j) = P(X = x_i) \quad \forall i \in I, j \in J$ ;
- (3)  $P(Y = y_j | X = x_i) = P(Y = y_j) \quad \forall i \in I, j \in J$ .
- (4)  $P(X \leq x, Y \leq y) = P(X \leq x) \cdot P(Y \leq y) \quad \forall x, y \in \mathbb{R}$ .

**Def. 13.**  $\mathbb{X} = (X_1, \dots, X_m)$  este un **vector aleator discret** dacă fiecare componentă a sa este o variabilă aleatoare discretă.

Dacă  $\mathbb{X}$  este un vector aleator discret care ia valori în mulțimea  $\mathbb{X}(\Omega) = \{\mathbb{x}_k : k \in K\} \subset \mathbb{R}^m$ , unde  $K \subseteq \mathbb{N}$  este o mulțime de indici, atunci

$$P(\mathbb{X} = \mathbb{x}_k) = P(\{\omega \in \Omega : \mathbb{X}(\omega) = \mathbb{x}_k\}), \quad k \in K,$$

determină **distribuția de probabilitate a vectorului aleator discret**  $\mathbb{X}$

$$\mathbb{X} \sim \left( \begin{array}{c} \mathbb{x}_k \\ P(\mathbb{X} = \mathbb{x}_k) \end{array} \right)_{k \in K}.$$

▷ Vectorii aleatori sunt caracterizați de distribuțiile lor! De exemplu, un vector aleator cu 2 componente:

$$\mathbb{X} = (X, Y) \sim \left( \begin{array}{c} (x_i, y_j) \\ p_{ij} \end{array} \right)_{(i,j) \in I \times J}$$

unde  $I, J \subseteq \mathbb{N}$  sunt mulțimi de indici,

$$p_{ij} = P((X, Y) = (x_i, y_j)) = P(\{X = x_i\} \cap \{Y = y_j\}), \quad p_{ij} > 0 \quad \forall i \in I, j \in J,$$

$$\text{iar } \sum_{(i,j) \in I \times J} p_{ij} = 1.$$

$X \backslash Y$	...	$y_j$	...
$\vdots$	$\vdots$	$\vdots$	$\vdots$
$x_i$	...	$p_{ij}$	...
$\vdots$	$\vdots$	$\vdots$	$\vdots$

▷ Uneori distribuția vectorului  $(X, Y)$  se dă sub formă tabelară:

Observație: Dacă  $X$  și  $Y$  sunt v.a. independente, atunci

$$(1) \quad p_{ij} = P(\{X = x_i\} \cap \{Y = y_j\}) = P(X = x_i)P(Y = y_j) \quad \forall i \in I, j \in J.$$

▷ Dacă  $X$  și  $Y$  sunt v.a. independente, și se știu distribuțiile lor, atunci distribuția vectorului aleator  $(X, Y)$  se determină pe baza formulei (1).

▷ Dacă se cunoaște distribuția vectorului aleator  $(X, Y)$  distribuțiile lui  $X$  și  $Y$  se determină astfel:

$$P(X = x_i) = \sum_{j \in J} p_{ij} \quad \forall i \in I$$

$$P(Y = y_j) = \sum_{i \in I} p_{ij} \quad \forall j \in J.$$

### Operații cu variabile aleatoare (numerice)

- Cunoscând distribuția vectorului  $(X, Y)$  cum se determină distribuția pentru  $X + Y$ ,  $X \cdot Y$ ,  $X^2 - 1$ ,  $2Y$ ?

**Exemplu:** Fie vectorul aleator discret  $(X_1, X_2)$  cu distribuția dată de următorul tabel:

$X_2 \backslash X_1$	0	1
1	$\frac{2}{16}$	$\frac{1}{16}$
2	$\frac{1}{16}$	$\frac{5}{16}$

Determinați:

- distribuțiile variabilelor aleatoare  $X_1$  și  $X_2$ ;
- distribuțiile variabilelor aleatoare  $X_1 + X_2$  și  $X_1 \cdot X_2$ ,  $X_1^2 - 1$ ;
- dacă variabilele aleatoare  $X_1$  și  $X_2$  sunt independente sau dependente.

a)  $X_1 \sim \begin{pmatrix} 1 & 2 \\ \frac{5}{16} & \frac{11}{16} \end{pmatrix}$  și  $X_2 \sim \begin{pmatrix} 0 & 1 & 2 \\ \frac{3}{16} & \frac{6}{16} & \frac{7}{16} \end{pmatrix}$ .

b)  $X_1 + X_2 \sim \begin{pmatrix} 1 & 2 & 3 & 4 \\ \frac{2}{16} & \frac{2}{16} & \frac{7}{16} & \frac{5}{16} \end{pmatrix}$  și  $X_1 \cdot X_2 \sim \begin{pmatrix} 0 & 1 & 2 & 4 \\ \frac{3}{16} & \frac{1}{16} & \frac{7}{16} & \frac{5}{16} \end{pmatrix}$ ,  $X_1^2 - 1 \sim \begin{pmatrix} 0 & 3 \\ \frac{5}{16} & \frac{11}{16} \end{pmatrix}$

c)  $X_1$  și  $X_2$  nu sunt independente, pentru că  $\frac{2}{16} = P(X_1 = 1, X_2 = 0) \neq P(X_1 = 1)P(X_2 = 0) = \frac{5}{16} \cdot \frac{3}{16}$ . ♡

- Cunoscând distribuțiile variabilelor aleatoare independente (discrete)  $X$  și  $Y$ , cum se determină distribuția pentru  $X + Y$ ,  $X \cdot Y$ ?

**Exercițiu:** Fie  $X, Y$  v.a. independente, având distribuțiile

$$X \sim \begin{pmatrix} 0 & 1 \\ \frac{1}{3} & \frac{2}{3} \end{pmatrix}, \quad Y \sim \begin{pmatrix} -1 & 0 & 1 \\ \frac{1}{2} & \frac{1}{4} & \frac{1}{4} \end{pmatrix}$$

- Care sunt distribuțiile v.a.  $2X + 1$ ,  $Y^2$ , dar distribuția vectorului aleator  $(X, Y)$ ?
- Care sunt distribuțiile v.a.  $X + Y$ ,  $X \cdot Y$ ,  $\max(X, Y)$ ,  $\min(X, Y^2)$ ? ♣

**Exercițiu:** Se aruncă două zaruri. a) Să se scrie distribuția de probabilitate pentru variabila aleatoare, care este suma celor două numere apărute. b) Să se scrie distribuția de probabilitate pentru variabila aleatoare, care este produsul celor două numere apărute. ♠

**Def. 14. Valoarea medie a unei variabile aleatoare discrete (numerice)  $X$ , care ia valorile  $\{x_i, i \in I\}$ , este**

$$E(X) = \sum_{i \in I} x_i P(X = x_i),$$

dacă  $\sum_{i \in I} |x_i| P(X = x_i) < \infty$ .

▷ Valoarea medie a unei variabile aleatoare caracterizează *tendința centrală* a valorilor acesteia.

**P. 9.** Fie  $X$  și  $Y$  v.a. discrete. Au loc proprietățile:

→  $E(aX + b) = aE(X) + b$  pentru orice  $a, b \in \mathbb{R}$ ;

→  $E(X + Y) = E(X) + E(Y)$ ;

→ Dacă  $X$  și  $Y$  sunt v.a. independente, atunci  $E(X \cdot Y) = E(X)E(Y)$ .

→ Dacă  $g : \mathbb{R} \rightarrow \mathbb{R}$  e o funcție astfel încât  $g(X)$  este v.a., atunci

$$E(g(X)) = \sum_{i \in I} g(x_i)P(X = x_i),$$

dacă  $\sum_{i \in I} |g(x_i)|P(X = x_i) < \infty$ .

Matlab/Octave: `mean(x)`

pentru  $x = [x(1), \dots, x(n)]$ , se calculează  $\text{mean}(x) = \frac{1}{n}(x(1) + \dots + x(n))$

**Exemplu:** Joc: Se aruncă un zar; dacă apare 6, se câștigă 3 u.m. (unități monetare), dacă apare 1 se câștigă 2 u.m., dacă apare 2,3,4,5 se pierde 1 u.m. În medie cât va câștiga sau pierde un jucător după 30 de repetiții ale jocului?

Răspuns: Fie  $X$  v.a. care indică venitul la un joc

$$X \sim \begin{pmatrix} -1 & 2 & 3 \\ \frac{4}{6} & \frac{1}{6} & \frac{1}{6} \end{pmatrix}$$

Pentru  $i \in \{1, \dots, 30\}$  fie  $X_i$  venitul la al  $i$ -lea joc;  $X_i$  are aceeași distribuție ca  $X$ . Venitul mediu al jucătorului după 30 de repetiții ale jocului este

$$E(X_1 + \dots + X_{30}) = E(X_1) + \dots + E(X_{30}) = 30 \cdot E(X) = 30 \cdot \frac{1}{6} \cdot (2 - 4 + 3) = 5 \text{ (u.m.)}.$$

Așadar jucătorul câștigă în medie 5 u.m.

**Exercițiu:**

Input: Fie  $A(1), \dots, A(200)$  un vector cu 200 de elemente, din care 50 sunt egale cu 0, 70 egale cu 1 și 80 sunt egale cu 2 (ordinea lor este necunoscută).

Output: Să se găsească un 0 în vector, alegând aleator un element din șir și verificând dacă acesta este 0.

**Întrebare:** În medie câte iterații sunt necesare înainte să apară primul 0?

```
clear all
A=[zeros(1,50), zeros(1,70)+1, zeros(1,80)+2];
index=randperm(length(A));
A=A(index);
c=0;
i=randi(length(A));
while A(i)~=0
```

```

c=c+1;
i=randi(length(A));
end
fprintf('nr. iteratii: %d \n',c)

clc
clear all
A=[zeros(1,50), zeros(1,70)+1,zeros(1,80)+2];
s=[];
N=100;
for j=1:N
index=randperm(length(A));
A=A(index);
c=0;
i=randi(length(A));
A(i)
while A(i)~=0
c=c+1;
i=randi(length(A));
end
s=[s,c];
end
mean(s)
fprintf('nr. mediu de iteratii: %4.3f \n',mean(s))

```

Probabilitatea să apară la orice iterație 0 este  $p = \frac{50}{200} = 0.25$ .

Notăm cu  $X$  v.a. care indică numărul de iterații necesare *înainte* să apară primul 0  $\Rightarrow X \sim \text{Geo}(p)$ .

*Numărul mediu* de iterații necesare *înainte* să apară primul 0 este  $E(X)$ . Să se calculeze această valoare medie! ▼

**Def. 15.** Fie  $X_1, \dots, X_n$  cu  $n \in \mathbb{N}$ ,  $n \geq 2$ , variabile aleatoare, care iau valori în mulțimile  $\mathcal{X}_1, \dots, \mathcal{X}_n$ .  $X_1, \dots, X_n$  sunt **variabile aleatoare independente**, dacă

$$P(X_1 = x_1, \dots, X_n = x_n) = P(X_1 = x_1) \cdot \dots \cdot P(X_n = x_n)$$

pentru fiecare  $x_1 \in \mathcal{X}_1, \dots, x_n \in \mathcal{X}_n$

**Notăție 3.** Fie  $\mathbb{U}$  un vector aleator discret care ia valori în  $\mathcal{U} \subset \mathbb{R}^m$  și fie  $\mathbb{V}$  un vector aleator discret care ia valori în  $\mathcal{V} \subset \mathbb{R}^n$ ,  $m, n \in \mathbb{N}^*$ . Notăm  $P[\mathbb{U}] : \mathcal{U} \rightarrow [0, 1]$ ,

$$P[\mathbb{U}](\mathfrak{u}) = P(\mathbb{U} = \mathfrak{u}) \forall \mathfrak{u} \in \mathcal{U}$$

$\hookrightarrow P[\mathbb{U}]$  este **distribuția de probabilitate** a vectorului aleator, dacă  $m > 1$  (a se vedea Def. 13), sau a v.a., dacă  $m = 1$  (a se vedea Def. 11).

