

# Curs 5: Distribuții clasice de probabilitate: Bernoulli, binomială, geometrică și Poisson.

## 1.1 Distribuția Bernoulli

Un experiment Bernoulli este un experiment aleator ce constă dintr-un număr precizat sau nu, de încercări independente. O încercare are doar două rezultate mutual exclusive, unul numit succes și celălalt eșec. Variabila aleatoare,  $X$ , ce înregistrează rezultatul unei încercări se numește variabilă aleatoare Bernoulli. Exemplul clasic de experiment Bernoulli este aruncarea monedei: succesul fiind, de exemplu stema, iar eșecul, banul. Înregistrând 1 dacă se produce un succes și 0 în caz de eșec, și notând cu  $p$  probabilitatea succesului, distribuția de probabilitate a variabilei aleatoare Bernoulli este:

$$X = \begin{pmatrix} 1 & 0 \\ p & 1 - p \end{pmatrix} \quad (1)$$

Valoarea medie a unei variabile aleatoare Bernoulli este

$$M(X) = 1 \cdot p + 0(1 - p) = p,$$

iar dispersia:

$$\sigma^2(X) = (1 - p)^2 p + (0 - p)^2 (1 - p) = (1 - p)^2 p - p(1 - p) = p(1 - p)$$

Distribuția Bernoulli se folosește în generarea de șiruri de biți aleatori, sau la alegerea unei alternative din două.

**Exemplul 1.** (Run-uri de biți) Având un șir de  $n$  biți rezultați din simularea de  $n$  ori a variabilei aleatoare  $X$  de mai sus, numim run de biți o succesiune de biți identici în șir. De exemplu în șirul de biți 0, 1, 1, 0, 0, 0, 1, 0, 1, 1, 1, 0, avem următoarele run-uri de biți 1:

$$11, \quad 1, \quad 111$$

În analiza șirurilor de biți aleatori folosiți în probleme de securitate (criptografie) este foarte util numărul mediu de run-uri de biți conținuți în șir.

Notăm cu  $N$  variabila aleatoare ce dă numărul de run-uri de biți 1, în șirul de  $n$  biți

$$b_1, b_2, \dots, b_n,$$

obținuți din simularea (observarea) variabilelor aleatoare independente  $X_1, X_2, \dots, X_n$ , Bernoulli distribuite, adică:

$$X_i = \begin{pmatrix} 1 & 0 \\ p & 1-p \end{pmatrix} \quad (2)$$

Practic simulăm aceeași variabilă aleatoare  $X$ , doar că îi asociem indicele  $1, 2, 3, \dots, n$  care indică al câtelea bit generăm.

Prin urmare fiecare bit  $b_i$  din șir este o "observație" asupra variabilei  $X_i$ ,  $i = \overline{1, n}$

Notăm cu  $Y_i$  variabila aleatoare ce ia valoarea 1 dacă în poziția  $i$  a șirului de biți începe un *run* de unu și 0 în caz contrar. Astfel variabila aleatoare  $N$ , este suma variabilelor  $Y_1, Y_2, \dots, Y_n$ ,  $N = \sum_{i=1}^n Y_i$ , iar numărul mediu de *run*-uri este  $M(N) = \sum_{i=1}^n M(Y_i)$ .

Dar conform definiției variabilelor  $Y_i$  avem că:

$$Y_1 = \begin{pmatrix} 1 & 0 \\ P(X_1 = 1) & P(X_1 = 0) \end{pmatrix}, \quad M(Y_1) = P(X_1 = 1) = p$$

adică un run de biți începe din poziția 1, dacă  $X_1 = 1$ . Pentru  $i \geq 2$ , un run de biți începe din poziția  $i$ , dacă bitul  $b_i = 1$ , dar bitul  $b_{i-1} = 0$ , adică pentru  $i = \overline{2, n}$ ,  $P(Y_i = 1) = P(X_{i-1} = 0, X_i = 1) = P(X_{i-1} = 0)P(X_i = 1) = (1-p)p$ . Notând  $\pi = (1-p)p$ , variabila aleatoare  $Y_i$  are distribuția de probabilitate:

$$Y_i = \begin{pmatrix} 1 & 0 \\ \pi & 1-\pi \end{pmatrix}$$

și media  $M(Y_i) = \pi = (1-p)p$ . Prin urmare numărul mediu de *run*-uri de 1, în șirul de biți  $b_1, b_2, \dots, b_n$ , este  $M(N) = p + (n-1)(1-p)p$ .

Analiza statistică a *run*-urilor dintr-un șir de biți se efectuează pentru a evalua caracterul aleator al acestora. Unul din testele folosite în criptografie este testul *run*, care numără *run*-urile de 0 și 1 dintr-un șir de biți produs de un generator de biți aleatori. Dacă numărul *run*-urilor de 1,2,3,4,5 biți și respectiv  $\geq 6$  biți, nu intră în intervale prescrise, șirul este considerat nesigur.

## 1.2 Distribuția binomială

O variabilă aleatoare binomială este asociată unui experiment Bernoulli ce constă din  $n$  încercări (repetiții). O încercare este o etapă a experimentului, ce are două rezultate mutual exclusive: unul numit *succes* și celălalt *eșec*.

- Încercările sunt independente, în sensul că rezultatul unei încercări nu influențează rezultatul celorlalte.

