

Laborator 4

Analiza și prelucrarea datelor cu erori aleatorii sau cu evoluție aleatorie

1. Noțiuni teoretice

Erorile aleatorii (întâmplătoare/accidentale) se caracterizează prin faptul că, la repetarea măsurărilor, în condiții identice, asupra aceluiași măsurand, apar diferențe atât ca sens cât și ca valoare, nefiind posibil să se stabilească reguli de predeterminare pentru ele; caracterizarea lor poate fi făcută numai în sens probabilistic.

Pentru caracterizarea erorilor aleatorii, conținute în datele experimentale, se **aplica metoda selecției (esantionului)**, în sensul că se consideră un sir de măsurări individuale V_1, V_2, \dots, V_n , în care n se numește **volumul selecției (o selecție de volum n)**.

Se parcurg următoarele etape:

A) Deoarece, în forma obținută prin experiment, datele constituie o mulțime dezordonată de valori, **se ordonează datele V_i în sens crescător, formându-se seria variațională**; rezultă că V_1 apare de n_1 ori, V_2 de n_2 ori, ..., V_k de n_k ori, cu:

$$V_1 < V_2 < \dots < V_k \quad \text{și} \quad \sum_{i=1}^k n_i = n$$

B) Se reprezintă rezultatele măsurărilor individuale printr-o **histogramă** sau printr-un **poligon de frecvențe**.

Histograma este o reprezentare în plan, pe ordonată considerându-se *frecvența absolută de apariție* n_i , sau *frecvența relativă de apariție* $f_i = n_i/n$, iar pe abscisă se împarte domeniul valorilor $[V_{\min}, V_{\max}]$ în intervale elementare de lungime Δ , denumite *intervale de grupare sau de clasă*, Δ fiind dat de **formula lui Sturges**:

$$\Delta = \frac{V_{\max} - V_{\min}}{1 + 3,22 \cdot \lg n}$$

Ca punct de plecare în construcția histogramei se consideră - în general - valoarea V_{\min} , astfel ca datele se grupează în intervalele $[V_{\min} + i\Delta; V_{\min} + (i+1)\Delta)$ deschise la dreapta.

C) *Repartitia empirică este prelucrată statistic* în vederea obținerii unor **valori tipice de selecție** ca: **media m_v** , modul (dominantă) Mo , mediana Me , dispersia de selecție μ , dispersia reală σ , **estimația $\hat{\sigma}$** , eroarea medie absolută θ etc, care reprezintă indicatori sintetici esențiali pentru evaluarea erorilor aleatorii.

Valoarea medie m_v este foarte importantă întrucât ea se consideră, în mod obișnuit, ca mărime de referință (valoare convențională). Se mai numește și valoarea cea mai probabilă.

$$m_v = \frac{1}{n} \sum_{i=1}^n V_i$$

Estimația erorii standard (deviația standard) reprezintă valoarea medie a patratelor erorilor aparente:

$$\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (V_i - m_v)^2$$

Dispersia de selecție μ^2 reprezintă media aritmetică a patratelor erorilor aparente:

$$\mu^2 = \frac{1}{n} \sum_{i=1}^n \Delta V_i^2 = \frac{1}{n} \sum_{i=1}^n (V_i - m_v)^2$$

Eroarea medie absoluta de selectie θ reprezinta media aritmetica a erorilor reale în

modul, dar în practică se folosește $\hat{\theta}$: adica:
$$\hat{\theta} = \frac{1}{n} \sum_{i=1}^n |V_i - m_v|$$

Mai multe valori tipice de selecție sunt definite în cursul nr. 3.

2. Realizarea unei histograme în MATLAB folosind generatorul de numere aleatoare

În cazul în care nu aveți la dispoziție date experimentale pentru realizarea unei histograme în MATLAB, puteți folosi funcția **random** pentru a genera numere aleatoare cu o distribuție de probabilitate dată. Acest mod de lucru este util în special pentru a realiza simulări. În acest caz vom simula un set de date pentru a testa realizarea histogramei.

Sintaxa pentru apelarea **funcției random** este:

- random(NUM, A) returnează un vector de numere aleatoare, alese din distribuția de probabilitate cu un parametru, specificat de NUM, cu valoarea parametrului A.
- random(NUM, A, B) sau random(NUM, A, B, C) returnează un vector de numere aleatoare, alese dintr-o distribuție de probabilitate cu valorile parametrilor A, B (și C).

Există multe tipuri de distribuții care pot fi alese la utilizarea acestei funcții. Introduceți **help random** în linia de comandă pentru a lansa o listă completă cu numele și parametrii de intrare.

| Distribuție | Parametrul de intrare A | Parametrul de intrare B |
|-------------------------|-------------------------|--|
| 'bino' sau 'Binomial' | n: numărul de încercări | p: probabilitatea de succes pentru fiecare încercare |
| 'exp' sau 'Exponential' | μ : media | - |
| 'norm' sau 'Normal' | μ : media | σ : deviația standard |

Odată creat setul de date cu ajutorul generatorului de numere aleatoare, puteți reprezenta grafic datele, folosind funcția histograma.