► **Observație: (1)** Def. 15 se transcrie mai compact astfel:

$X_1, \dots, X_n$  sunt **variabile aleatoare independente**

$$\iff P[X_1, \dots, X_n] = P[X_1] \cdot \dots \cdot P[X_n].$$

(2) Dacă  $U_1$  and  $U_2$  sunt 2 v.a. discrete atunci:

▷  $U_1$  au  $U_2$  aceeași distribuție  $\iff P[U_1] = P[U_2]$ .

▷  $U_1$  și  $U_2$  sunt v.a. independente  $\iff P[U_1, U_2] = P[U_1]P[U_2] \iff P[U_1|U_2] = P[U_1] \iff P[U_2|U_1] = P[U_2]$ .

**P. 10.** Dacă  $X_1, \dots, X_n$  sunt variabile aleatoare independente, atunci pentru orice indici diferiți  $i_1, \dots, i_k \subset \{1, \dots, n\}$   $X_{i_1}, \dots, X_{i_k}$  sunt variabile aleatoare independente, adică

$$P[X_{i_1}, \dots, X_{i_k}] = P[X_{i_1}] \cdot \dots \cdot P[X_{i_k}].$$

### Clasificarea naivă Bayes

În învățarea automată, clasificatorii bayesieni naivi sunt o familie de clasificatori probabilistici simpli, bazați pe aplicarea formulei lui Bayes (a se vedea P.5) cu ipoteze “naive” de independență condiționată între atribute (în engl. features), cunoscând clasificarea. Pentru unele tipuri de modele de probabilitate, clasificatorii bayesieni naivi pot fi antrenați foarte eficient. În aplicații practice pentru modelele bayesiene naive se folosește *metoda probabilității maxime*. Noțiunea folosită în acest context este condițional independența între evenimente, respectiv v.a.

Fie  $(\Omega, \mathcal{K}, P)$  un spațiu de probabilitate. De asemenea considerăm că toate probabilitățile condiționate sunt definite (adică condiționarea se face în raport cu un eveniment a cărui probabilitate nu este 0).

**Def. 16.** Evenimentele  $A, B \in \mathcal{K}$  sunt **condițional independente**, cunoscând evenimentul  $C \in \mathcal{K}$ , dacă și numai dacă

$$P(A \cap B|C) = P(A|C)P(B|C).$$

**P. 11.** Au loc echivalențele:

$$P(A \cap B|C) = P(A|C)P(B|C) \iff P(A|B \cap C) = P(A|C) \iff P(B|A \cap C) = P(B|C).$$

*Demonstrație:* • Pentru prima echivalență: “ $\Rightarrow$ ”

$$P(A|B \cap C) = \frac{P(A \cap B \cap C)}{P(B \cap C)} = \frac{P(A \cap B|C)P(C)}{P(B|C)P(C)} = \frac{P(A|C)P(B|C)}{P(B|C)} = P(A|C).$$

“ $\Leftarrow$ ”

$$P(A \cap B|C) = \frac{P(A \cap B \cap C)}{P(C)} = \frac{P(A|B \cap C)P(B \cap C)}{P(C)} = \frac{P(A|C)P(B \cap C)}{P(C)} \\ = P(A|C)P(B|C).$$

•  $P(A \cap B|C) = P(A|C)P(B|C) \Leftrightarrow P(B|A \cap C) = P(B|C)$  se demonstrează analog. ■

**Exemplu:** Fie  $Z$  o v.a. care indică rezultatul aruncării unui zar. Considerăm evenimentele:  $A = (Z \in \{1, 2\})$ ,  $B = (Z \in \{2, 4, 6\})$  și  $C = (Z \in \{1, 4\})$ . Să se arate că:

a)  $A$  și  $B$  sunt independente;

b)  $A$  și  $B$  nu sunt condițional independente, cunoscând evenimentul  $C$ .

R.: a)  $P(A \cap B) = P(Z = 2) = \frac{1}{6} = \frac{1}{3} \cdot \frac{1}{2} = P(Z \in \{1, 2\})P(Z \in \{2, 4, 6\}) = P(A)P(B) \Rightarrow A$  și  $B$  sunt independente.

b)  $P(A \cap B|C) = P(Z = 2|Z \in \{1, 4\}) = 0$ ,  $P(A|C) = P(Z = 1|Z \in \{1, 4\}) = \frac{1}{2} = P(Z = 4|Z \in \{1, 4\}) = P(B|C) \Rightarrow A$  și  $B$  nu sunt condițional independente, cunoscând  $C$ . ▲

**Exercițiu:** Într-o cutie sunt 2 zaruri. La primul zar 3 apare cu probabilitatea  $\frac{1}{6}$ , iar la celălalt zar (care e măsluit) 3 apare cu probabilitatea  $\frac{5}{6}$ . Se alege aleator un zar, care este apoi aruncat de 2 ori. Considerăm evenimentele

$A_i$ : “zarul ales indică 3 la aruncarea  $i$ ”,  $i \in \{1, 2\}$

$Z_j$ : “se alege zarul  $j$ ”,  $j \in \{1, 2\}$ .

Sunt  $A_1$  și  $A_2$  condițional independente, cunoscând  $Z_1$ ? Sunt  $A_1$  și  $A_2$  independente?

R.: Dacă se cunoaște tipul zarului ales, atunci aruncările sunt în mod evident independente:

$P(A_1 \cap A_2|Z_1) = \frac{1}{36} = P(A_1|Z_1) \cdot P(A_2|Z_1)$ . Din formula probabilității totale P.5 avem:

$P(A_2) = P(A_1) = P(A_1|Z_1)P(Z_1) + P(A_1|Z_2)P(Z_2) = \frac{1}{6} \cdot \frac{1}{2} + \frac{5}{6} \cdot \frac{1}{2} = \frac{1}{2}$  și  $P(A_1 \cap A_2) = P(A_1 \cap A_2|Z_1)P(Z_1) + P(A_1 \cap A_2|Z_2)P(Z_2) = \frac{1}{6} \cdot \frac{1}{6} \cdot \frac{1}{2} + \frac{5}{6} \cdot \frac{5}{6} \cdot \frac{1}{2} = \frac{13}{36}$ . Deci  $P(A_2|A_1) = \frac{P(A_1 \cap A_2)}{P(A_1)} = \frac{13}{18}$ .  $A_1$  și  $A_2$  nu sunt independente, pentru că  $P(A_2|A_1) \neq P(A_2)$ . \*

**Def. 17.** Fie  $X, Y, Z$  v.a. discrete care iau valori în mulțimile  $\mathcal{X}, \mathcal{Y}, \mathcal{Z}$ . V.a.  $X$  este **condițional independentă** de  $Y$ , cunoscând (știind) v.a.  $Z$ , dacă pentru fiecare  $x \in \mathcal{X}, y \in \mathcal{Y}, z \in \mathcal{Z}$ , are loc

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

**Notăție 4.** Notăm cu  $P[\mathbb{U}|\mathbb{V}] : \mathcal{U} \times \mathcal{V} \rightarrow [0, 1]$  distribuția de probabilitate condiționată

$$P[\mathbb{U}|\mathbb{V}](u, v) = P(\mathbb{U} = u|\mathbb{V} = v) \forall u \in \mathcal{U}, v \in \mathcal{V}.$$

► **Observație:** Def. 17 se transcrie astfel:

$$(2) \quad P[X, Y|Z] = P[X|Z]P[Y|Z].$$

Folosind **P.11**, rezultă:

**P. 12.** V.a.  $X$  și  $Y$  sunt condițional independente, cunoscând  $Z \Leftrightarrow P[X, Y|Z] = P[X|Z]P[Y|Z] \Leftrightarrow P[X|Y, Z] = P[X|Z] \Leftrightarrow P[Y|X, Z] = P[Y|Z]$ .

Vom introduce noțiunea de condițional independență pentru mai multe v.a. discrete.

**Def. 18.** Fie  $X, Y_1, \dots, Y_m, Z_1, \dots, Z_n$  v.a. discrete. V.a.  $X$  este **condițional independentă** de  $Y_1, \dots, Y_m$ , știind (cunoscând) v.a.  $Z_1, \dots, Z_n$ , dacă are loc

$$P[X, Y_1, \dots, Y_m | Z_1, \dots, Z_n] = P[X | Z_1, \dots, Z_n] P[Y_1, \dots, Y_m | Z_1, \dots, Z_n].$$

**P. 13.** Fie  $X$  o v.a. discretă condițional independentă de v.a. discrete  $Y_1, \dots, Y_m$ , cunoscând v.a. discrete  $Z_1, \dots, Z_n$ . Dacă  $i_1, \dots, i_k \in \{1, \dots, m\}$  sunt indici diferiți, atunci are loc

$$P[X | Y_{i_1}, \dots, Y_{i_k}, Z_1, \dots, Z_n] = P[X | Z_1, \dots, Z_n],$$

$$P[Y_{i_1}, \dots, Y_{i_k} | X, Z_1, \dots, Z_n] = P[Y_{i_1}, \dots, Y_{i_k} | Z_1, \dots, Z_n].$$

### Exemplu de clasificare naivă Bayes

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

- **vreme**  $V$ : ploaie  $p$ , zăpadă  $z$ , senin  $s$ , înnorat  $\hat{i}$  (dar nu plouă și nu ninge) ;
- **timp**  $Ti$ : dimineață  $di$ , amiază  $am$ , seară  $se$ , noapte  $no$ .

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

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

Considerăm evenimentul următor, denumit *vector de attribute*:

$$E = (V = p) \cap (Ti = am).$$

Se caută o clasă pentru  $E$ , stabilind care din următoarele probabilități este mai mare:  $P(\mathbf{T} = \mathbf{a}|E)$  sau  $P(\mathbf{T} = \mathbf{r}|E)$ ; aceasta este **metoda de probabilitate maximă**. Știind că vremea este ploioasă și este amiază, ce *previziune* se poate face despre trafic?

Se face următoarea **presupunere naivă**: **atributele sunt condițional independente**, dacă se **dă clasificarea**. De asemenea, presupunem că toate probabilitățile condiționate sunt definite (adică condiționarea este în raport cu un eveniment a cărui probabilitate nu este 0).

Presupunerea *naivă* în acest exemplu este:  $V$  și  $Ti$  sunt condițional independente, cunoscând (știind)  $\mathbf{T}$  (a se vedea Def. 17, respectiv relația (2)), adică

$$(3) \quad P[V, Ti|\mathbf{T}] = P[V|\mathbf{T}]P[Ti|\mathbf{T}].$$

Această notație se transcrie prin

$$(4) \quad P(V = v, Ti = ti|\mathbf{T} = \mathbf{t}) = P(V = v|\mathbf{T} = \mathbf{t})P(Ti = ti|\mathbf{T} = \mathbf{t}),$$

pentru fiecare  $v \in \{p, z, s, \hat{i}\}$ ,  $ti \in \{di, am, se, no\}$ ,  $\mathbf{t} \in \{\mathbf{a}, \mathbf{r}\}$ . De exemplu, avem:

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

altfel spus: *probabilitatea să plouă, știind că traficul este aglomerat, nu depinde de timp (adică este independentă de timp)*.

