

Curs 10

Elemente de statistică

10.1 Populații statistice

Statistica se poate defini ca fiind știința care se ocupă cu colectarea, prezentarea, clasificarea, analiza și interpretarea cantitativă a datelor și cu aplicarea teoriei probabilităților în analiza și estimarea caracteristicilor (parametrilor) populației.

Definiția 10.1.1 Populația statistică este o mulțime de elemente supusă cercetării statistice. Un element al acestei populații se numește unitate statistică.

Definiția 10.1.2 O proprietate comună tuturor unităților statistice se numește caracteristică statistică a populației.

Caracteristicile pot fi:

- **cantitative**, dacă pot fi cuantificate printr-un număr,
- **calitative**, dacă nu pot fi cuantificate printr-un număr, ci prin aprecieri de tipul bun, foarte bun, mult, puțin.

Definiția 10.1.3 Valoarea numerică a unei caracteristici cantitative, care se modifică de la o unitate statistică la alta, se numește variabilă statistică.

Definiția 10.1.4 O submulțime a unei populații statistice se numește selecție sau eșantion.

De exemplu, dacă se solicită opinia unei populații într-o problemă oarecare este dificil să fie consultată întreaga populație și se recurge la extragerea unui eșantion.

Un eșantion ale cărui unități au fost alese la întâmplare se numește *eșantion reprezentativ* (*aleator*). Extragerea unui eșantion reprezentativ se realizează astfel încât elementele să aibă șanse egale de a fi extrase. În multe situații se folosesc numere aleatoare. De exemplu dacă avem o listă a populației inițiale x_1, x_2, \dots, x_N se aleg numere aleatoare k_1, k_2, \dots, k_n între 1 și N și un eșantion aleator este $x_{k_1}, x_{k_2}, \dots, x_{k_n}$.

Variabilele statistice pot fi *discrete* (ex: numărul de becuri care se ard după 1000 de ore de întrebuințare) sau *continue* (ex: timpii de defectare a unui număr fixat de n becuri). Aceste rezultate experimentale se prezintă sub forma diverselor scheme grafice care indică o anumită comportare a variabilei statistice. În cazul variabilelor discrete rezultatele sunt prezentate sub forma unui tablou de forma următoare:

X	x_1	x_2	x_3	x_4	\dots	x_k
n	n_1	n_2	n_3	n_4	\dots	n_k

(10.1)

cu $n = \sum_{j=1}^k n_j$. Numărul n se numește *volumul selecției*. Presupunem că în tabel am aranjat astfel încât $x_1 < x_2 < x_3 < \dots < x_k$.

Frecvența absolută a unei valori x_j , $j = \overline{1, k}$ a unei caracteristici este numărul n_j și reprezintă numărul de unități statistice ale populației statistice care corespund valorii x_j a variabilei statistice X .

Frecvența relativă a unei valori x_j , $j = \overline{1, k}$ a unei caracteristici este raportul dintre frecvența absolută și volumul selecției, adică $f_j = \frac{n_j}{n}$, $j = \overline{1, k}$. Suma frecvențelor relative este egală cu 1.

Frecvența absolută cumulată crescător este suma frecvențelor absolute valorilor mai mici sau egale cu valoarea x_j : $\sum_{i=1}^j n_i$, $j \leq k$.

Frecvența relativă cumulată crescător a unei valori x_j , $j = \overline{1, k}$ a unei caracteristici este suma frecvențelor relative corespunzătoare valorilor mai mici sau egale cu x_j : $\frac{1}{n} \sum_{i=1}^j n_i$, $j \leq k$.

Definiția 10.1.5 Se numește funcție de repartiție empirică asociată unei v. a. X și unei selecții (x_1, x_2, \dots, x_n) *Funcția*

$$F_n^* : \mathbb{R} \rightarrow \mathbb{R}, F_n^*(x) = \frac{\text{număr de valori } x_j < x}{n} = \frac{k_x}{n}.$$

Teorema de mai jos pune în evidență că funcțiile de repartiție empirice aproximează oricât de bine funcția reală de repartiție.

Teorema 10.1.6 Fie o populație statistică și X variabila aleatoare atașată ei cu funcția de repartiție F . Pentru o selecție de volum n : (x_1, x_2, \dots, x_n) construim ca mai sus funcția de repartiție empirică F_n^* . Atunci

$$\forall \varepsilon > 0 : P \{ |F(x) - F_n^*| \geq \varepsilon \} \xrightarrow{n \rightarrow \infty} 0, \quad \varepsilon \text{ fixat.}$$

Demonstrație. Să notăm cu $p = P(X < x) = F(x)$, și cu $F_n^*(x) = \frac{k_x}{n}$ funcția de repartiție empirică. Definim v. a. Y_1, Y_2, \dots, Y_n , astfel încât Y_j are valoarea 1 dacă $x_j < x$ și 0 în caz contrar. Variabilele Y_1, Y_2, \dots, Y_n sunt independente (selecțiile fiind independente) și au distribuția