- Pentru fiecare încercare probabilitatea succesului este aceeași,  $p$ .

**Exemple de experimente Bernoulli:**

• Exemplul clasic este aruncarea de  $n$  ori a unei monede; succesul: cade banul, eșecul, cade stema;

• Observarea a  $n$  execuții consecutive a grupului de instrucțiuni:

```
if (B) execută I1;
else execută I2;
```

Succesul = "expresia B este adevărată", iar eșecul = "expresia B este falsă".

• Transmiterea unei succesiuni de  $n$  biți printr-un canal de comunicație. Succesul = "bit recepționat corect", eșec = "bit recepționat incorect".

O variabilă aleatoare  $X$  asociată unui experiment Bernoulli, ce înregistrează numărul succeselor din  $n$  încercări se numește *variabilă aleatoare binomială*.

Mulțimea valorilor unei variabile aleatoare de tip binomial este  $D = \{0, 1, 2, \dots, n\}$ , deoarece în  $n$  încercări se pot înregistra, 0 succese, un succes, ...,  $n$  succese.

Să determinăm distribuția de probabilitate a unei variabile aleatoare binomiale  $X$ , în funcție de numărul  $n$  al încercărilor și probabilitatea  $p$  a unui succes într-o încercare, adică să calculăm probabilitatea ca în  $n$  încercări să înregistrăm  $k$  succese:  $P(X = k)$ ,  $k \in \{0, 1, 2, \dots, n\}$ .

În acest scop remarcăm că rezultatul fiecărei încercări este o observație asupra unei variabile aleatoare Bernoulli. Notăm cu  $X_1, X_2, \dots, X_n$  variabilele aleatoare Bernoulli, identic distribuite:

$$X_i = \begin{pmatrix} 1 & 0 \\ p & 1-p \end{pmatrix}, \quad i = \overline{1, n}$$

ce sunt "observate" în încercările  $1, 2, \dots, n$ . Conform caracterizării experimentului Bernoulli, aceste variabile sunt independente și deci vectorul aleator  $(X_1, X_2, \dots, X_n)$  are distribuția de probabilitate dată de:

$$\begin{aligned} P((X_1, X_2, \dots, X_n) = (i_1, i_2, \dots, i_n)) &= P(X_1 = i_1, X_2 = i_2, \dots, X_n = i_n) \\ &= P(X_1 = i_1)P(X_2 = i_2) \cdots P(X_n = i_n) \\ &= p^k(1-p)^{n-k}, \end{aligned} \quad (3)$$

unde  $k$  este numărul de valori 1 în  $n$ -lista de biți  $(i_1, i_2, \dots, i_n)$ . Numărul de succese în  $n$  încercări este  $X = X_1 + X_2 + \cdots + X_n$ . Evident că suma variabilelor Bernoulli  $X_1 + X_2 + \cdots + X_n$  ia valori în mulțimea  $\{0, 1, 2, \dots, n\}$ . Probabilitatea ca în  $n$  încercări să avem  $k$  succese este

$$\begin{aligned} P(X_1 + X_2 + \cdots + X_n = k) &= P\left(\bigcup_{\substack{(i_1, i_2, \dots, i_n) \\ i_1 + i_2 + \cdots + i_n = k}} (X_1 = i_1, X_2 = i_2, \dots, X_n = i_n)\right) \\ &= \sum_{\substack{(i_1, i_2, \dots, i_n) \\ i_1 + i_2 + \cdots + i_n = k}} P(X_1 = i_1, X_2 = i_2, \dots, X_n = i_n), \end{aligned}$$

unde suma se calculează după toate  $n$ -listele de biți  $(i_1, i_2, \dots, i_n)$  pentru care suma biților este  $k$  (există  $C_n^k$  astfel de liste!). Deci,

$$P(X_1 + X_2 + \cdots + X_n = k) = C_n^k p^k (1-p)^{n-k}$$

Ca să înțelegeți mai bine luăm un exemplu simplu  $n = 3$  și  $k = 2$ . În acest caz avem:  
 $P(X_1 + X_2 + X_3 = 2) = P((X_1 = 1, X_2 = 1, X_3 = 0) \cup (X_1 = 1, X_2 = 0, X_3 = 1) \cup (X_1 = 0, X_2 = 1, X_3 = 1)) = C_3^2 p^2 (1 - p)$

Conform formulei binomului lui Newton avem:  $\sum_{k=0}^n C_n^k p^k (1-p)^{n-k} = (p + (1-p))^n = 1$ .

În concluzie o variabilă aleatoare binomială este o variabilă discretă având distribuția de probabilitate:

$$X = \left( C_n^k p^k (1-p)^{n-k} \right), \quad k = 0, 1, \dots, n \quad (4)$$

Numele de variabilă aleatoare binomială vine deci de la faptul că  $P(X = k)$  este termenul  $k + 1$  al dezvoltării binomiale  $(p + (1-p))^n$ .

Notăm cu  $\text{Bin}(n, p)$  clasa variabilelor aleatoare binomiale, asociate unor experimente Bernoulli cu  $n$  încercări și probabilitatea  $p$  a unui succes într-o încercare, iar  $X \sim \text{Bin}(n, p)$  indică că  $X$  este o variabilă aleatoare ce are distribuția de probabilitate binomială.  $n$  și  $p$  se numesc parametrii distribuției binomiale.