Din (3) și P. 12 se deduce că  $P[V|Ti, \mathbf{T}] = P[V|\mathbf{T}]$  și  $P[Ti|V, \mathbf{T}] = P[Ti|\mathbf{T}]$ , adică  $P(V = v|Ti = ti, \mathbf{T} = \mathbf{t}) = P(V = v|\mathbf{T} = \mathbf{t})$  și  $P(Ti = ti|V = v, \mathbf{T} = \mathbf{t}) = P(Ti = ti|\mathbf{T} = \mathbf{t})$  pentru fiecare  $v \in \{p, z, s, \hat{i}\}$ ,  $ti \in \{di, am, se, no\}$ ,  $\mathbf{t} \in \{\mathbf{a}, \mathbf{r}\}$ .

► Folosind datele din tabel, determinăm mai întâi probabilitățile claselor și probabilitățile condiționate ale atributelor, cunoscând clasa.

$\mathbf{T} = \mathbf{a}$	$\mathbf{T} = \mathbf{r}$	$P(\mathbf{T} = \mathbf{a})$	$P(\mathbf{T} = \mathbf{r})$
9	6	$\frac{9}{15}$	$\frac{6}{15}$

$V$	$\mathbf{T} = \mathbf{a}$	$\mathbf{T} = \mathbf{r}$	$P(V = \dots \mathbf{T} = \mathbf{a})$	$P(V = \dots \mathbf{T} = \mathbf{r})$
$p$	2	1	$\frac{2}{9}$	$\frac{1}{6}$
$z$	3	1	$\frac{3}{9}$	$\frac{1}{6}$
$s$	1	3	$\frac{1}{9}$	$\frac{3}{6}$
$\hat{i}$	3	1	$\frac{3}{9}$	$\frac{1}{6}$

$Ti$	$\mathbf{T} = \mathbf{a}$	$\mathbf{T} = \mathbf{r}$	$P(Ti = \dots \mathbf{T} = \mathbf{a})$	$P(Ti = \dots \mathbf{T} = \mathbf{r})$
$di$	3	1	$\frac{3}{9}$	$\frac{1}{6}$
$am$	2	1	$\frac{2}{9}$	$\frac{1}{6}$
$se$	3	1	$\frac{3}{9}$	$\frac{1}{6}$
$no$	1	3	$\frac{1}{9}$	$\frac{3}{6}$



► Pe baza formulei lui Bayes P. 5 și a ipotezei de independență condiționată, deducem că:

$$\begin{aligned} P(\mathbf{T} = \mathbf{a}|E) &= \frac{P(E|\mathbf{T} = \mathbf{a})P(\mathbf{T} = \mathbf{a})}{P(E)} = \frac{P(V = p, Ti = am|\mathbf{T} = \mathbf{a})P(\mathbf{T} = \mathbf{a})}{P(E)} \\ &= \frac{P(V = p|\mathbf{T} = \mathbf{a})P(Ti = am|\mathbf{T} = \mathbf{a})P(\mathbf{T} = \mathbf{a})}{P(E)} = \frac{\frac{2}{9} \cdot \frac{2}{9} \cdot \frac{9}{15}}{P(E)} = \frac{1}{P(E)} \cdot \frac{4}{135} \end{aligned}$$

și

$$\begin{aligned} P(\mathbf{T} = \mathbf{r}|E) &= \frac{P(E|\mathbf{T} = \mathbf{r})P(\mathbf{T} = \mathbf{r})}{P(E)} = \frac{P(V = p, Ti = am|\mathbf{T} = \mathbf{r})P(\mathbf{T} = \mathbf{r})}{P(E)} \\ &= \frac{P(V = p|\mathbf{T} = \mathbf{r})P(Ti = am|\mathbf{T} = \mathbf{r})P(\mathbf{T} = \mathbf{r})}{P(E)} = \frac{\frac{1}{6} \cdot \frac{1}{6} \cdot \frac{6}{15}}{P(E)} = \frac{1}{P(E)} \cdot \frac{1}{90}. \end{aligned}$$

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

$$E = (V = p) \cap (Ti = am) \text{ clasa } \mathbf{T} = \mathbf{a}.$$

Această concluzie se poate completa la datele din tabelul de date inițial!

► În plus, putem determina  $P(E) = P(V = p, Ti = am)$  astfel: Scriem

$$1 = P(\mathbf{T} = \mathbf{a}|E) + P(\mathbf{T} = \mathbf{r}|E) = \frac{1}{P(E)} \left( \frac{4}{135} + \frac{1}{90} \right)$$

și deducem  $P(E) = P(V = p, Ti = am) = \frac{11}{270}$ .

★

## Rețele Bayes

Rețeaua Bayes este un graf orientat aciclic (i.e. nu conține niciun drum orientat închis).

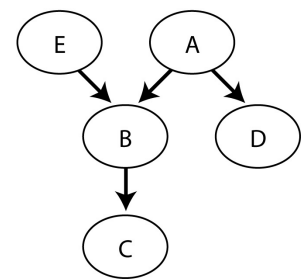
► Nodul  $Y$  este **părinte** pentru nodul  $X$ , dacă există o muchie orientată de la  $Y$  la  $X$ . Mulțimea părinților lui  $X$  se notează cu  $p(X)$ . Dacă  $X$  este nod rădăcină, atunci  $p(X) = \emptyset$ . De exemplu:  $p(B) = \{E, A\}$ ,  $p(D) = \{A\}$ ,  $p(C) = \{B\}$ ,  $p(E) = p(A) = \emptyset$ .

► Nodul  $Y$  este **descendent** al nodului  $X$ , dacă există un drum orientat de la  $X$  la  $Y$ . Mulțimea descendenților lui  $X$  se notează cu  $de(X)$ . De exemplu:  $d(E) = \{B, C\}$ ,  $d(A) = \{B, C, D\}$ ,  $d(B) = \{C\}$ ,  $d(D) = \emptyset$ .

Într-o rețea în care există o structură cauzală, nodurile din  $p(X)$  reprezintă *cauzele* pentru  $X$ , iar nodurile din  $de(X)$  sunt *efectele* nodului  $X$ .

► Nodul  $Y$  este **nondescendent** al nodului  $X$ , dacă nu este descendent al nodului  $X$ . Mulțimea nondescendenților lui  $X$  se notează cu  $nd(X)$ . De exemplu:  $nd(E) = \{A, D\}$ ,  $nd(A) = \{E\}$ ,  $nd(B) = \{E, A, D\}$ ,  $nd(D) = \{E, A, B, C\}$ ,  $nd(C) = \{E, A, B, D\}$ .

► Fiecare nod  $X_1, \dots, X_n$  din rețea este identificat cu o variabilă aleatoare și este definit pe



același spațiu de probabilitate  $(\Omega, \mathcal{K}, P)$ ; probabilitățile  $P[X_j|p(X_j)]$ ,  $j = \overline{1, n}$  sunt date; are loc convenția  $P[X_j|p(X_j)] = P[X_j]$ , dacă  $X_j$  este nod rădăcină ( $P[X_j]$  este distribuția de probabilitate a lui  $X_j$ , a se vedea Notăția 3, iar  $P[X_j|p(X_j)]$  este distribuția de probabilitate condiționată, a se vedea Notăția 4).

► **Proprietatea rețelei Bayes:** orice nod  $X$  și nondescendenții săi  $nd(X)$  sunt *condițional independenți*, dacă se cunosc valorile părinților  $p(X)$ ; dacă  $p(X) = \emptyset$ , atunci  $X$  și  $nd(X)$  sunt *independenți*.

**Exemplul 1:** Se dă rețeaua Bayes din figura alăturată, în care  $X_1, \dots, X_6$  sunt variabile aleatoare binare.

▷ Au loc proprietățile:

• Mulțimile de noduri corespunzătoare părinților, descendenților, nondescendenților sunt:

$$p(X_1) = \emptyset, p(X_2) = \{X_1\}, p(X_3) = \{X_1, X_2\},$$

$$p(X_4) = p(X_5) = \{X_3\}, p(X_6) = \{X_4, X_5\}$$

$$de(X_1) = \{X_2, X_3, X_4, X_5, X_6\},$$

$$de(X_2) = \{X_3, X_4, X_5, X_6\},$$

$$de(X_3) = \{X_4, X_5, X_6\},$$

$$de(X_4) = de(X_5) = \{X_6\}, de(X_6) = \emptyset,$$

$$nd(X_2) = \{X_1\}, nd(X_3) = \{X_1, X_2\},$$

$$nd(X_4) = \{X_1, X_2, X_3, X_5\}$$

$$nd(X_5) = \{X_1, X_2, X_3, X_4\},$$

$$nd(X_6) = \{X_1, X_2, X_3, X_4, X_5\};$$

• probabilitățile (asociate nodurilor), care definesc rețeaua Bayes sunt:

$$P[X_1], P[X_2|X_1], P[X_3|X_1, X_2], P[X_4|X_3], P[X_5|X_3], P[X_6|X_4, X_5];$$

• independențe condiționate:

▷  $X_4$  este condițional independentă de  $X_1, X_2, X_5$ , cunoscând  $X_3$

$$\Rightarrow P[X_4|X_1, X_2, X_3, X_5] = P[X_4|X_3],$$

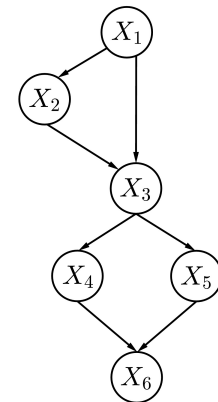
▷  $X_5$  este condițional independentă de  $X_1, X_2, X_4$ , cunoscând  $X_3$

$$\Rightarrow P[X_5|X_1, X_2, X_3, X_4] = P[X_5|X_3],$$

▷  $X_6$  este condițional independentă de  $X_1, X_2, X_3$ , cunoscând  $X_4, X_5$

$$\Rightarrow P[X_6|X_1, X_2, X_3, X_4, X_5] = P[X_6|X_4, X_5];$$

• (exemplu de calcul în rețea Bayes) se știe  $P(X_1=1)=0.5$ ,  $P(X_2=1|X_1=1)=0.6$ ,  $P(X_3=1|X_1=1, X_2=1)=0.5$ ,  $P(X_4=1|X_3=1)=0.4$ ,  $P(X_4=1|X_3=0)=0.3$ , atunci să se cal-



**Rețea Bayes**

culeze  $P(X_4=1, X_2=1, X_1=1)$ :

$$\begin{aligned}
 &P(X_4=1, X_2=1, X_1=1) \\
 &= P(X_4=1, X_3=1, X_2=1, X_1=1) + P(X_4=1, X_3=0, X_2=1, X_1=1) \\
 &= P(X_1=1)P(X_2=1|X_1=1)P(X_3=1|X_1=1, X_2=1)P(X_4=1|X_3=1) \\
 &\quad + P(X_1=1)P(X_2=1|X_1=1)P(X_3=0|X_1=1, X_2=1)P(X_4=1|X_3=0) \\
 &= 0.105.
 \end{aligned}$$



**Exemplul 2** Se dă rețeaua Bayes din figura alăturată, în care  $X_1, \dots, X_5$  sunt variabile aleatoare binare. Se știu probabilitățile:

$$P(X_1 = 0) = 0.4, P(X_2 = 0|X_1 = 0) = 0.2,$$

$$P(X_2 = 0|X_1 = 1) = 0.5, P(X_3 = 0|X_1 = 0) = 0.3,$$

$$P(X_3 = 0|X_1 = 1) = 0.4,$$

$$P(X_4 = 0|X_2 = 0) = 0.2, P(X_4 = 0|X_2 = 1) = 0.5,$$

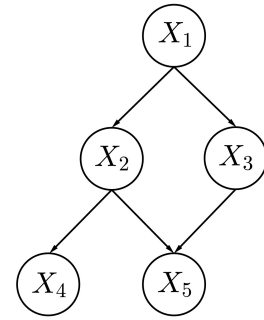
$$P(X_5 = 0|X_2 = 0, X_3 = 0) = 0.5,$$

$$P(X_5 = 0|X_2 = 0, X_3 = 1) = 0.2,$$

$$P(X_5 = 0|X_2 = 1, X_3 = 0) = 0.7,$$

$$P(X_5 = 0|X_2 = 1, X_3 = 1) = 0.4.$$

a) Să se calculeze



Rețea Bayes

$$P(X_3 = 1|X_2 = 1), P(X_1 = 0, X_3 = 1), P\left(\bigcap_{i=1}^5 \{X_i = 1\}\right).$$

b) Să se scrie distribuția de probabilitate a variabilei aleatoare  $X_3$ .

