

Material suplimentar-curs 8.: Histograma observațiilor asupra unei variabile aleatoare.

Generatori de numere pseudo-aleatoare. Generatorul liniar congruențial

0.1 Histograma observațiilor asupra unei variabile aleatoare

Pentru a evidenția aspectele practice legate de variabilele aleatoare discrete/continue, considerăm graficul unei densități de probabilitate și divizăm domeniul său de definiție prin puncte echidistante x_i , cu pasul h fixat. Probabilitatea ca variabila aleatoare X să ia valori în intervalul $[x_i, x_{i+1})$ este aria trapezului curbiliniu de baze segmentul $[x_i, x_{i+1})$ și arcul de grafic deasupra acestui segment. Aria trapezului o putem aproxima cu aria dreptunghiului de bază $[x_i, x_{i+1})$ și înălțime $f((x_i + x_{i+1})/2)$ (Fig.1, sus), fie cu aria dreptunghiului de bază $[x_i, x_{i+1})$ și înălțime $f(x_i)$ (Fig.1, jos). Deci $P(x_i \leq X < x_{i+1}) = \int_{x_i}^{x_{i+1}} f(x) dx$ = aria trapezului curbiliniu, este aproximativ aria dreptunghiului menționat. Reuniunea tuturor dreptunghiurilor astfel construite se numește histogramă asociată densității de probabilitate. Histograma evidențiază aproximativ distribuția valorilor variabilei aleatoare în intervalele $[x_i, x_{i+1})$.

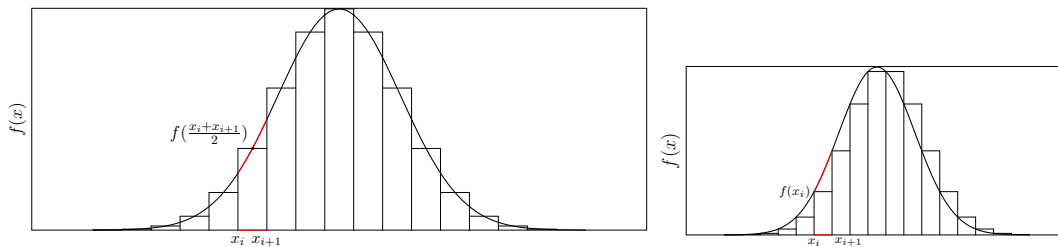


Fig.1: Histograma asociată unei densități de probabilitate

Precizarea distribuției de probabilitate a unei variabile aleatoare continue, se bazează de obicei pe istoricul observațiilor asupra valorilor sale. În experimentele de laborator sau experimentele virtuale, de simulare, se înregistrează valorile observate, măsurate sau generate, ale unei variabile aleatoare, adică o listă de numere reale x_1, x_2, \dots, x_N . Cel mai adesea nu se cunoaște distribuția de probabilitate a variabilei studiate. Informația primară se obține asociind seriei de date înregistrate o histogramă cu m bare, m fixat apriori. Și anume:

- se determină valoarea minimă, $xmin$, și valoarea maximă, $xmax$, a seriei de date.
- Se divide segmentul $[xmin, xmax]$ prin puncte echidistante cu pasul,

$$h = \frac{xmax - xmin}{m}. \text{ Notăm cu } y_j \text{ punctele de diviziune, } y_j = xmin + j * h, j = \overline{0, m}.$$

- Se calculează numărul de valori, n_j , ale seriei de date, care aparțin intervalului

$$I_j = [y_j, y_{j+1}), \quad j = \overline{0, m-2},$$

respectiv intervalului $I_{m-1} = [y_{m-1}, xmax]$

- Deoarece din cele N valori ale seriei de date, n_j cad în intervalul I_j , rezultă că am obținut informația că probabilitatea ca variabila observată să ia valori în intervalul I_j este aproximativ:

$$P(y_j \leq X < y_{j+1}) \approx \frac{n_j}{N}$$

Dar această probabilitate este comparând cu cazul teoretic de mai sus aria dreptunghiului (a barei) ce are baza segmentul $[y_j, y_{j+1})$. • Rezultă astfel că deasupra intervalului I_j desenăm un dreptunghi de arie $A = n_j/N$. Dar aria este baza ori înălțimea dreptunghiului, H_j , și deci din $A = Baza \times H_j$, rezultă că înălțimea dreptunghiului (a barei) este $H_j = A/B$. Cunoscând lungimea bazei ca fiind $h = y_{j+1} - y_j$, rezultă că înălțimea dreptunghiului este $H_j = \frac{n_j}{h N}$, $j = \overline{0, m-1}$.