**Exemplul 2.** După ce un virus a pătruns în sistem, se verifică starea a 20 de fișiere. Știind că probabilitatea ca un fișier să fi fost afectat de virus este 0.2, independent de celelalte fișiere, care este probabilitatea ca cel mult cinci din cele 20 de fișiere să fi fost deteriorate ?

Variabila aleatoare  $X$  ce dă numărul fișierelor infectate din cele 20, are distribuția de probabilitate  $\text{Bin}(n = 20, p = 0.2)$ . Probabilitatea cerută este

$$P(X \leq 5) = \sum_{k=0}^5 C_{20}^k (0.2)^k (0.8)^{20-k}$$

**Exemplul 3.** Presupunem că există 10 sateliți GPS (Global Positioning System) pe orbită. O unitate GPS de la sol este activă, dacă cel puțin 4 sateliți GPS funcționează (pot fi contactați). Știind că încercările de contactare ale celor 10 sateliți sunt independente și că probabilitatea ca GPS-ul de la sol să eșueze în contactarea oricărui satelit din cei 10, este aceeași și egală cu  $p = 0.75$ , să se calculeze probabilitatea ca unitatea GPS de la sol să fie activă.

Remarcăm că încercările de contactare a celor  $n = 10$  sateliți GPS definesc un experiment Bernoulli, cu  $p = 0.75$  (în acest context succesul nu are sensul uzual, ci  $p = 0.75$  este probabilitatea ca GPS-ul de la sol să eșueze în contactarea unui satelit).

Variabila aleatoare  $X$ , ce dă numărul de sateliți ce nu pot fi contactați de la sol, are distribuția de probabilitate  $\text{Bin}(n = 10, p = 0.75)$ . Astfel probabilitatea ca unitatea GPS de la sol să funcționeze este probabilitatea evenimentului ( $X \leq 6$ ), adică probabilitatea ca cel mult 6 sateliți să nu poată fi contactați (deci cel puțin patru să poată fi contactați):

$$P(X \leq 6) = \sum_{k=0}^6 C_{10}^k (0.75)^k (0.25)^{10-k}$$

**Propoziția 1.2.1** Fie  $X$  o variabilă aleatoare binomială  $X \sim \text{Bin}(n, p)$ . Media lui  $X$  este  $M(X) = np$ , iar dispersia,  $\sigma^2(X) = np(1 - p)$ .

**Observația 1.2.1** Experimentul Bernoulli se asimilează cu experimentul extragerii succesive a câte unei bile dintr-o urnă cu bile albe și negre și returnarea bilei după extragere. Asimilăm extragerea unei bile albe cu succesul și a unei bile negre cu eșecul. Returnarea bilei după extragere, asigură independența încercărilor. Dacă bila nu se returnează, încercările nu mai sunt independente, deoarece conținutul urnei se modifică după fiecare extragere.

### 1.3 Distribuția geometrică

Distribuția geometrică se definește în contextul unui experiment Bernoulli în care numărul de încercări nu este fixat apriori. Presupunem că se efectuează o succesiune de încercări independente ce pot avea ca rezultat, succes sau eșec. Fie  $X_i$  rezultatul celei de-a  $i$ -a încercări.  $X_i$  ia valoarea 1, respectiv 0 după cum în a  $i$ -a încercare rezultatul este un succes, respectiv eșec. Notăm cu  $Y$  variabila aleatoare ce dă numărul de încercări până la primul succes, inclusiv.

- Mulțimea valorilor variabilei  $Y$  este  $\mathbb{N}^* = \{1, 2, \dots, n, \dots\}$ ;
- Să determinăm distribuția de probabilitate, adică să calculăm  $P(Y = k)$ ,  $k \in \mathbb{N}^*$ :  
Evenimentul  $(Y = k)$  se exprimă astfel

$$(Y = k) = (X_1 = 0, X_2 = 0, \dots, X_{k-1} = 0, X_k = 1)$$

Dacă notăm cu  $p$  probabilitatea succesului, atunci datorită independenței variabilelor  $X_1, X_2, \dots, X_k$ , avem:

$$\begin{aligned} P(Y = k) &= P(X_1 = 0)P(X_2 = 0) \cdots P(X_{k-1} = 0)P(X_k = 1) \\ &= (1 - p)^{k-1}p \end{aligned} \quad (5)$$

Să arătăm că suma seriei  $\sum_{k=1}^{\infty} (1 - p)^{k-1}p = 1$ . Notând  $q = 1 - p$ ,  $m = k - 1$  avem seria geometrică cu rația  $q$  ori  $p$ :