*Rezolvare:* Se calculează:  $P(X_1 = 1) = 1 - P(X_1 = 0) = 0.6$

$$P(X_2 = 1|X_1 = 0) = 1 - P(X_2 = 0|X_1 = 0) = 0.8;$$

$$P(X_2 = 1|X_1 = 1) = 1 - P(X_2 = 0|X_1 = 1) = 0.5;$$

$$P(X_3 = 1|X_1 = 0) = 1 - P(X_3 = 0|X_1 = 0) = 0.7;$$

$$P(X_3 = 1|X_1 = 1) = 1 - P(X_3 = 0|X_1 = 1) = 0.6;$$

$$P(X_4 = 1|X_2 = 0) = 1 - P(X_4 = 0|X_2 = 0) = 0.8;$$

$$P(X_4 = 1|X_2 = 1) = 1 - P(X_4 = 0|X_2 = 1) = 0.5;$$

$$P(X_5 = 1|X_2 = 1, X_3 = 1) = 1 - P(X_5 = 0|X_2 = 1, X_3 = 1) = 0.6.$$

a) Are loc:

$$P(X_3 = 1|X_2 = 1) = \frac{P(X_3 = 1, X_2 = 1)}{P(X_2 = 1)}.$$

Folosind formula probabilităților totale și proprietățile rețelelor Bayes ( $X_2$  este condițional independentă de  $X_3$ , cunoscând  $X_1$ )<sup>1</sup>:

$$\begin{aligned}
 & \bullet P(X_3 = 1, X_2 = 1) = P(X_3 = 1, X_2 = 1|X_1 = 0)P(X_1 = 0) \\
 & \quad + P(X_3 = 1, X_2 = 1|X_1 = 1)P(X_1 = 1) = \\
 & = P(X_3 = 1|X_1 = 0)P(X_2 = 1|X_1 = 0)P(X_1 = 0) \\
 & \quad + P(X_3 = 1|X_1 = 1)P(X_2 = 1|X_1 = 1)P(X_1 = 1) \\
 & \bullet P(X_2 = 1) = P(X_2 = 1|X_1 = 0)P(X_1 = 0) + P(X_2 = 1|X_1 = 1)P(X_1 = 1).
 \end{aligned}$$

Are loc

$$P(X_1 = 0, X_3 = 1) = P(X_3 = 1|X_1 = 0)P(X_1 = 0).$$

Folosind regula de înmulțire și proprietățile rețelelor Bayes ( $X_2$  este condițional independentă de  $X_3$ , cunoscând  $X_1$ ;  $X_4$  este condițional independentă de  $X_1, X_3$ , cunoscând  $X_2$ ;  $X_5$  este condițional independentă de  $X_1, X_4$ , cunoscând  $X_2, X_3$ )

$$\begin{aligned}
 & P(X_1 = 1, X_2 = 1, X_3 = 1, X_4 = 1, X_5 = 1) \\
 & = P(X_1 = 1)P(X_2 = 1|X_1 = 1)P(X_3 = 1|X_1 = 1, X_2 = 1) \cdot \\
 & \quad \cdot P(X_4 = 1|X_1 = 1, X_2 = 1, X_3 = 1)P(X_5 = 1|X_1 = 1, X_2 = 1, X_3 = 1, X_4 = 1) = \\
 & = P(X_1=1)P(X_2=1|X_1=1)P(X_3 = 1|X_1 = 1)P(X_4 = 1|X_2 = 1)P(X_5 = 1|X_2 = 1, X_3 = 1).
 \end{aligned}$$

$$\begin{aligned}
 \text{b) } P(X_3 = 0) & = P(X_3 = 0|X_1 = 0)P(X_1 = 0) + P(X_3 = 0|X_1 = 1)P(X_1 = 1) \\
 & = 0.12 + 0.24 = 0.36 \Rightarrow P(X_3 = 1) = 0.64
 \end{aligned}$$

$$\Rightarrow X_3 \sim \begin{pmatrix} 0 & 1 \\ 0.36 & 0.64 \end{pmatrix}.$$



### Variabile aleatoare continue

V.a. continuă: ia un număr infinit și nenumărabil de valori într-un interval sau reuniune de intervale (v.a. poate lua orice valoare din intervalul considerat);

▷ v.a. continue pot modela caracteristici fizice precum timp (de ex. timp de instalare, timp de așteptare), greutate, lungime, poziție, volum, temperatură (de ex.  $X$  e v.a. care indică durata de funcționare a unui dispozitiv până la prima defectare;  $X$  e v.a. care indică temperatura într-un oraș la ora amiezii)

▷ ea este caracterizată de o funcție de densitate.

**Def. 19. Funcția de densitate**  $f : \mathbb{R} \rightarrow \mathbb{R}$  a unei v.a. continue este funcția pentru care are loc

$$P(X \leq x) = \int_{-\infty}^x f(t)dt, \forall x \in \mathbb{R}.$$

<sup>1</sup>Orice nod  $X$  și nondescendenții săi  $nd(X)$  sunt condițional independenți, dacă se dau valorile părinților  $p(X)$ .

**P. 14.** Fie  $f$  funcția de densitate a unei v.a. continue  $X$ . Au loc proprietățile:

(1)  $f(t) \geq 0$  pentru orice  $t \in \mathbb{R}$ ;

(2)  $\int_{-\infty}^{\infty} f(t) dt = 1$ ;

(3)  $P(a < X \leq b) = \int_a^b f(t) dt \forall a, b \in \mathbb{R}, a < b$ ;

(4)  $P(X = a) = 0 \forall a \in \mathbb{R}$ ;

(5) pentru  $\forall a < b, a, b \in \mathbb{R}$  au loc

$$P(a \leq X \leq b) = P(a < X \leq b) = P(a \leq X < b) = P(a < X < b) = \int_a^b f(t) dt.$$

**Observație:** Orice funcție  $f : \mathbb{R} \rightarrow \mathbb{R}$ , care are proprietățile (1), (2) din **P.14** este o funcție de densitate.

### Exemple de distribuții clasice continue

➔ **Distribuția uniformă pe un interval  $[a, b]$ :**  $X \sim \text{Unif}[a, b]$ ,  $a, b \in \mathbb{R}, a < b$

• funcția de densitate este

$$f(t) = \begin{cases} \frac{1}{b-a}, & \text{pentru } t \in [a, b] \\ 0, & \text{pentru } t \in \mathbb{R} \setminus [a, b] \end{cases}$$

Matlab/Octave:

▷ când  $a = 0, b = 1$  `rand(M, N)` returnează o matrice  $M \times N$  cu valori aleatoare din  $[0, 1]$

▷ `unifrnd(a, b, M, N)`, respectiv  $(b-a) \text{ rand}(M, N) + a$  returnează o matrice  $M \times N$  cu valori aleatoare din  $[a, b]$



**Friedrich Gauss și legea normală  $N(m, \sigma^2)$  (bancnota de 10 DM)**

➡ **Distribuția normală (Gauss):**  $X \sim N(m, \sigma^2)$ ,  $m \in \mathbb{R}$ ,  $\sigma > 0$

- funcția de densitate este

$$f(t) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(t-m)^2}{2\sigma^2}\right), t \in \mathbb{R}.$$

- Pentru  $m = 0, \sigma = 1$ :  $N(0, 1)$  se numește *distribuția standard normală*.
- Distribuția normală se aplică în: măsurarea erorilor (de ex. termenul eroare în analiza regresională), în statistică (teorema limită centrală, teste statistice) etc.

Matlab/Octave: `normrnd(m, sigma, ...)`

➡ **Distribuția exponențială:**  $X \sim \text{Exp}(\lambda)$ ,  $\lambda > 0$

- funcția de densitate este

$$f(t) = \begin{cases} \lambda e^{-\lambda t}, & \text{pentru } t > 0 \\ 0, & \text{pentru } t \leq 0 \end{cases}$$

Matlab/Octave: `exprnd(1/lambda, ...)`

➡ **Distribuția Student:**  $X \sim St(n)$ ,  $n \in \mathbb{N}^*$

- distribuția Student cu  $n \in \mathbb{N}^*$  grade de libertate are funcția de densitate

$$f(t) = \frac{\Gamma\left(\frac{n+1}{2}\right)}{\sqrt{n\pi}\Gamma\left(\frac{n}{2}\right)} \left(1 + \frac{t^2}{n}\right)^{-\frac{n+1}{2}} \quad t \in \mathbb{R}$$

unde

$$\Gamma(a) = \int_0^{\infty} v^{a-1} \exp(-v) dv, \quad a > 0$$

este funcția Gamma.

Matlab/Octave: `trnd(n, ...)`

➡ **Distribuția Chi-pătrat:**  $X \sim \chi^2(n)$ ,  $n \in \mathbb{N}^*$

- distribuția  $\chi^2$  cu  $n \in \mathbb{N}^*$  grade de libertate are funcția de densitate

$$f(t) = \begin{cases} 0 & \text{dacă } x \leq 0 \\ \frac{1}{\Gamma\left(\frac{n}{2}\right)2^{\frac{n}{2}}} x^{\frac{n}{2}-1} \exp\left(-\frac{x}{2}\right) & \text{dacă } x > 0, \end{cases}$$

Matlab/Octave: `chi2rnd(n, ...)`

**Exemplu:** Fie  $X \sim \text{Exp}(0.5)$  v.a. care indică timpul de funcționare a unei baterii (câte luni funcționează bateria). Folosind simulări, să se estimeze a)  $P(2 \leq X \leq 4)$ ; b)  $P(X > 3)$  și să se compare rezultatele obținute cu rezultatele teoretice.

```
X=exprnd(2,1,10000);
p=mean((2<=X) & (X<=4))
q=mean(X>3)
> p = 0.23280
> q= 0.22060
```

$$P(2 \leq X \leq 4) = \int_2^4 0.5e^{-0.5t} dt = -e^{0.5t} \Big|_2^4 = e^{-1} - e^{-2} \approx 0.23254$$

$$P(X > 3) = \int_3^\infty 0.5e^{-0.5t} dt = -e^{0.5t} \Big|_3^\infty = e^{-1.5} \approx 0.22313$$



**Def. 20. Funcția de repartiție**  $F : \mathbb{R} \rightarrow [0, 1]$  a unei variabile aleatoare  $X$  (discrete sau continue) este

$$F(x) = P(X \leq x) \quad \forall x \in \mathbb{R}.$$

**P. 15.** Funcția de repartiție  $F$  a unei variabile aleatoare  $X$  (discrete sau continue) are următoarele proprietăți:

(1)  $F$  este monoton crescătoare, adică pentru orice  $x_1 < x_2$  rezultă  $F(x_1) \leq F(x_2)$ .

(2)  $\lim_{x \rightarrow \infty} F(x) = 1$  și  $\lim_{x \rightarrow -\infty} F(x) = 0$ .

(3)  $F$  este continuă la dreapta, adică  $\lim_{x \searrow x_0} F(x) = F(x_0) \quad \forall x_0 \in \mathbb{R}$ .