Obiectul grafic rezultat din desenarea celor m dreptunghiuri (bare) se numește *histograma seriei de date* sau *histograma distribuției de frecvențe*, deoarece ea ilustrează modul în care datele sunt distribuite în intervalele $[y_j, y_{j+1})$.

În Fig.2 sunt ilustrate histogramele a două serii de date.

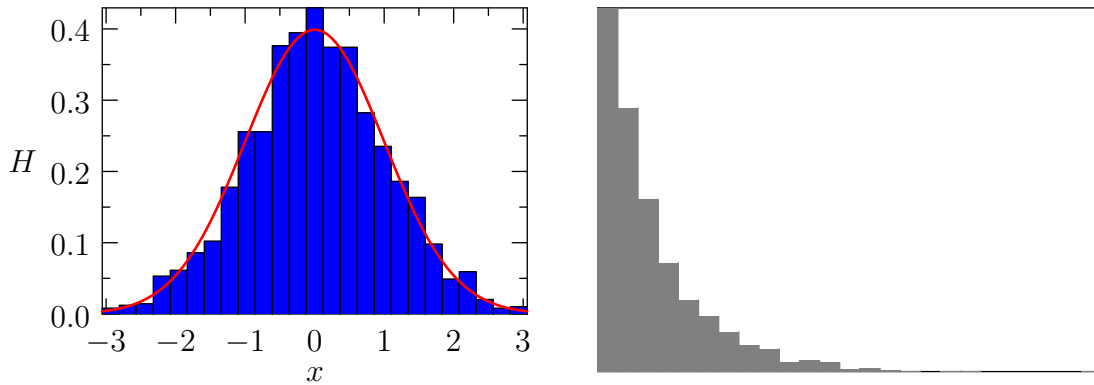


Fig.2: Histograme asociate la două serii de date numerice, rezultate din simularea unor variabile aleatoare ce au densitatea ilustrată în roșu.

0.2 Generatori de numere pseudo-aleatoare. Generatorul liniar congruențial

A simula o variabilă aleatoare, X , revine la a rula un algoritm determinist care produce sau generează un șir de numere, x_1, x_2, \dots, x_N , ce au proprietățile unui șir de valori de observații/măsurători independente asupra variabilei X .

Numerele astfel generate sunt pseudo-aleatoare și se folosesc în probleme de simulare a proceselor din sisteme complexe, cu intrări aleatoare. De asemenea, un generator de numere pseudo-aleatoare este invocat în diverse etape ale algoritmilor probabilisti. Pe de alte parte, numerele pseudo-aleatoare se generează în scopul asigurării securității proceselor într-un sistem de operare sau pentru asigurarea securității transmiterii informației pe un canal de comunicație.

Modalitatea de generare a numerelor pseudo-aleatoare depinde de scopul pentru care acestea sunt folosite. Noi abordăm doar problematica generării numerelor aleatoare pentru prima clasă de probleme enunțată, nu și a celor folosite în criptografie.

Baza oricărei simulări a unui fenomen sau proces aleator o constituie numerele uniform distribuite pe intervalul $[0, 1)$.

Aleatorul este în general greu de definit, dar cel mai adesea un fenomen este considerat aleator dacă este imprevizibil și nereproductibil. Cum un algoritm determinist nu poate genera șiruri de numere cu aceste calități, numerele generate se numesc pseudo-aleatoare, nu aleatoare. Algoritmii folosiți în simularea variabilelor aleatoare se numesc *generatori de numere pseudo-aleatoare*.

Metoda de generare de numere aleatoare uniform distribuite, cel mai frecvent folosită până de curând, este metoda liniar congruențială. Generatorul liniar congruențial a fost introdus de Lehmer, fost profesor la Universitatea Berkeley, care a fost și unul dintre fondatorii teoriei computaționale a numerelor.