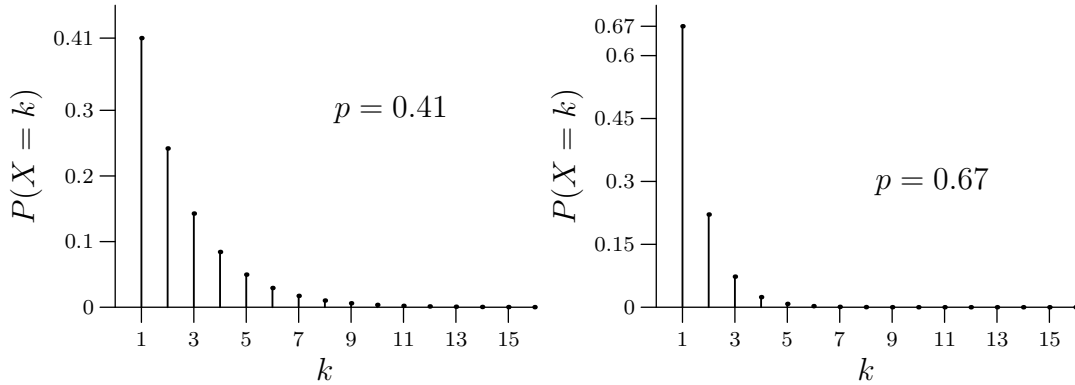
$$p \sum_{m=0}^{\infty} q^m = \frac{p}{1 - q} = 1$$

Deci, o variabilă aleatoare  $Y$  ce are distribuția geometrică, dă numărul de încercări într-un proces Bernoulli, până se obține primul succes, inclusiv.  $Y$  ia o mulțime numărabilă de valori.

Notăm cu  $\text{Geom}(p)$  clasa variabilelor aleatoare ce au distribuția geometrică de parametru  $p$ .

Dacă  $X \sim \text{Geom}(p)$ ,  $X$  are distribuția de probabilitate:

$$X = \left( \binom{k}{(1 - p)^{k-1}p} \right), \quad k \in \mathbb{N}^*, p \in (0, 1) \quad (6)$$



**Fig.1:** Ilustrarea distribuției de probabilitate geometrică pentru  $p = 0.41$ , respectiv  $p = 0.67$ .

Distribuția de probabilitate geometrică,  $\text{Geom}(p = 0.41)$ , și  $\text{Geom}(p = 0.67)$ , este ilustrată în Fig.1.

**Propoziția 1.3.1** Dacă  $X \sim \text{Geom}(p)$ , atunci valoarea medie și dispersia variabilei aleatoare,  $X$ , este:

$$M(X) = \frac{1}{p}, \text{ iar } \sigma^2(X) = \frac{1-p}{p^2}. \quad (7)$$

**Exemplul 4.** Considerăm segmentul de program:

```
do{
  bloc de instructiuni I;
}while{!B}
```

Presupunem că expresia booleană  $B$  ia valoarea **true** cu probabilitatea  $p$  și valoarea **false** cu probabilitatea  $1 - p$ . Dacă rezultatele testelor succesive asupra lui  $B$  sunt independente ce distribuție de probabilitate are numărul de execuții ale grupului de instrucțiuni  $I$ ?

Experimentul constă în parcurgerea a cel puțin o dată a blocului  $I$  și după fiecare parcurgere se înregistrează eșec dacă  $B$  este **false** (deci  $!B$  este **true**), respectiv succes în caz contrar. Variabila aleatoare  $N$  ce dă numărul de parcurgeri ale buclei **do-while** are distribuția geometrică de parametru  $p$ . Într-adevăr, dacă de  $k - 1$  ori consecutiv, rezultatul testului este  $!B$  este **true**, iar după a  $k$  execuție  $!B$  este **false**, atunci  $P(N = k) = (1 - p)^{k-1}p$ .

**Exemplul 5.** Un algoritm generează la întâmplare un întreg fără semn, pe  $k$  biți ( $k=32$ , de exemplu). Notăția  $s \leftarrow_R \{0, 1\}^k$  semnifică că algoritmul generează stringul de biți  $s$ . La întâmplare înseamnă că fiecare string de  $k$  biți poate fi generat cu aceeași probabilitate și rezultatul fiecărei generări este independent de restul generărilor.

Dacă  $s = (b_{k-1}b_{k-2} \dots b_1b_0)_2$ , este reprezentarea în baza 2 a unui întreg fără semn, atunci numărul zecimal asociat este  $I(s) = \sum_{i=0}^{k-1} b_i 2^i$ . Considerăm algoritmul probabilist:

```

do{
    s ←R {0, 1}k;
    n ← I(s);
}while (n ≥ M);
return n;

```

Să se determine distribuția de probabilitate a variabilei aleatoare  $N$ , ce ia ca valori numerele pe  $k$  biți returnate de algoritm și distribuția de probabilitate a variabilei,  $X$ , ce dă numărul de parcurgeri ale buclei **do-while**.

**Rezolvare:** Notăm cu  $Y$  variabila aleatoare ce are ca valori numerele întregi  $n$  ce pot fi generate de secvența:

```

s ←R {0, 1}k
n ← I(s)

```

Mulțimea valorilor lui  $Y$  în binar este  $\{\underbrace{(00 \dots 0)}_k, \dots, \underbrace{(11 \dots 1)}_k\}$ , iar în baza 10:

$$D = \{0, 1, 2 \dots 2^k - 1\}$$

Cardinalul lui  $D$  este  $|D| = 2^k$ . Deci probabilitatea de generare a unui număr întreg  $m$ , fără semn, pe  $k$  biți, este  $P(Y = m) = 1/2^k$ ,  $\forall m \in D$ . Variabila aleatoare  $N$  este variabila  $Y$  condiționată de evenimentul  $A = (Y < M)$ , notat  $N = (Y|Y < M)$ .