(4)  $P(a < X \leq b) = F(b) - F(a) \quad \forall a, b \in \mathbb{R}, a < b$ .

### Observație importantă:

▷ Orice funcție  $F : \mathbb{R} \rightarrow \mathbb{R}$ , care are proprietățile (1), (2), (3) din **P.15** este o funcție de repartiție.

Matlab/Octave:

Distribuția v.a. discrete $X$	Generare valori aleatoare	Funcția de repartiție $F_X(x)$	Probabilitate $P(X = x)$
$Bino(n, p)$	<code>binornd(n, p)</code>	<code>binocdf(x, n, p)</code>	<code>binopdf(x, n, p)</code>
$Unid(n)$	<code>unidrnd(n)</code>	<code>unidcdf(x, n)</code>	<code>unidpdf(x, n)</code>
$Hyge(n, n_1, n_2)$	<code>hygernd(n<sub>1</sub>+n<sub>2</sub>, n<sub>1</sub>, n)</code>	<code>hygecdf(x, n<sub>1</sub>+n<sub>2</sub>, n<sub>1</sub>, n)</code>	<code>hygepdf(x, n<sub>1</sub>+n<sub>2</sub>, n<sub>1</sub>, n)</code>
$Geo(p)$	<code>geornd(p)</code>	<code>geocdf(x, p)</code>	<code>geopdf(x, p)</code>

Distribuția v.a. continue $X$	Generare valori aleatoare	Funcția de repartiție $F_X(x)$	Funcția de densitate $f_X(x)$
$Unif[a, b]$	<code>unifrnd(a, b)</code>	<code>unifcdf(x, a, b)</code>	<code>unifpdf(x, a, b)</code>
$N(m, \sigma^2)$	<code>normrnd(m, sigma)</code>	<code>normcdf(x, m, sigma)</code>	<code>normpdf(x, m, sigma)</code>
$Exp(\lambda)$	<code>exprnd(1/lambda)</code>	<code>expcdf(x, 1/lambda)</code>	<code>exppdf(x, 1/lambda)</code>

**V.a. discretă**

- caracterizată de distribuția de probabilitate discretă

$$X \sim \left( \begin{matrix} x_i \\ P(X = x_i) \end{matrix} \right)_{i \in I}$$

- $\sum_{i \in I} P(X = x_i) = 1$
- funcția de repartiție  $F(x) = P(X \leq x) \forall x \in \mathbb{R}$
- $F(x) = \sum_{i \in I: x_i \leq x} P(X = x_i) \forall x \in \mathbb{R}$
- $F$  este funcție continuă la dreapta
- $F$  este discontinuă în punctele  $x_i, \forall i \in I$
- $\forall a < b, a, b \in \mathbb{R}$   

$$P(a \leq X \leq b) = \sum_{\substack{i \in I \\ a \leq x_i \leq b}} P(X = x_i)$$

- $P(X = a) = 0$  dacă  $a \notin \{x_i : i \in I\}$

**V.a. continuă**

- caracterizată de funcția de densitate  $f$

$$P(X \leq x) = \int_{-\infty}^x f(t) dt$$

- $\int_{-\infty}^{\infty} f(t) dt = 1$
- funcția de repartiție  $F(x) = P(X \leq x) \forall x \in \mathbb{R}$
- $F(x) = \int_{-\infty}^x f(t) dt \quad \forall x \in \mathbb{R}$
- $F$  este funcție continuă în orice punct  $x \in \mathbb{R}$

- $\forall a < b, a, b \in \mathbb{R}$

$$P(a \leq X \leq b) = P(a < X \leq b)$$

$$= P(a \leq X < b) = P(a < X < b) = \int_a^b f(t) dt$$

- $P(X = a) = \int_a^a f(t) dt = 0 \quad \forall a \in \mathbb{R}$

- dacă  $F$  este derivabilă în punctul  $x$   
 $\Rightarrow F'(x) = f(x).$

**Exemplu:** Fie  $X$  v.a. care indică timpul de funcționare neîntreruptă (în ore) până la prima defectare a unui aparat, pentru care  $P(X > x) = 2^{-x}, x > 0$  și  $P(X > x) = 1, x \leq 0$ . Să se determine  $f_X$  și  $P(2 < X < 3)$ .

**Exemplu:** V.a.  $X$  urmează legea uniformă pe  $[a, b]$  cu  $0 < a < b$ . Să se arate că  $P(X > 0) = 1$  și să se determine funcția densitate de probabilitate a variabilei aleatoare  $Y = \ln \left( \frac{1}{X} \right)$ .

**Soluție:** Funcția de densitate a lui  $X$  este:

$$f_X(x) = \begin{cases} \frac{1}{b-a}, & \text{pentru } x \in [a, b] \\ 0, & \text{pentru } x \in \mathbb{R} \setminus [a, b] \end{cases}.$$

Atunci are loc

$$P(X > 0) = 1 - P(X \leq 0) = 1 - \int_{-\infty}^0 f_X(x) dx = 1.$$



Scriem succesiv

$$F_Y(y) = P(Y < y) = P\left(\ln\left(\frac{1}{X}\right) < y\right) = 1 - F_X(e^{-y}).$$

Derivăm în raport cu  $y$

$$f_Y(y) = f_X(e^{-y})e^{-y}.$$

Folosind definiția lui  $f_X$ , obținem

$$f_Y(y) = \begin{cases} \frac{e^{-y}}{b-a}, & y \in [-\ln b, -\ln a] \\ 0, & \text{altfel.} \end{cases}$$

△

### Generarea de numere pseudo-aleatoare ce urmează o distribuție discretă dată (metoda inversei)

Se dau  $(x_1, \dots, x_n)$  (valorile) și  $(p_1, \dots, p_n)$  (probabilitățile lor). Realizați un program care generează  $N$  numere pseudo-aleatoare, care urmează *distribuția discretă*

$$X \sim \begin{pmatrix} x_1 & x_2 & \dots & x_n \\ p_1 & p_2 & \dots & p_n \end{pmatrix},$$

folosind numere aleatoare uniform distribuite pe  $[0,1]$ .

*Procedeul de generare al numerelor aleatoare  $Y(i)$ ,  $i = \overline{1, N}$ , este:*

- Se citesc valorile  $x_1, x_2, \dots, x_n$  și probabilitățile corespunzătoare  $p_1, p_2, \dots, p_n$ , precum și numărul  $N$ . Fie  $p_0 = 0$ .
- Se generează  $N$  numere aleatoare uniform distribuite pe  $[0,1]$ :  $U(i)$ ,  $i = \overline{1, N}$ .
- Pentru fiecare  $i = \overline{1, N}$ :  $Y(i) = x_k$  dacă și numai dacă

$$p_0 + p_1 + \dots + p_{k-1} < U(i) \leq p_0 + p_1 + \dots + p_k, \quad k \in \{1, \dots, n\}.$$

- Se returnează:  $Y(i)$ ,  $i = \overline{1, N}$ .

Verificarea procedurii: deoarece  $U$  urmează legea uniformă, avem pe baza procedurii de mai sus:

$P(\text{"se generează } x_k") = P(p_0 + p_1 + \dots + p_{k-1} < U \leq p_0 + p_1 + \dots + p_k) = p_k$ ,  $k = 1, \dots, n$ , deci numerele generate urmează legea de distribuție discretă dată.

**Problemă:** Conform statisticilor medicale 46% din oameni au grupa sanguină **O**, 40% au grupa sanguină **A**, 10% au grupa sanguină **B** și 4% au grupa sanguină **AB**. Simulați de  $N(= 100, 1000)$  ori stabilirea grupei sanguine a unei persoane alese aleator și afișați frecvența de apariție a fiecărei grupe sanguine. Comparați rezultatele obținute cu cele teoretice.

```

function Y=ivtdiscret(x,p,N) %inverse transform method
Y=zeros(1,N);
q=cumsum(p);
for i=1:N
    U=unifrnd(0,1);
    Y(i)=x(find(U<=q,1));
end
%Aplicatia
clc
clear all
close all
N=1000;%numarul de simulari
x=[0,1,2,3];% valorile variabilei aleatoare
p=[0.46,0.4,0.1,0.04];%probabilitatile
Y=ivtdiscret(x,p,N);%generarea celor N numere
ny=hist(Y,length(unique(Y)));%determinarea frecventei de aparitie
[p' ny'/N]

```

## Vector aleator continuu

**Def. 21.**  $(X_1, \dots, X_n)$  este un **vector aleator continuu** dacă fiecare componentă a sa este o variabilă aleatoare continuă.

**Def. 22.**  $f_{(X,Y)} : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}_+$  este **funcția de densitate a vectorului aleator continuu**  $(X, Y)$ , dacă

$$P(X \leq x, Y \leq y) = \int_{-\infty}^x \int_{-\infty}^y f_{(X,Y)}(s, t) ds dt \quad \forall x, y \in \mathbb{R}.$$

**Def. 23.**  $F_{(X,Y)} : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}_+$  este **funcția de repartiție a vectorului aleator**  $(X, Y)$  (discret sau continuu), dacă

$$F_{(X,Y)}(x, y) = P(X \leq x, Y \leq y) \quad \forall x, y \in \mathbb{R}.$$

**Exemplu:** Vectorul aleator  $(X_1, X_2)$  este dat prin următorul tabel de contingență:

$X_1 \backslash X_2$	0	3
-2	0.4	0.3
4	0.2	0.1

$\Rightarrow (X_1, X_2)$  are funcția de repartiție  $F_{(X_1, X_2)} : \mathbb{R} \times \mathbb{R} \rightarrow [0, 1]$

%Matlab

clear all; clf; hold on; grid on;

xl= [-3, 6]; y=[-1, 4];  
view(35, 30)

x=-3:0.1:6; y=-1:0.15:4;

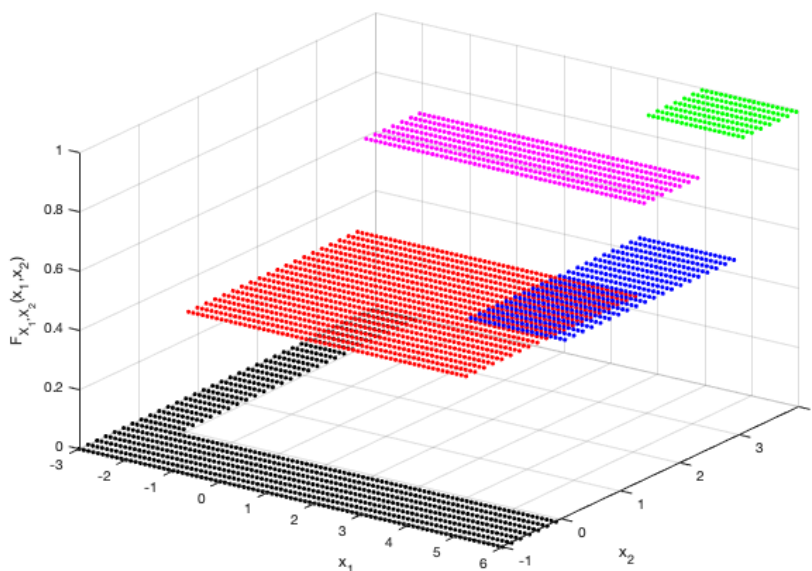
for i=1:length(x)

for j=1:length(y)

if x(i)<-2 || y(j)<0

f(i, j)=0;

$$f(x_1, x_2) = P(X_1 \leq x_1, X_2 \leq x_2) = \begin{cases} 0, & \text{dacă } x_1 < -2 \text{ sau } x_2 < 0 \\ 0.4, & \text{dacă } -2 \leq x_1 < 4 \text{ și } 0 \leq x_2 < 3 \\ 0.7, & \text{dacă } -2 \leq x_1 < 4 \text{ și } 3 \leq x_2 \\ 0.6, & \text{dacă } 4 \leq x_1 \text{ și } 0 \leq x_2 < 3 \\ 1, & \text{dacă } 4 \leq x_1 \text{ și } 3 \leq x_2. \end{cases}$$



Funcția de repartiție  $F_{(X_1, X_2)}$

```

    plot3(x(i),y(j),f(i,j),'k.')
end
if -2<=x(i) && x(i)<4 && 0<=y(j)&& y(j)<3
    f(i,j)=0.4;
    plot3(x(i),y(j),f(i,j),'r.')
end
if -2<=x(i) && x(i)<4 && 3<=y(j)
    f(i,j)=0.7;
    plot3(x(i),y(j),f(i,j),'m.')
end
if 4<=x(i) && 0<=y(j)&& y(j)<3
    f(i,j)=0.6;
    plot3(x(i),y(j),f(i,j),'b.')
end
if 4<=x(i) && 3<=y(j)
    f(i,j)=1;
    plot3(x(i),y(j),f(i,j),'g.')
end
end
end
xlabel('x_1'); ylabel('x_2'); zlabel('F_{X_1,X_2}(x_1,x_2)');

```



**Observație:** Cunoscând distribuția vectorului aleator continuu  $(X, Y)$ , cum determinăm distribuția

v.a. continue  $X$  respectiv  $Y$ ?