O astfel de metodă generează numere în inelul, \mathbb{Z}_p , al claselor de resturi modulo p , unde p este un număr natural fixat, numit *modul*. Dacă $n \in \mathbb{Z}$, $n \pmod{p}$ este restul împărțirii lui n la p . Teorema împărțirii cu rest din aritmetică, asigură că pentru $n \in \mathbb{Z}$, $p \in \mathbb{N} \setminus \{0\}$ există $q \in \mathbb{Z}$ și $r \in \mathbb{N}$ astfel încât $n = qp + r$, $r \in \{0, 1, 2, \dots, p-1\}$. Prin urmare inelul \mathbb{Z}_p conține clasele de resturi $\{0, 1, 2, \dots, p-1\}$.

Generatorul liniar congruențial produce un șir de numere x_1, x_2, \dots, x_N din \mathbb{Z}_p printr-o formulă recursivă:

$$x_n = ax_{n-1} + b \pmod{p} \quad (1)$$

unde parametrii a, b sunt fixați în \mathbb{Z}_p , iar x_0 se numește valoare inițială sau *seed*. Cu alte cuvinte x_n este restul împărțirii numărului $ax_{n-1} + b$ la modulul p .

Cum elementele șirului (x_n) aparțin mulțimii $\{0, 1, 2, \dots, p-1\}$, șirul asociat (u_n) , cu $u_n = \frac{x_n}{p}$, este un șir de numere din intervalul $[0, 1)$.

Are șirul (u_n) , $n = \overline{1, N}$, atributele unui șir de valori de observație asupra unui și de variabile i.i.d (U_n) , $U_n \sim \text{Unif}[0, 1)$, $n = \overline{1, n}$? Pentru a răspunde acestei întrebări, caracterizăm șirul $(x_n) \subset \mathbb{Z}_p$.

• Șirul (x_n) definit în mod recursiv pornind de la valoarea inițială x_0 este reproductibil (repetabil) și reproductibilitatea nu este un atribut al aleatorului. Opțiunea pentru astfel de generatori se explică prin faptul că permit verificarea și *debugging*-ul codului implicat în simulare, folosind de mai multe ori același șir de numere produse de un generator.

• Șirul (x_n) , fiind un șir de elemente dintr-o mulțime finită, este periodic, adică există un T astfel încât $x_{k+T} = x_k$, $\forall k \in \mathbb{N}$.

De exemplu, generatorul de modul $p = 16$, $a = 3$, $b = 4$ și $x_0 = 0$ conduce la șirul periodic $0, 4, 0, 4, \dots$. Luând $x_0 = 1$ obținem șirul:

$$7, 9, 15, 1, 7, 9, 15, 1, \dots$$

ce are perioada 4 (secvența 7, 9, 15, 1 se repetă). Luând $p = 16$, $a = 5$, $b = 3$, valoarea inițială nu are importanță, deoarece generatorul produce șirul periodic cu secvența repetitivă:

$$0, 3, 2, 13, 4, 7, 6, 1, 8, 11, 10, 5, 12, 15, 14, 9$$

Evident că un șir periodic nu poate fi considerat aleator. Dacă însă perioada este suficient de mare în raport cu numărul de termeni ce se folosesc în simulare, atunci el este acceptabil dacă șirul (u_n) , cu $u_n = x_n/p$ trece anumite teste.

Sunt considerați generatori buni, cei pentru care lungimea perioadei este aceeași pentru orice valoare inițială, și mai mult perioada este maximum posibil.

Studii experimentale îndelungate, pe diferite tipuri de calculatoare, recomandă următorii generatori liniar congruențiali ai căror parametri (p, a, b, x_0) sunt, respectiv:

$$\begin{aligned} &(2147483647, 16807, 0, 1), \\ &(2147483647, 950706376, 0, 1), \\ &(2147483647, 630360016, 0, 1), \\ &(2147483648, 452807053, 0, 1), \\ &(2147483647, 1078318381, 0, 1), \end{aligned} \tag{2}$$

$p = 2147483647 = 2^{31} - 1$, este număr prim de tip Mersenne, adică un număr prim de forma $2^q - 1$. Astfel generatorii asociați au perioada maximă $p - 1 = 2147483646$. Prin urmare se poate genera un șir de peste 2 miliarde de numere distincte x din \mathbb{Z}_p .