Evenimentul  $(Y < M) = (Y = 0) \cup (Y = 1) \cup \dots \cup (Y = M - 1)$  are probabilitatea  $P(Y < M) = P(Y = 0) + P(Y = 1) + \dots + P(Y = M - 1) = M/2^k$ . Variabila aleatoare  $N$  ia valorile  $\{0, 1, 2 \dots, M - 1\}$  și

$$P(N = i) = P(Y = i|Y < M) = \frac{P(Y = i, Y < M)}{P(Y < M)} = \frac{1/2^k}{M/2^k} = 1/M,$$

$\forall i = 0, 1, \dots, M - 1$ . Deci  $N$  are distribuția uniformă pe mulțimea  $\{0, 1, 2, \dots, M - 1\}$ .

Variabila aleatoare,  $X$ , ce dă numărul de parcurgeri ale buclei este o variabilă aleatoare de tip geometric, cu probabilitatea de succes  $p = P(Y < M) = M/2^k$ .

## 1.4 Distribuția Poisson

O variabilă aleatoare  $X$ , discretă, ce dă numărul de produceri ale unui eveniment rar, într-un interval de timp fixat se numește variabilă aleatoare de distribuție Poisson.

Se presupune că evenimentele ce se pot produce în orice moment al intervalului de timp fixat sunt independente. Dacă evenimentul se produce cu "intensitate" constantă astfel încât, în medie, se produc  $\lambda$  evenimente în intervalul fixat de timp, atunci distribuția de probabilitate a variabilei aleatoare Poisson este:

$$X = \left( \frac{k}{e^{-\lambda} \lambda^k} \right) \quad k = 0, 1, 2, \dots, n, \dots, \quad (8)$$

adică probabilitatea ca evenimentul rar,  $A$ , să se producă de  $k$  ori într-un interval de timp fixat, este:  $P(X = k) = \frac{e^{-\lambda} \lambda^k}{k!}$ . Să verificăm că:

$$\sum_{k=0}^{\infty} P(X = k) = 1. \quad (9)$$

Într-adevăr,  $\sum_{k=0}^{\infty} \frac{e^{-\lambda} \lambda^k}{k!} = e^{-\lambda} \sum_{k=0}^{\infty} \frac{\lambda^k}{k!}$ . Se știe însă că  $e^x = \sum_{k=0}^{\infty} \frac{x^k}{k!}$ . Prin urmare,

$$\sum_{k=0}^{\infty} \frac{\lambda^k}{k!} = e^{\lambda} \quad (10)$$

și deci  $\sum_{k=0}^{\infty} P(X = k) = 1$ .

Notăm cu  $\text{Pois}(\lambda)$  clasa variabilelor aleatoare Poisson de parametru  $\lambda$ .

**Propoziția 1.4.1** *Dacă  $X$  este o variabilă aleatoare de tip Poisson, cu parametrul  $\lambda$ , atunci media sa este  $M(X) = \lambda$  și dispersia  $\sigma^2(X) = \lambda$ .*

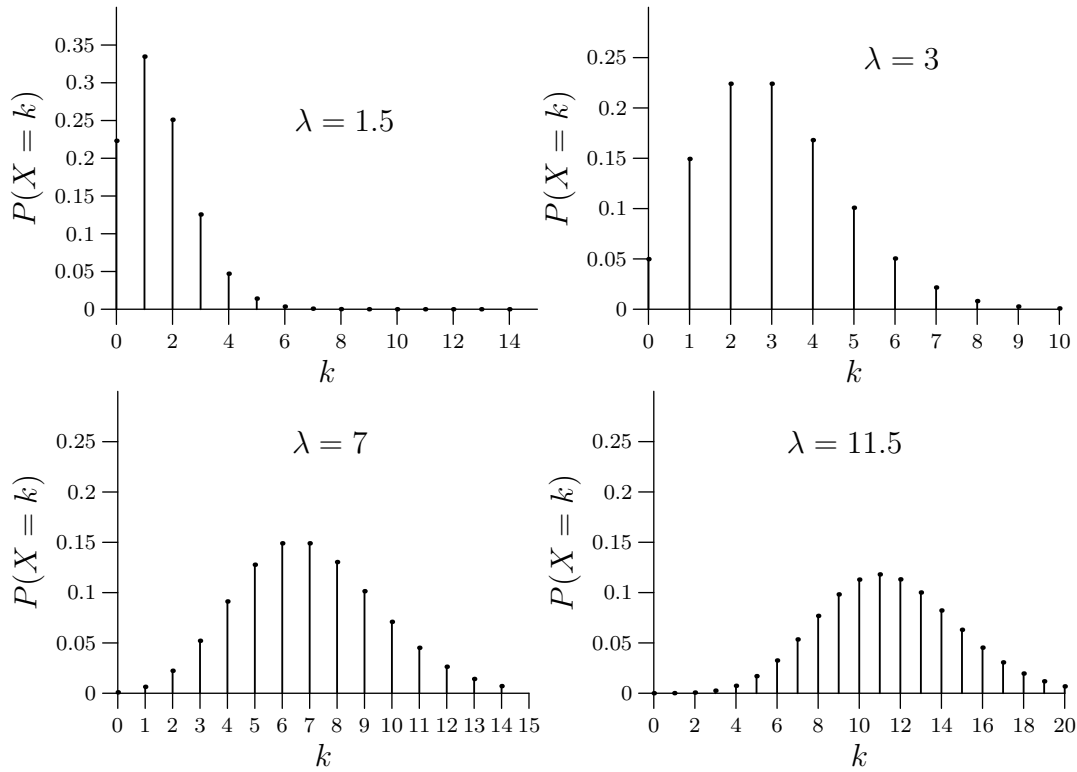
Variabilele aleatoare Poisson sunt singurele cunoscute, pentru care media este egală cu dispersia.