▷ dacă se cunoaște  $f_{(X,Y)}$ :  $f_X$ , respectiv  $f_Y$ , se determină cu

$$f_X(x) = \int_{-\infty}^{\infty} f_{(X,Y)}(x, y) dy, \quad \forall x \in \mathbb{R}, \quad f_Y(y) = \int_{-\infty}^{\infty} f_{(X,Y)}(x, y) dx, \quad \forall y \in \mathbb{R};$$

▷ dacă se cunoaște  $F_{(X,Y)}$ :  $F_X$ , respectiv  $F_Y$ , se determină cu

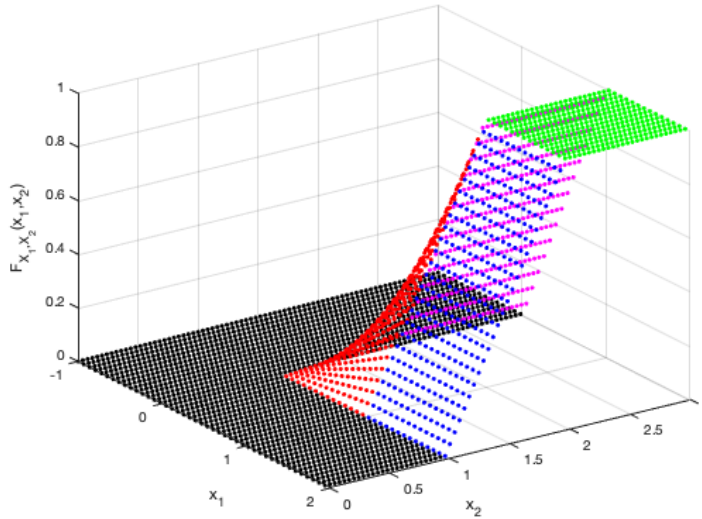
$$F_X(x) = \lim_{y \rightarrow \infty} F_{(X,Y)}(x, y), \quad F_Y(y) = \lim_{x \rightarrow \infty} F_{(X,Y)}(x, y).$$

**Exemplu:** Funcția de repartiție a vectorului aleator  $(X_1, X_2)$  este  $F_{(X_1, X_2)} : \mathbb{R} \times \mathbb{R} \rightarrow [0, 1]$

$$F_{(X_1, X_2)}(x_1, x_2) = \begin{cases} 0, & \text{dacă } x_1 < 0 \text{ sau } x_2 < 1 \\ x_1(x_2 - 1), & \text{dacă } 0 \leq x_1 < 1 \text{ și } 1 \leq x_2 < 2 \\ x_1, & \text{dacă } 0 \leq x_1 < 1 \text{ și } 2 \leq x_2 \\ x_2 - 1, & \text{dacă } 1 \leq x_1 \text{ și } 1 \leq x_2 < 2 \\ 1, & \text{dacă } 1 \leq x_1 \text{ și } 2 \leq x_2. \end{cases}$$

Ce distribuție au  $X_1$ , respectiv  $X_2$ ?

R.: Se determină  $F_{X_1}, F_{X_2}$  și  $f_{X_1}, f_{X_2} \implies X_1 \sim Unif[0, 1], X_2 \sim Unif[1, 2]$ .



Funcția de repartiție  $F_{(X_1, X_2)}$

**Def. 24.**  $X_1, \dots, X_n$  sunt ***n* variabilele aleatoare independente** (discrete sau continue), dacă

$$P(X_1 \leq x_1, \dots, X_n \leq x_n) = P(X_1 \leq x_1) \cdot \dots \cdot P(X_n \leq x_n) \quad \forall x_1, \dots, x_n \in \mathbb{R}.$$

**Observație** ( $n = 2$  în definiția de mai sus):  $X$  și  $Y$  sunt **două variabilele aleatoare independente**, dacă

$$P(X \leq x, Y \leq y) = P(X \leq x) \cdot P(Y \leq y) \quad \forall x, y \in \mathbb{R},$$

adică

$$F_{(X,Y)}(x, y) = F_X(x) \cdot F_Y(y) \quad \forall x, y \in \mathbb{R}.$$

**P. 16.** Variabilele aleatoare continue  $X$  (cu funcția de densitate  $f_X$ ) și  $Y$  (cu funcția de densitate  $f_Y$ ) sunt independente, dacă și numai dacă

$$f_X(x)f_Y(y) = f_{(X,Y)}(x,y) \quad \forall (x,y) \in \mathbb{R}^2,$$

unde  $f_{(X,Y)}$  este funcția de densitate a vectorului aleator  $(X, Y)$ .

**Exemplu:**  $(X_1, X_2)$  are distribuție uniformă pe  $I = [a_1, b_1] \times [a_2, b_2]$ , cu  $a_1, a_2, b_1, b_2 \in \mathbb{R}$ ,  $a_1 < b_1, a_2 < b_2$  dacă

$$f(x_1, x_2) = \begin{cases} \frac{1}{(b_1 - a_1)(b_2 - a_2)} & \text{dacă } (x_1, x_2) \in I \\ 0 & \text{dacă } (x_1, x_2) \notin I. \end{cases}$$

Funcția de repartiție are expresia

$$F(x_1, x_2) = \int_{-\infty}^{x_1} \int_{-\infty}^{x_2} f(t_1, t_2) dt_1 dt_2 = \left( \frac{x_1 - a_1}{b_1 - a_1} \right)_* \cdot \left( \frac{x_2 - a_2}{b_2 - a_2} \right)_*,$$

unde

$$u_* = \begin{cases} 0 & \text{dacă } u < 0 \\ u & \text{dacă } 0 \leq u \leq 1 \\ 1 & \text{dacă } 1 < u. \end{cases}$$



**P. 17.** Pentru un vector aleator continuu  $(X, Y)$  au loc proprietățile:

1.  $\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{(X,Y)}(u, v) du dv = 1.$
2.  $F_{(X,Y)}$  este funcție continuă.
3. Dacă  $F_{(X,Y)}$  este derivabilă parțial în  $(x, y)$ , atunci are loc:

$$\frac{\partial^2 F_{(X,Y)}(x, y)}{\partial x \partial y} = f_{(X,Y)}(x, y).$$

4.  $P((X, Y) \in A) = \underbrace{\int \int_A}_{A} f_{(X,Y)}(u, v) du dv, \quad A \subset \mathbb{R}^2 \text{ (măsurabilă)}.$

**Exemplu:** Fie  $(X, Y)$  vector aleator continuu, având funcția de repartiție

$$F_{(X,Y)}(x, y) = \begin{cases} (1 - e^{-x})(1 - e^{-2y}) & \text{dacă } x > 0 \text{ și } y > 0 \\ 0 & \text{în rest} \end{cases}$$

Sunt  $X$  și  $Y$  v.a. independente? Să se calculeze  $P(1 \leq X \leq 2 \leq Y \leq 3)$ .

R.: Se calculează  $F_X(x) = 1 - e^{-x}$  pentru  $x > 0$  și  $F_X(x) = 0$  pentru  $x \leq 0$ , precum și  $F_Y(y) = 1 - e^{-2y}$  pentru  $y > 0$  și  $F_Y(y) = 0$  pentru  $y \leq 0$ . Se verifică

$$F_{(X,Y)}(x, y) = F_X(x) \cdot F_Y(y) \quad \forall x, y \in \mathbb{R}.$$

Deci,  $X$  și  $Y$  sunt v.a. independente.

$$P(1 \leq X \leq 2 \leq Y \leq 3) = \int_1^2 \int_2^3 f_X(u) f_Y(v) du dv = (e^{-1} - e^{-2})(e^{-4} - e^{-6}) \approx 0.00368.$$

♡

### Operații cu v.a. continue

**Proprietate:** Fie  $(X, Y)$  vector aleator continuu cu funcția de densitate  $f_{(X,Y)}$ . Atunci  $X + Y$  și  $X \cdot Y$  sunt v.a. continue, având funcțiile de densitate:

$$f_{X+Y}(z) = \int_{-\infty}^{\infty} f_{(X,Y)}(u, z-u) du \quad \forall z \in \mathbb{R},$$

$$f_{X \cdot Y}(z) = \int_{-\infty}^{\infty} \frac{1}{|u|} f_{(X,Y)}\left(u, \frac{z}{u}\right) du \quad \forall z \in \mathbb{R}.$$

Dacă  $X$  și  $Y$  sunt v.a. independente, atunci

$$f_{X+Y}(z) = \int_{-\infty}^{\infty} f_X(u) f_Y(z-u) du \quad \forall z \in \mathbb{R},$$

$$f_{X \cdot Y}(z) = \int_{-\infty}^{\infty} \frac{1}{|u|} f_X(u) f_Y\left(\frac{z}{u}\right) du \quad \forall z \in \mathbb{R}.$$

**Def. 25.** Valoarea medie a unei v.a. continue  $X$ , care are funcția de densitate  $f$ , este

$$E(X) = \int_{-\infty}^{\infty} t f(t) dt$$

$$\text{dacă } \int_{-\infty}^{\infty} |t| f(t) dt < \infty.$$

▷ Valoarea medie a unei variabile aleatoare caracterizează tendința centrală a valorilor acesteia.

**P. 18.** *Proprietăți ale valorii medii; fie  $X, Y$  v.a. continue:*

→  $E(aX + b) = aE(X) + b$  for all  $a, b \in \mathbb{R}$ ;

→  $E(X + Y) = E(X) + E(Y)$ ;

→ Dacă  $X$  și  $Y$  sunt variabile aleatoare independente, atunci  $E(X \cdot Y) = E(X)E(Y)$ .

→ Dacă  $g : \mathbb{R} \rightarrow \mathbb{R}$  e o funcție, astfel încât  $g(X)$  este o v.a. continuă, atunci

$$E(g(X)) = \int_{-\infty}^{\infty} g(x)f_X(x)dx,$$

dacă  $\int_{-\infty}^{\infty} |g(x)|f_X(x)dx < \infty$ .

**Exemplu:** Durata drumului parcurs de un elev dimineața de acasă până la școală este o v.a. uniform distribuită între 20 și 26 minute. Dacă elevul pornește la 7:35 (a.m.) de acasă și are ore de la 8 (a.m.), care este probabilitatea ca elevul să ajungă la timp la școală? În medie cât durează drumul elevului până la școală?