După ce perioada maximă a fost asigurată, generatorul este supus unor teste de uniformitate.

ISO C pune la dispoziție funcția din `stdlib.h`:

```
int rand(void);
```

ce implementează generatorul liniar congruențial de modul $p = 2^{31}$. Funcția returnează numere de tip `int` (pe 32 biți) din mulțimea $\{0, 1, 2, \dots, \text{RAND_MAX}\}$, unde constanta `RAND_MAX` este $2^{31} - 1$.

Setarea valorii inițiale se realizează cu funcția:

```
void srand(unsigned seed);
```

Astfel codul:

```

#include<stdlib.h>
#define N 1000

int main()
{
    int i, r;
    double u[N];
    srand(231);

    for(i=0;i<N;i++)
    {
        r=rand();
        u[i]=(double)r/RAND_MAX;
    }
    return 0;
}

```

generează un șir de 1000 de numere pseudo-aleatoare "uniform distribuite" pe $[0, 1)$.

De ce s-a ales pentru `rand()` un modul, putere a lui 2, $p = 2^m$? Motivul este că restul împărțirii unui număr întreg pozitiv, x , la 2^m este numărul binar constând din ultimii m biți ai lui x , care se poate calcula foarte rapid.

Într-adevăr, dacă x se exprimă în binar prin: $x = (b_{31} \dots b_m b_{m-1} \dots b_1 b_0)_2$, atunci

$$\begin{aligned}
 x &= b_{31}2^{31} + \dots + b_m 2^m + b_{m-1} 2^{m-1} + \dots + b_1 2 + b_0 2^0 \\
 &= 2^m (b_{31} 2^{31-m} + \dots + b_m 2^0) + b_{m-1} 2^{m-1} + \dots + b_0 2^0
 \end{aligned}$$

Prin urmare restul împărțirii lui x la 2^m este în binar numărul $b_{m-1} 2^{m-1} + \dots + b_0 2^0 = (0 \dots 0 b_{m-1} \dots b_1 b_0)_2$.

Pentru calculul rapid a numărului $(0 \dots 0 b_{m-1} \dots b_1 b_0)_2$ se ține seama că

- $p = 2^m = (00 \dots 0 \underbrace{1}_{c_m} 0 \dots 0)_1$

- $p - 1 = (00 \dots 00 \underbrace{1 \dots 1}_{c_{m-1} \dots c_1 c_0});$

- Practic se definește `masca=p-1` efectuând operații pe biți, astfel:

`masca=(1<<m)-1;`

- Atunci restul împărțirii lui `int x` la $p = 2^m$ este `r=masca & x;`

Testul de k -uniformitate: Elementele unui șir de numere pseudo-aleatoare (u_n) , $n = \overline{1, k}$, $u_n \in [0, 1)$, trebuie să aibă atributele unor valori de observație asupra unui șir de variabile aleatoare U_1, U_2, \dots, U_k , independente și uniform distribuite pe intervalul $[0, 1)$.

Pentru a prezenta unul din cele mai folosite teste de uniformitate și independență, reamintim că dacă U este o variabilă aleatoare uniform distribuită pe intervalul $[0, 1)$, atunci probabilitatea ca U să ia valori într-un subinterval $[a, b] \subset [0, 1)$ este egală cu lungimea subintervalului: $P(U \in [a, b]) = b - a$, iar dacă U_1, \dots, U_k , sunt variabile

aleatoare independente și identic distribuite după legea uniformă pe $[0, 1)$, atunci probabilitatea ca vectorul aleator (U_1, U_2, \dots, U_k) să ia valori în k -paralelipipedul $[a_1, b_1] \times [a_2, b_2] \times \dots \times [a_k, b_k] \subset [0, 1)^n$ este egală cu volumul paralelipiedului:

$$\begin{aligned} P(a_1 \leq U_1 \leq b_1, \dots, a_k \leq U_k \leq b_k) &= \\ P(a_1 \leq U_1 \leq b_1) \cdots P(a_k \leq U_k \leq b_k) &= \\ (b_1 - a_1) \cdots (b_k - a_k), \quad \forall a_i, b_i, 0 \leq a_i < b_i < 1, i = \overline{1, k} \end{aligned} \quad (3)$$