În Fig.2 se poate observa modul în care se schimbă distribuția de probabilitate Poisson pe măsură ce se modifică parametrul  $\lambda$ .

**Exemplul 6.** Un server bază de date primește în medie 30 de cereri de acces pe secundă. Știind că cererile de acces sunt independente, ce distribuție de probabilitate are variabila aleatoare, ce dă numărul acestor cereri/sec? Să se calculeze probabilitatea ca serverul să primească cel puțin o cerere în 6 milisecunde.

Observăm că variabila aleatoare  $X$  care dă numărul cererilor de acces pe secundă se asociază experimentului de observare a evenimentelor rare (deci care nu se produc în mod continuu) de accesare a bazei de date. Prin urmare variabila  $X$  are distribuția Poisson de parametru  $\lambda = 30$  cereri/sec. O milisecundă este  $1/1000$  secunde. Dacă în medie există într-o secundă 30 de cereri, într-o milisecundă numărul mediu de cereri este  $30/1000 = 0.03$ , iar în 6 milisecunde de  $\lambda' = 6(0.03) = 0.18$ . Privind  $X$  ca o variabilă aleatoare ce are distribuția  $\text{Poiss}(0.18)$ , avem de calculat  $P(X \geq 1)$ . Dar cum  $X$  ia valori în  $\mathbb{N}$ , evenimentul  $(X \geq 1) = \mathbb{C}(X < 1) = \mathbb{C}(X = 0)$  și deci  $P(X \geq 1) = 1 - P(X = 0) = 1 - e^{-0.18} \frac{\lambda^0}{0!} = 1 - e^{-0.18}$





**Fig.2:** Vizualizarea distribuției Poisson pentru diferite valori ale parametrului  $\lambda$ .

### Optional

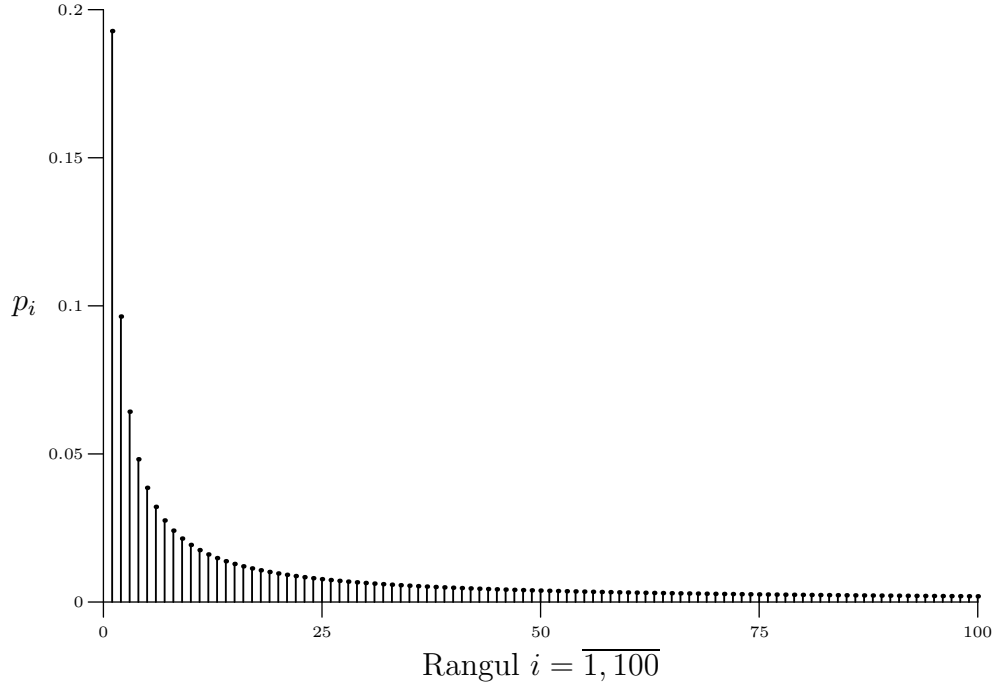
## 1.5 Distribuția Zipf

Distribuția Zipf este o distribuție discretă, ce a fost identificată în numeroase probleme de analiză a cererii de conținut din WWW sau internet. De exemplu, distribuția popularității paginilor WEB, distribuția popularității fișierelor multimedia, descărcate de pe un server, sunt doar câteva exemple.