Răspuns: fie  $X$  (v.a.) = durata drumului parcurs până la școală (în minute)  $\Rightarrow X \sim Unif[20, 26]$

$$P(\text{“elevul ajunge la timp la școală”}) = P(X \leq 25) = \frac{25 - 20}{26 - 20} = \frac{5}{6}.$$

$$E(X) = \int_{20}^{26} x \frac{1}{26 - 20} dx = \frac{20 + 26}{2} = 23 \text{ (minute).}$$



**Def. 26.** *Varianța (dispersia) unei variabile aleatoare  $X$  (discrete sau continue) este*

$$V(X) = E\left((X - E(X))^2\right),$$

(dacă valoarea medie  $E\left((X - E(X))^2\right)$  există). Valoarea  $\sqrt{V(X)}$  se numește **deviația standard** a lui  $X$ .

▷ Varianța unei variabile aleatoare caracterizează împrăștierea (dispersia) valorilor lui  $X$  în jurul valorii medii  $E(X)$ .

**P. 19.** *Proprietăți ale varianței:*

→  $V(X) = E(X^2) - E(X)^2$ .

→  $V(aX + b) = a^2V(X) \forall a, b \in \mathbb{R}$ .

→ Dacă  $X$  și  $Y$  sunt variabile aleatoare independente, atunci  $V(X + Y) = V(X) + V(Y)$ .

**Exemplu:** 1)  $X \sim \text{Bino}(n, p) \implies E(X) = np, V(X) = np(1 - p)$ .

2) Dacă  $X \sim N(m, \sigma^2) \implies E(X) = m, V(X) = \sigma^2$ .



► Matlab/Octave: *mean, var, std*

**Def. 27.**  $(X_n)_n$  este **șir de v.a. independente**, dacă  $\forall \{i_1, \dots, i_k\} \subset \mathbb{N}$  v.a.  $X_{i_1}, \dots, X_{i_k}$  sunt independente, adică

$$P(X_{i_1} \leq x_{i_1}, \dots, X_{i_k} \leq x_{i_k}) = P(X_{i_1} \leq x_{i_1}) \cdots P(X_{i_k} \leq x_{i_k})$$

$\forall x_{i_1}, \dots, x_{i_k} \in \mathbb{R}$ .

**Exemplu:** a)  $X_n =$  v.a. care indică numărul apărut la a  $n$ -aruncare a unui zar  $\implies (X_n)_n$  șir de v.a. independente

b) Se aruncă o monedă

$$X_n = \begin{cases} 0 & : \text{la a } n\text{-a aruncare a apărut cap,} \\ 1 & : \text{la a } n\text{-a aruncare a apărut pajură.} \end{cases}$$

$\implies (X_n)_n$  șir de v.a. independente.

**Def. 28.** Șirul de v.a.  $(X_n)_n$  **converge aproape sigur** la v.a.  $X$ , dacă

$$P(\omega \in \Omega : \lim_{n \rightarrow \infty} X_n(\omega) = X(\omega)) = 1.$$

**Notăție:**  $X_n \xrightarrow{\text{a.s.}} X$

► Cu alte cuvinte, convergența a.s.  $X_n \xrightarrow{\text{a.s.}} X$  impune ca  $(X_n(\omega))_n$  să convergă la  $X(\omega)$  pentru fiecare  $\omega \in \Omega$ , cu excepția unei mulțimi “mici” de probabilitate nulă; dacă  $X_n \xrightarrow{\text{a.s.}} X$  atunci evenimentul

$$M = \{\omega \in \Omega : (X_n(\omega))_n \text{ nu converge la } X(\omega)\} \text{ are } P(M) = 0.$$

**Exemplu:**

$$X_n \sim \begin{pmatrix} -\frac{1}{n} & \frac{1}{n} \\ 0.5 & 0.5 \end{pmatrix} \implies X_n \xrightarrow{\text{a.s.}} 0.$$

**Exemplu:**

$\Omega := [0, 1]$  spațiul de selecție, fie  $P$  probabilitatea pe  $[0, 1]$  indusă de măsura Lebesgue pe  $[0, 1]$ , adică pentru  $\forall \alpha < \beta$  din  $[0, 1]$  are loc

$$P([\alpha, \beta]) = P([\alpha, \beta)) = P((\alpha, \beta]) = P((\alpha, \beta)) := \beta - \alpha \text{ (lungimea intervalului)}$$

1)  $X_n(\omega) = \omega + \omega^n, \omega \in [0, 1], n \geq 1 \implies X_n \xrightarrow{\text{a.s.}} ???$

R.:

$$\lim_{n \rightarrow \infty} X_n(\omega) = \begin{cases} \omega & \text{pentru } \omega \in [0, 1) \\ 2 & \text{pentru } \omega = 1. \end{cases}$$



Fie  $X(\omega) = \omega$  pentru fiecare  $\omega \in \Omega$

$$\begin{aligned} &\Rightarrow \{\omega \in \Omega : \lim_{n \rightarrow \infty} X_n(\omega) = \omega\} = [0, 1) \\ &\Rightarrow P(\omega \in \Omega : \lim_{n \rightarrow \infty} X_n(\omega) = \omega) = P([0, 1)) = 1. \\ &X_n \xrightarrow{a.s.} X. \end{aligned}$$

2)  $X_n(\omega) = (-1)^n \omega$ ,  $\omega \in [0, 1]$ ,  $n \geq 1$ ; converge  $(X_n)_n$  a.s.?

R.:  $(X_n)_n$  nu converge a.s. spre o v.a.; șirul  $(X_n(\omega))_n$  este convergent doar în  $\omega = 0$ . ▲

**Frecvențe relative și absolute** (a se vedea Def.2): Fie  $A$  un eveniment asociat unei experiențe, repetăm experiența de  $n$  ori (în aceleași condiții date) și notăm cu  $r_n$  numărul de realizări ale evenimentului  $A$ ; **frecvența relativă** a evenimentului  $A$  este numărul

$$f_n(A) = \frac{r_n(A)}{n}$$

$r_n(A)$  este **frecvența absolută** a evenimentului  $A$ .

Experiment: Se aruncă o monedă de  $n$  ori;  $A$ : se obține *pajură*

$n$	frecvență absolută $r_n(A)$	frecvență relativă $f_n(A)$
100	48	0.48
1000	497	0.497
10000	5005	0.5005

$$f_n(A) \xrightarrow{a.s.} \frac{1}{2} \text{ (a se vedea P.21, LTNM)}$$

### Legea tare a numerelor mari (LTNM)

*Legea numerelor mari* se referă la descrierea rezultatelor unui experiment repetat de foarte multe ori. Conform acestei legi, rezultatul mediu obținut se apropie tot mai mult de valoarea așteptată, cu cât experimentul se repetă de mai multe ori. Aceasta se explică prin faptul că abaterile aleatoare se compensează reciproc.

Legea numerelor mari are două formulări: **legea slabă a numerelor mari (LSNM)** și **legea tare a numerelor mari (LTNM)**.

△ **Scurt istoric:** Jacob Bernoulli (1655 -1705) a formulat LSNM pentru frecvența relativă a unui experiment și a dat răspunsul la întrebarea “*Putem aproxima empiric probabilitățile?*” (în opera publicată postum, în 1713, *Ars conjectandi*):

▷ Teorema lui Bernoulli afirmă: “*Frecvențele relative converg în probabilitate la probabilitatea teoretică.*”

▷ În cadrul unui experiment poate să apară evenimentul  $A$  (*succes*) sau  $\bar{A}$  (*insucces*).

•  $X_i = 0 \Leftrightarrow$  dacă  $\bar{A}$  apare în a  $i$ -a repetiție a experimentului

- $X_i = 1 \Leftrightarrow$  dacă  $A$  apare în a  $i$ -a repetiție a experimentului  
 $\Rightarrow X_i \sim \text{Bernoulli}(p)$  cu  $p := P(A)$

$$X_i \sim \begin{pmatrix} 0 & 1 \\ 1 - P(A) & P(A) \end{pmatrix}$$

- $X_1, \dots, X_n$  sunt v.a. independente
- frecvența relativă de apariție a lui  $A$  este

$$f_n(A) = \frac{1}{n}(X_1 + \dots + X_n); \quad f_n(A) \text{ este o v.a.}$$

- $(X_n)_n$  verifică LSNM, adică

$$\lim_{n \rightarrow \infty} P\left(\left|f_n(A) - P(A)\right| > \varepsilon\right) = 0 \quad \forall \varepsilon > 0.$$

△



**Fig. 5.** Jacob Bernoulli (timbru emis în 1994 cu ocazia Congresului Internațional al Matematicienilor din Elveția)

**Def. 29.** Șirul de v.a.  $(X_n)_n$  cu  $E|X_n| < \infty \quad \forall n \in \mathbb{N}$  verifică **legea tare a numerelor mari (LTNM)** dacă

$$P\left(\left\{\omega \in \Omega : \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n (X_k(\omega) - E(X_k)) = 0\right\}\right) = 1,$$

adică

$$\frac{1}{n} \sum_{k=1}^n (X_k - E(X_k)) \xrightarrow{a.s.} 0.$$

**P. 20.** Fie  $(X_n)_n$  șir de v.a. independente având aceeași distribuție (și există  $m = E(X_n) \quad \forall n \in \mathbb{N}$ )  $\Rightarrow (X_n)_n$  verifică **LTNM**, adică

$$\frac{1}{n}(X_1 + \dots + X_n) \xrightarrow{a.s.} m.$$

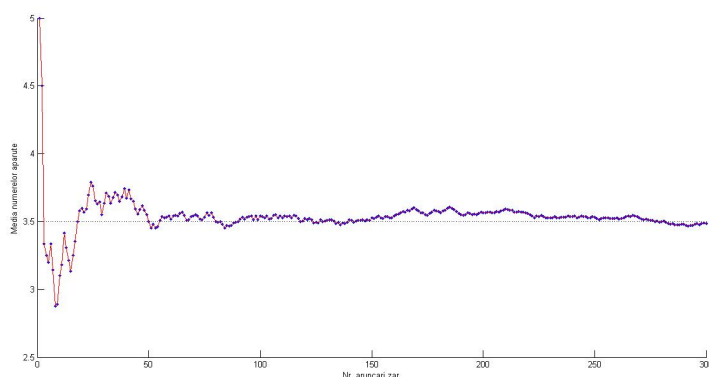
**În simulări:**  $\frac{1}{n}(X_1 + \dots + X_n) \approx m$ , dacă  $n$  este suficient de mare.

**Aplicație Matlab/Octave:** Simulare LTNM

```

clear all
clf
hold on
n=300;
x=unidrnd(6,1,n);
for i=1:n
    s(i)=sum(x(1:i))/i;
    y(i)=i;
    plot(y(i),s(i),'b.')
    plot(y(i),3.5,'k-')
end
plot(y,s,'r-')
xlabel('Nr. aruncari zar')
ylabel('Media numerelor aparute')

```



**Fig. 4. Simulare LTNM**

**P. 21.** Fie  $A$  un eveniment asociat unei experiențe, repetăm experiența de  $n$  ori (în aceleași condiții date). LTNM: cu cât repetăm mai des un experiment ( $n \rightarrow \infty$ ), cu atât mai bine aproximează frecvența relativă  $f_n(A)$  a evenimentului  $A$  probabilitatea sa (teoretică) de apariție  $P(A)$ :

$$f_n(A) \xrightarrow{a.s.} P(A), \text{ dacă } n \rightarrow \infty.$$

**În simulări:**  $f_n(A) \approx P(A)$ , dacă  $n$  este suficient de mare.

## Statistică matematică

► Statistica matematică este o ramură a matematicii aplicate, care se ocupă de *colectarea, gruparea, analiza și interpretarea datelor* referitoare la anumite fenomene în scopul obținerii unor previziuni;

- statistica descriptivă: metode de colectare, organizare, sintetizare, prezentare și descriere a datelor numerice (sau nenumere) într-o formă convenabilă
- statistica inferențială: metode de interpretare a rezultatelor obținute prin metodele statisticii

descriptive, utilizate apoi pentru luarea deciziilor.