Având un șir de numere (u_n) din intervalul $[0, 1)$, se pune problema să decidem dacă acest șir este "aproximativ" uniform distribuit. În acest scop se consideră vectorii constituiți din k -numere consecutive ale șirului (u_n) :

$$(u_0, u_1, \dots, u_{k-1}), (u_1, u_2, \dots, u_k), (u_2, u_3, \dots, u_{k+1}), \dots, \quad (4)$$

Pentru orice k -paralelipiped $D = [a_1, b_1] \times \dots \times [a_k, b_k]$, $0 \leq a_i < b_i < 1$, $i = \overline{1, k}$, se definește funcția caracteristică:

$$\mathbf{1}_D(u) = \begin{cases} 1 & \text{dacă } y \in D \\ 0 & \text{dacă } y \notin D \end{cases} \quad (5)$$

Suma $\sum_{i=0}^{n-1} \mathbf{1}_D(u_i, u_{i+1}, \dots, u_{i+k-1})$ dă numărul de k -vectori, dintre cei n constituiți ca în (4), ce aparțin paralelipipedului D .

Definiția 0.2.1 Șirul $(u_n) \subset [0, 1)$ se numește șir k -uniform, $k \geq 2$, dacă pentru orice k -paralelipiped $D = [a_1, b_1] \times \dots \times [a_k, b_k]$, $0 \leq a_i < b_i < 1$, $i = \overline{1, k}$,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=0}^{n-1} \mathbf{1}_D(u_i, u_{i+1}, \dots, u_{i+k-1}) = (b_1 - a_1) \cdots (b_k - a_k) = \text{vol}(D). \quad (6)$$

Practic, definiția spune că dacă estimăm probabilitatea evenimentului ca k -vectorii constituiți din șirul generat să cadă într-un k -paralelipiped, ca limita frecvențelor experimentale de producere a evenimentului, atunci șirul este k -uniform dacă din k -vectorii constituiți, proporția celor ce cad într-un k -paralelipiped este aproximativ egală cu volumul paralelipipedului.

Testele de k -uniformitate asupra generatorilor liniar congruențiali au adus numeroase surprize. Să analizăm, de exemplu, generatorul **randu** care a fost inclus în biblioteca științifică a mainframe-urilor IBM (IBM 360/370) un număr mare de ani și folosit pentru simulări în numeroase proiecte științifice de anvergură. Parametrii generatorului **randu** sunt:

$$p = 2^{31}, a = 65539 = 2^{16} + 3, b = 0$$

După mulți ani de folosire, Marsaglia, profesor la Universitatea din Florida, a observat următoarea deficiență a generatorului **randu**:

$$\begin{aligned} x_{n+2} &= (2^{16} + 3)x_{n+1} = (2^{16} + 3)^2 x_n \\ &= (2^{32} + 6 \cdot 2^{16} + 9)x_n = \underbrace{2 \cdot 2^{31}}_{0 \bmod 2^{31}} x_n + (6 \cdot 2^{16} + 18 - 9)x_n \\ &= [6 \cdot (2^{16} + 3) - 9]x_n = 6(2^{16} + 3)x_n - 9x_n = 6x_{n+1} - 9x_n, \end{aligned}$$

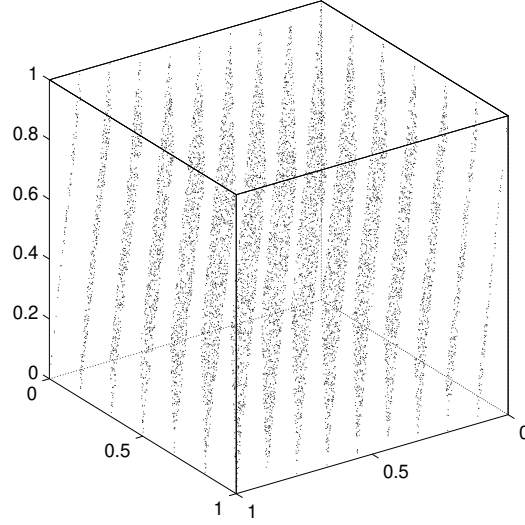


Fig.3: Ilustrarea dispunerii pe plane paralele a tripletelor de numere (u_n, u_{n+1}, u_{n+2}) asociate unui șir (u_n) , generat de **randu**.