Pentru a înțelege distribuția Zipf luăm ca exemplu modalitatea de analiză a spamurilor, în scopul proiectării unui filtru spam. Se analizează mai multe emailuri spam și se descompune textul acestora în unități, adică în termeni, cuvinte. Se determină apoi frecvența,  $f_i$ , a termenilor:

$$f_i = \frac{\text{numărul de câte ori apare termenul } t_i}{\text{numărul total de termeni identificați în toate spamurile analizate}}$$

Se ordonează descrescător lista termenilor, în raport cu frecvența, adică termenul  $t_1$  are frecvența maximă și  $t_n$  are frecvența minimă. Indicele  $i$  în această ierarhizare se numește rangul termenului. S-a observat că distribuția de probabilitate a variabilei aleatoare  $X$  ce asociază unui termen, rangul său, are distribuția Zipf, adică probabilitatea ca rangul termenului să fie  $i$  este invers proporțională cu  $i$ :

**Fig.3:** Distribuția Zipf

$$P_X(i) = \frac{c}{i}, \quad i \in \{1, 2, \dots, n\},$$

unde  $c$  o constantă de normalizare, ce se determină din condiția:

$$\sum_{i=1}^n c/i = 1 \quad \Leftrightarrow \quad c = \frac{1}{\sum_{i=1}^n 1/i} = \frac{1}{H_n}$$

Deci variabila aleatoare  $X$  are distribuția de probabilitate:

$$X = \begin{pmatrix} 1 & 2 & \dots & i & \dots & n \\ \frac{1}{H_n} & \frac{1}{2H_n} & \dots & \frac{1}{iH_n} & \dots & \frac{1}{nH_n} \end{pmatrix}$$

Numărul  $H_n = \sum_{i=1}^n \frac{1}{i}$  se numește numărul armonic, de ordin  $n$ .

Astfel modelul probabilist al compoziției spam-urilor analizate este definit de variabila aleatoare:

$$X = \begin{pmatrix} t_1 & t_2 & \dots & t_{n-1} & \dots & t_n \\ \frac{1}{H_n} & \frac{1}{2H_n} & \dots & \frac{1}{iH_n} & \dots & \frac{1}{nH_n} \end{pmatrix}$$

unde  $t_i$ ,  $i = \overline{1, n}$  sunt termenii cei mai frecvenți din spam-uri, ordonați în ordinea descrescătoare a frecvenței de apariție.

În jurul anului 2000 s-au făcut primele investigații statistice relativ la WWW, în scopul de a deduce distribuția de probabilitate a numărului de accesări ale paginilor WEB. Au fost monitorizați userii AOL din punctul de vedere al paginilor pe care le accesează într-o perioadă dată. S-a constatat că există puține site-uri ce sunt vizitate de foarte mulți surferi, în timp ce majoritatea sunt accesate foarte rar. Analizând datele rezultate din monitorizare s-a concluzionat că probabilitatea cererii de accesare a celei de-a  $i$ -a pagini, din ordinea descrescătoare a popularității, este invers proporțională cu  $i$ . Nivelul de popularitate a fost atribuit conform frecvenței de accesare.

Interpretând reprezentarea grafică din Fig. 3 ca indicând probabilitățile de cerere de accesare a primelor 100 de site-uri monitorizate, se observă disproporționalitatea exagerată în atragerea traficului, chiar la paginile 1-5.

Alte observații au condus la concluzia că cererea de accesare/descărcare de conținut de pe internet are distribuția Zipf generalizată, de parametru  $\alpha \in (0, 5, 1]$ . Numărul de accesări/descărcări ale fiecărui obiect depinde de popularitatea acestuia. Presupunem că obiectele sunt indexate descrescător în raport cu popularitatea,  $O_1, O_2, \dots, O_n$ . Probabilitatea ca obiectul  $O_i$  să fie accesat/descărcat este  $P(X = i) = \frac{C}{i^\alpha}$ , adică cererea este caracterizată de o variabilă ce are distribuția Zipf de parametru  $\alpha$ :

$$X = \left( \frac{1}{\frac{C}{1^\alpha}} \quad \frac{2}{\frac{C}{2^\alpha}} \quad \dots \quad \frac{k}{\frac{C}{k^\alpha}} \quad \dots \quad \frac{n}{\frac{C}{n^\alpha}} \right), \quad C = \frac{1}{H_\alpha(n)},$$

unde

$$H_\alpha(n) = \frac{1}{1^\alpha} + \frac{1}{2^\alpha} + \dots + \frac{1}{n^\alpha}$$