Exemplul 1

Realizați o rutină MATLAB pentru a genera o selecție aleatoare de 1000 de puncte de date dintr-o distribuție Gaussiană cu $m_v = 1$ și $\hat{\sigma} = 0.5$.

```
media=1;
dispersia=0.5;
N=1000;
date=random('Norm',media,dispersia,1,N);
figure(1);
hist(date)
title('Histograma cu media=1, dispersia=0.5');
```

Completați rutina Matlab cu instrucțiunile aferente pentru a obține alte 2 grafice în care se modifică, pe rând, valorile pentru medie și apoi pentru estimatie. Care sunt efectele modificării m_v ? Dar la mărirea valorii $\hat{\sigma}$?

3. Determinarea histogramei pentru un set de date

Exemplul 2

Se consideră un set de date care constituie notele obținute de studenți la un examen. Aceste date se găsesc în fisierul "Student/Note Stud.xls".

Se dorește găsirea valorilor tipice de selecție media m_v și estimația $\hat{\sigma}$. În vederea determinării acestor valori, datele din fișier trebuie importate ca și vector.

Folosiți Help-ul Matlab-ului pentru a înțelege sintaxa instrucțiunilor "xlsread", "mean", "std" pentru **determinarea** m_v , $\hat{\sigma}$. Realizați **histograma** aferentă acestui set de date, utilizând instrucțiunea "hist". Aceasta trebuie să rezulte ca în **figura 1**.

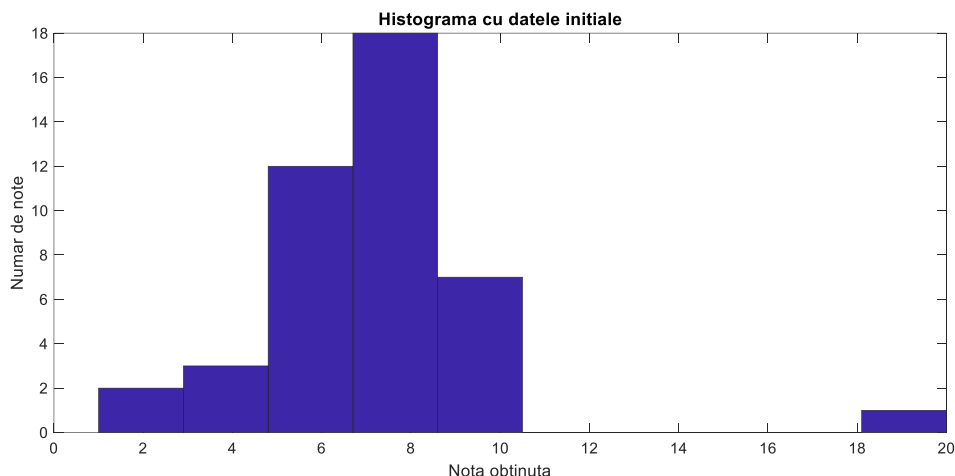


Fig.1. Histograma Exemplul 2

4. Exerciții propuse

4.1. La măsurarea dozelor anuale de radiații din incinta unei facilități de cercetare în domeniul nuclear sunt utilizate 10 aparate de măsură (dozimetre) situate în locații diferite. Doza maximă acceptată de radioactivitate pentru o persoană obișnuită, aprobată de Comisia Națională pentru Controlul Activităților Nucleare și de Organizația Mondială a Sănătății este de 1mSv (miliSievert). Datele înregistrate de cele 10 dozimetre se găsesc în tabelul următor:

| Dozimetru | Nivel de radioactivitate [mSv] |
|-----------|--------------------------------|
| 1 | $0.121 \cdot 10^{-3}$ |
| 2 | $0.217 \cdot 10^{-3}$ |
| 3 | $0.325 \cdot 10^{-3}$ |
| 4 | $0.517 \cdot 10^{-3}$ |
| 5 | $0.476 \cdot 10^{-3}$ |
| 6 | $0.192 \cdot 10^{-3}$ |
| 7 | $0.088 \cdot 10^{-3}$ |
| 8 | $0.235 \cdot 10^{-3}$ |
| 9 | $0.217 \cdot 10^{-3}$ |
| 10 | $0.112 \cdot 10^{-3}$ |

Să se determine următoarele valori tipice de selecție pentru măsurătorile obținute de cele 10 dozimetre:

- media aritmetica $mv1$;
- mediana de selecție me ;
- dispersia de selecție $\mu^2 niu$;
- estimația $\hat{\sigma}$ $sig1$
- eroarea medie absoluta de selecție $teta$.

4.2. Numărul de mesaje trimise pe oră printr-o rețea de calculatoare este exprimat prin frecvența de apariție f_i și are următoarea distribuție:

| | | | | | | |
|----------------------------|------|------|-----|-----|-----|------|
| $x=\text{numar de mesaje}$ | 10 | 11 | 12 | 13 | 14 | 15 |
| $f_i(x)=n_i/n$ | 0.08 | 0.15 | 0.3 | 0.2 | 0.2 | 0.07 |

Calculați media aritmetică m_{v2} și estimația standard $sig2$ pentru numărul de mesaje trimise orar.

4.3. Se consideră setul de date "Date_temp.xls". Să se determine:

- media m_v m_{v3} și estimația $\hat{\sigma}$ $sig3$
- histograma