adică $9x_n - 6x_{n+1} + x_{n+2} = 0 \pmod{2^{31}}$. Relația $9x_n - 6x_{n+1} + x_{n+2} = 0 \pmod{2^{31}}$ este echivalentă cu $9x_n - 6x_{n+1} + x_{n+2} = k \cdot 2^{31}$, $k \in \mathbb{Z}$. Prin împărțire la 2^{31} rezultă că pentru orice n , tripletele (u_n, u_{n+1}, u_{n+2}) aparțin unor plane de ecuație $9x - 6y + z = k$:

$$9u_n - 6u_{n+1} + u_{n+2} = k, \quad u_n = x_n/2^{31}.$$

Dintre toate planele paralele de ecuație $9x - 6y + z = k$, $k \in \mathbb{Z}$, intersectează cubul unitate doar cele corespunzătoare lui $k \in \{-5, -4, \dots, 9\}$. Prin urmare tripletele de numere generate, (u_n, u_{n+1}, u_{n+2}) , sunt dispuse în cubul $[0, 1]^3$ pe 15 plane paralele (Fig.3).

Această particularitate ilustrează că punctele (u_n, u_{n+1}, u_{n+2}) nu sunt uniform dispersate în cub, și deci șirul (u_n) nu are atributele unui șir uniform distribuit.

După depistarea acestei deficiențe a generatorului **randu** s-a demonstrat că orice generator liniar congruențial are acest defect de regularitate în anumite dimensiuni k , adică k -punctele:

$$(u_0, u_1, \dots, u_{k-1}), (u_1, u_2, \dots, u_k), \dots$$

constituite din elemente ale șirului (u_n) sunt dispuse pe un număr redus de hiperplane din hipercubul $[0, 1]^k$, în loc să fie dispersate în hipercub.

În Fig.4 sunt reprezentați vectorii (u_k, u_{k+1}) , $k = \overline{1, 3000}$, constituiți din perechi de numere aleatoare consecutive, produse de generatorul liniar congruențial de parametri $a = 65$, $b = 1$, $p = 2048$ (stânga), respectiv $a = 3$, $b = 0$, $p = 2048$ (dreapta).

Toate limbajele comune de programare **C**, **C++**, **Java**, conțin funcții ce implementează un generator liniar congruențial. Știind că au acest defect, funcțiile respective nu sunt indicate în probleme serioase de simulare.

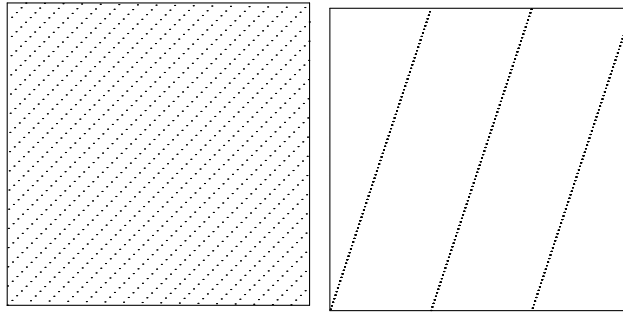


Fig.4: Structura regulată a numerelor aleatoare produse de generatori liniar congruențiali.

Proprietățile de regularitate ale șirului generat prin metoda liniar congruențială au condus la dezvoltarea unor metode ne-congruențiale de generare de numere pseudo-aleatoare uniform distribuite pe $[0, 1)$.

Cel mai performant generator existent la ora actuală este generatorul numit, Mersenne-Twister, dezvoltat de M. Matsumoto și T. Nishimura

Perioada șirului generat de Mersenne-Twister este $2^{19937} - 1$. Șirurile generate au trecut teste de k -uniformitate pentru orice $k \leq 623$.

Acest generator este implementat în Python, MATLAB, PHP, Ruby, etc.