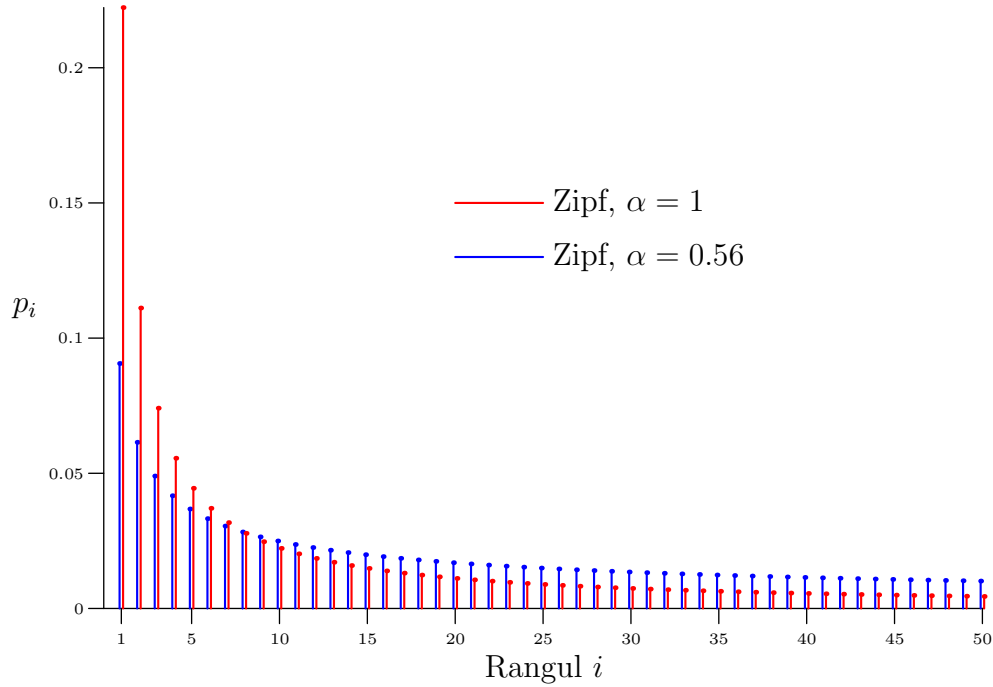
Pe baza acestor observații distribuția Zipf generalizată este exploatată în proiectarea *cache*-ului (alegerea celor mai solicitate pagini WEB, documente, fișiere multimedia, pentru stocare în baza numită *cache*). Și anume, dacă  $N$  este numărul mediu al accesărilor celor  $n$  obiecte de un anumit tip disponibile, într-un interval de timp dat, atunci numărul de accesări ale obiectului  $O_i$  este  $N \cdot P(X = i) = \frac{N}{i^\alpha H_\alpha(n)}$ . Numărul total de accesări a celor mai populare  $k < n$  obiecte este:

$$\sum_{i=1}^k N \cdot P(X = i) = N \sum_{i=1}^k \frac{1}{i^\alpha H_\alpha(n)} = N \frac{H_\alpha(k)}{H_\alpha(n)}$$

Raportul dintre numărul de accesări ale primelor  $k$ , cele mai populare obiecte și numărul de accesări ale tuturor celor  $n$  obiecte este:

$$h = \frac{N \frac{H_\alpha(k)}{H_\alpha(n)}}{N} = \frac{H_\alpha(k)}{H_\alpha(n)}$$

și se numește *hit ratio*. Proiectarea bazei *cache* se face în așa fel încât ea să aibă un raport  $h$  prescris. Din  $h$  și  $\alpha$  cunoscuți se deduce  $H_\alpha(k) = N \cdot H_\alpha(n)$ , și apoi  $k$ .  $k$ , astfel determinat, indică până la ce ordin de popularitate să fie inclus un obiect în *cache*.



**Fig.4:** Distribuție Zipf generalizată comparată cu cea clasică. Sunt trasate câte două segmente corespunzătoare aceleiași valori, unul pentru Zipf clasic (segmentul din dreapta) și unul pentru Zipf generalizat (segmentul din stânga).

Pentru a putea proiecta motoare de căutare, de generație viitoare, este necesar să fie cunoscută în detaliu structura rețelei WWW. WWW este o rețea imensă, în continuă creștere. Aceasta creștere este aleatoare și neregulată: fiecare individ sau organizație putând să-și creeze site-uri WEB, cu un număr nelimitat de documente și linkuri. WWW la un moment,  $t$ , este un graf  $G$  definit de nodurile,  $V = \{1, 2, \dots, n\}$ , reprezentate de paginile WEB și o multime de arce orientate  $E \subset V \times V$ , reprezentate de link-urile dintre pagini. Numărul de pagini (noduri în graf) se numește ordinul grafului, numărul de linkuri din pagina  $i \in V$ , se numește ordinul de ieșire din pagină, iar numărul de linkuri către pagina  $i$ , ordinul de intrare în acea pagină.

Topologia (structura) grafului WWW este dinamică, adică evoluează în timp, prin adăugarea și/sau ștergerea unor pagini și link-uri. Astfel caracteristicile numerice, ca ordinul grafului și gradele de intrare-ieșire din noduri (pagini) se modifică rapid. De aceea s-a ajuns la concluzia că cel mai potrivit model pentru graful WEB este cel de graf aleator.

Pentru a identifica caracteristicile unui astfel de graf, Barabási și Albert au făcut observații asupra unui subgraf din WEB, și anume asupra domeniului `nd.edu`, al Universității Notre Dame, Indiana, ce conținea 325729 documente și 1469680 linkuri. Experimentele au constatat în crawling a subrețelei, la intervale aleatoare de timp. Principala caracteristică observată a fost că distribuția de probabilitate a gradului de intrare,  $G_i$ ,

respectiv iesire,  $G_o$ , într-o/dintr-o pagină, este una de tip Zipf generalizat, cu parametrii  $\alpha = 2.1$ , pentru distribuția gradului de intrare, respectiv  $\alpha = 2.45$ , pentru  $G_o$ .

Pornind de la această observație experimentală s-a născut o nouă direcție de studiu a grafurilor aleatoare, și anume grafuri aleatoare de tip WEB. În Fig. 3 este ilustrată distribuția Zipf clasică, iar în Fig.4 distribuția Zipf generalizată, comparativ cu cea clasică.