► O *colectivitate* sau *populație statistică*  $\mathcal{C}$  este o mulțime de elemente care au anumite însușiri comune ce fac obiectul analizei statistice. Numărul elementelor populației se numește *volumul populației*.

Exemple de populații statistice: mulțimea persoanelor dintr-o anumită țară, localitate, zonă etc. într-un anumit an; mulțimea gospodăriilor din România la un moment dat; mulțimea consumatorilor unui anumit produs; mulțimea societăților care produc un anumit produs; angajații unei societăți; studenții unei facultăți.

► *Eșantionul*  $\mathcal{E}$  reprezintă o submulțime a unei populații statistice  $\mathcal{E} \subset \mathcal{C}$ , constituită după criterii bine stabilite:

- a) să fie aleatoare;
- b) toate elementele colectivității să aibe aceeași șansă de a fi alese în eșantion;
- c) eșantionul să fie reprezentativ (structura eșantionului să fie apropiată de structura populației);
- d) volumul eșantionului să fie suficient de mare.

► *Unitatea statistică* (indivizii) este elementul, entitatea de sine stătătoare a unei populații statistice, care posedă o serie de trăsături caracteristice ce-i conferă apartenența la populația studiată.

De exemplu: *unitatea statistică simplă*: un salariat, un student, un agent economic, o trăsătură, o părere; *unitatea statistică complexă*: o grupă de studenți sau o echipă de salariați, o familie sau o gospodărie, o categorie de mărfuri.

► *Variabila statistică* sau *caracteristica* reprezintă o însușire, o proprietate măsurabilă a unei unități statistice, întâlnită la toate unitățile care aparțin aceleiași colectivități și care prezintă variabilitate de la o unitate statistică la alta. Caracteristica sau variabila statistică corespunde unei variabile aleatoare.

Exemple de caracteristici: vârsta, salariul, preferințele politice, prețul unui produs, calitatea unor servicii, nivelul de studii.

a) variabile (caracteristici) continue → iau un număr infinit și nenumărabil de valori într-un interval sau reuniune de intervale (de ex.: greutatea, înălțimea, valoarea glicemiei, temperatura aerului)

b) variabile (caracteristici) discrete → iau număr finit sau infinit dar numărabil de valori discrete (de ex.: numărul elevi ai unei școli, numărul liceelor existente într-un oraș, valoarea IQ)

▷ caracteristicile de la a) și b) sunt variabile numerice (cantitative)

c) variabile (caracteristici) nominale (de ex.: culoarea ochilor, ramura de activitate, religia)

d) variabile (caracteristici) nominale ordinale (de ex.: starea de sănătate / calitatea unor servicii - precară, mai bună, bună, foarte bună)

e) variabile (caracteristici) dihotomiale (binare) (de ex.: stagiul militar - satisfăcut/nesatisfăcut, starea civilă - căsătorit/necăsătorit)

▷ caracteristicile de la c), d), e) sunt variabile calitative

▷ variabilele nominale mai sunt numite variabile categoriale

► *Datele statistice* reprezintă observațiile rezultate dintr-o cercetare statistică, sau ansamblul valorilor colectate în urma unei cercetări statistice.

De exemplu: un angajat al unei companii are o vechime de 6 ani în muncă. Angajatul reprezintă unitatea statistică, vechimea în muncă este caracteristica (variabila) cercetată, iar 6 este valoarea acestei caracteristici.

O *colectivitate* (populație)  $\mathcal{C}$  este cercetată din punctul de vedere al caracteristicii (variabilei statistice)  $X$ .

Distribuția caracteristicii  $X$  poate fi

1) complet specificată (de ex.:  $X \sim \text{Exp}(3)$ ,  $X \sim \text{Bin}(10, 0.3)$ ,  $X \sim N(0, 1)$ )

2) specificată, dar depinzând de unul sau mai mulți parametri necunoscuți (de ex.:  $X \sim \text{Exp}(\lambda)$ ,  $X \sim \text{Bin}(10, p)$ ,  $X \sim N(m, \sigma^2)$ )

3) necunoscută:  $X \sim ?$

• în cazul 2) parametrii sunt necunoscuți, iar în cazul 3) distribuția este necunoscută

→ se estimează → teoria estimăției

→ se testează → teste statistice

► Fie  $\mathcal{E} \subset \mathcal{C}$  un eșantion. Se numesc **date de selecție** relative la caracteristica  $X$  datele statistice  $x_1, \dots, x_n$  obținute prin cercetarea indivizilor care fac parte din eșantionul  $\mathcal{E}$ .

► Datele de selecție  $x_1, \dots, x_n$  pot fi considerate ca fiind valorile unor variabile aleatoare  $X_1, \dots, X_n$ , numite **variabile de selecție** și care se consideră a fi variabile aleatoare independente și având aceeași distribuție ca  $X$ .

► Fie  $x_1, \dots, x_n$  datele statistice pentru caracteristica cercetată  $X$ , notăm cu  $X_1, \dots, X_n$  variabilele de selecție corespunzătoare. Fie  $g : \mathbb{R}^n \rightarrow \mathbb{R}$  o funcție astfel încât  $g(X_1, \dots, X_n)$  este o variabilă aleatoare.

$g(X_1, \dots, X_n)$  se numește **funcție de selecție** sau **estimator**

$g(X_1, \dots, X_n)$  se numește valoarea funcției de selecție sau **valoarea estimatorului**.

• **Exemple de estimatori (funcții de selecție)** sunt: media de selecție, dispersia de selecție, funcția de repartiție empirică (de selecție).

▷ Estimatorii (funcțiile de selecție) se folosesc în statistică pentru estimarea punctuală a unor parametri necunoscuți, pentru obținerea unor intervale de încredere pentru parametri necunoscuți, pentru verificarea unor ipoteze statistice.

Fie  $x_1, \dots, x_n$  datele statistice pentru caracteristica cercetată  $X$ , notăm cu  $X_1, \dots, X_n$  variabilele de selecție corespunzătoare:

► **media de selecție**

$$\bar{X}_n = \frac{1}{n} (X_1 + \dots + X_n)$$

► valoarea mediei de selecție

$$\bar{x}_n = \frac{1}{n} (x_1 + \cdots + x_n)$$

► varianța (dispersia) de selecție

$$\tilde{S}_n^2 = \frac{1}{n-1} \sum_{k=1}^n (X_k - \bar{X}_n)^2$$

► valoarea varianței (dispersiei) de selecție

$$\hat{s}_n^2 = \frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x}_n)^2$$

► abaterea standard de selecție

$$\tilde{S}_n = \left( \frac{1}{n-1} \sum_{k=1}^n (X_k - \bar{X}_n)^2 \right)^{\frac{1}{2}}$$

► valoarea abaterii standard de selecție

$$\hat{s}_n = \left( \frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x}_n)^2 \right)^{\frac{1}{2}}$$

► funcția de repartiție empirică  $\hat{F}_n : \mathbb{R} \times \Omega \rightarrow \mathbb{R}$

$$\hat{F}_n(x) = \frac{\#\{i \in \{1, \dots, n\} : X_i \leq x\}}{n}, x \in \mathbb{R}$$

► valoarea funcției de repartiție empirice

$$\hat{F}_n(x) = \frac{\#\{i \in \{1, \dots, n\} : x_i \leq x\}}{n}, x \in \mathbb{R}$$

**Def. 30.**  $g(X_1, \dots, X_n)$  este *estimator nedeplasat* pentru parametrul necunoscut  $\theta$ , dacă

$$E(g(X_1, \dots, X_n)) = \theta.$$

$g(X_1, \dots, X_n)$  este *estimator consistent* pentru parametrul necunoscut  $\theta$ , dacă

$$g(X_1, \dots, X_n) \xrightarrow{a.s.} \theta.$$

Fie  $g_1(X_1, \dots, X_n)$  și  $g_2(X_1, \dots, X_n)$  estimatori nedeplasați pentru parametrul necunoscut  $\theta$ .  
 $g_1(X_1, \dots, X_n)$  este *mai eficient* decât  $g_2(X_1, \dots, X_n)$ , dacă  $V(g_1) < V(g_2)$ .

**Observații:**

- 1) Media de selecție  $\bar{X}_n$  este estimator nedeplasat și consistent pentru media teoretică  $E(X)$  a caracteristicii  $X$ ; în practică  $E(X) \approx \bar{x}_n$ .
- 2) Varianța (dispersia) de selecție  $\tilde{S}_n^2$  este estimator nedeplasat și consistent pentru varianța teoretică  $V(X)$  a caracteristicii  $X$ ; în practică  $V(X) \approx \tilde{s}_n^2$ .
- 3) Funcția de repartiție de selecție calculată în  $x \in \mathbb{R}$ :  $\hat{F}_n(x)$  este estimator nedeplasat și consistent pentru  $F_X(x)$ , care este valoarea funcției de repartiție teoretice calculată în  $x$ ; în practică  $F_X(x) \approx \hat{F}_n(x)$ .

**Metoda momentelor pentru estimarea parametrilor necunoscuți  $\theta = (\theta_1, \dots, \theta_r)$  pentru distribuția caracteristicii cercetate  $X$**

de exemplu:

$X \sim \text{Exp}(\lambda)$  parametrul necunoscut:  $\theta = \lambda$

$X \sim N(m, \sigma^2)$  parametri necunoscuți:  $(\theta_1, \theta_2) = (m, \sigma)$

$X \sim \text{Unif}[a, b]$  parametri necunoscuți:  $(\theta_1, \theta_2) = (a, b)$

Fie  $x_1, \dots, x_n$  datele statistice pentru caracteristica cercetată  $X$  și fie  $X_1, \dots, X_n$  variabilele de selecție corespunzătoare.

Se rezolvă sistemul

$$\begin{cases} E(X^k) = \frac{1}{n} \sum_{i=1}^n X_i^k, \\ k = \{1, \dots, r\} \end{cases}$$

cu necunoscutele  $\theta_1, \dots, \theta_r$ .

Soluția sistemului  $\hat{\theta}_1, \dots, \hat{\theta}_r$  este estimatorul pentru parametrii necunoscuți ai distribuției caracteristicii  $X$ .

**Exemplu:** Folosind metoda momentelor, să se estimeze parametrul necunoscut  $\theta := a$  pentru  $X \sim \text{Unif}[0, a]$ ; se dau datele statistice: 0.1, 0.3, 0.9, 0.49, 0.12, 0.31, 0.98, 0.73, 0.13, 0.62.

Avem cazul:  $r = 1$ , calculăm  $E(X) = \frac{a}{2}$ ,  $n = 10$ ,  $\bar{x}_n = 0.468$ . Se rezolvă

$$E(X) = \frac{1}{n} \sum_{i=1}^n X_i \implies \frac{a}{2} = \frac{1}{n} \sum_{i=1}^n X_i.$$

Estimatorul pentru parametrul necunoscut  $a$  este

$$\hat{a}(X_1, \dots, X_n) = \frac{2}{n} \sum_{i=1}^n X_i,$$

unde  $X_1, \dots, X_n$  sunt variabilele de selecție. Valoarea estimatorului este

$$\hat{a}(x_1, \dots, x_n) = \frac{2}{n} \sum_{i=1}^n x_i = 0.936.$$

Parametrul necunoscut  $a$  este estimat cu valoarea 0.936.

► Este  $\hat{a}(X_1, \dots, X_n)$  un estimator nedeplasat pentru parametrul  $a$ ?