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

Avem $M[Y_j] = p$, $D[Y_j] = p(1-p)$. Evident v. a. $Z_n = \frac{1}{n} \sum_{i=1}^n Y_i$ are $M[Z_n] = p$ și dispersia

$$D[Z_n] = \frac{p(1-p)}{n}.$$

Aplicăm inegalitatea lui Cebășev lui Z_n și găsim

$$P \{ |F(x) - F_n^*| \geq \varepsilon \} = P \{ |Z_n - p| \geq \varepsilon \} \leq \frac{D[Z_n]}{\varepsilon^2} = \frac{p(1-p)}{n} \xrightarrow{n \rightarrow \infty} 0. \quad \square$$

10.2 Reprezentarea grafică a datelor statistice

O reprezentare a datelor statistice este sub forma histogramei. *Histograma* se construiește sub forma unor dreptunghiuri lipite, cu baza pe Ox , mărimea lor fiind egală la bază cu

mărimea intervalului de variație respectiv. Înălțimea dreptunghiului va fi dată de frecvența corespunzătoare fiecărui interval de variație. Histograma arată forma de repartiție, densitatea de repartiție a frecvențelor și gradul de asimetrie.

Exemplul 10.2.1 Punctele obținute de studenți care au promovat examenul de matematică și care cuantifică cunoștințele lor sunt:

{64, 62, 76, 82, 66, 76, 72, 71, 74, 72, 71, 73, 70, 75, 77, 84, 92, 86, 62, 58, 78, 80, 79, 84, 83, 82, 66, 68, 68, 82, 84, 78, 76, 69, 77, 58, 62, 82, 85, 58, 78, 84, 94, 88, 77, 78, 88, 91, 70, 71, 78, 58, 65, 53, 60, 49, 68, 74, 71, 66, 68, 71, 73, 70, 85, 78, 65, 54, 51, 78, 89, 66, 68, 95, 94, 99, 81, 81, 92, 88, 99, 81, 81}

O reprezentare a acestor date este:

X	49	51	53	54	58	60	62	64	65	66	68	69
n_i	1	1	1	1	4	1	3	1	2	4	5	1
X	70	71	72	73	74	75	76	77	78	79	80	81
n_i	3	5	2	2	2	1	3	3	7	1	1	4
X	82	83	84	85	86	88	89	91	92	94	95	99
n_i	4	1	4	2	1	3	1	1	2	2	1	2

În cazul variabilelor continue, datele se grupează într-un tablou de forma:

X	$[t_1, t_2)$	$[t_2, t_3)$	\dots	$[t_m, t_{m+1})$
n	n_1	n_2	\dots	n_m

În acest caz X este variabila statistică, iar n_j este numărul valorilor variabilei statistice aflate în intervalul $[t_j, t_{j+1})$ și se numește *frecvența absolută*. La fel $n = \sum_{j=1}^m n_j$ este *volumul selecției*.

Frecvența relativă corespunzătoare intervalului $[t_j, t_{j+1})$ este $\frac{n_j}{n}$ și reprezintă numărul valorilor variabilei statistice aflate în intervalul $[t_j, t_{j+1})$ împărțit la volumul selecției.

Analog se introduc și celelalte noțiuni de frecvență absolută cumulată, frecvență relativă cumulată etc.

10.2.1 Gruparea pe clase

În urma oricărei selecții de volum n dintr-o populație de numere se obține un șir finit de n numere numit *serie statistică* (de volum n). Cum construim o densitate de probabilitate empirică? Pentru a răspunde la această întrebare grupăm termenii unei serii statistice în intervale disjuncte: I_1, I_2, \dots, I_r , după criterii mai mult sau mai puțin subiective.

Dacă avem o selecție x_1, x_2, \dots, x_n de volum n obținute în urma măsurării unor mărimi fizice, tehnice, economice de care depinde evoluția unor procese sau fenomene, notăm cu $m = \min_i x_i$ și $M = \max_i x_i$. Organizăm datele selecției în grupe sau subintervale obținute divizând intervalul $[m, M]$ într-un număr r de subintervale congruente (de aceeași lungime). Pentru determinarea numărului de subintervale de grupare poate fi folosită *formula lui H. A. Sturges*:

$$r = [1 + 3.322 \lg n]. \quad (10.2)$$

Definiția 10.2.2 Se numește pasul de histogramă numărul $h = \frac{M - m}{r}$.

Intervalele de grupare sunt astfel

$$I_1 = [m, m + h), \quad I_2 = [m + h, m + 2h), \dots, \quad I_r = [m + (r - 1)h, M].$$

Formula lui Sturges (10.2) inițial a fost scrisă sub forma echivalentă $r = 1 + \log_2 n$. Într-adevăr, avem $\log_2 n = \frac{\lg n}{\lg 2} = \log_2 10 \lg n \approx 3.322 \lg n$.

Notăm cu n_1 numărul acelor valori x_i care sunt situate în intervalul I_1 , n_2 numărul acelor valori x_i care sunt situate în intervalul I_2 , \dots , n_r numărul acelor valori x_i care sunt situate în intervalul I_r . Histograma asociată selecției x_1, x_2, \dots, x_m este reuniunea dreptunghiurilor $D_k = I_k \times [0, n_k]$ (produs cartezian), pentru $1 \leq k \leq r$. *Moda* selecției este subgrupa (intervalul I_k) ce corespunde dreptunghiului cel mai înalt.

Deci datele selecției x_1, x_2, \dots, x_n se grupează în r grupe și există posibilitatea de vizualizare: se alege un sistem ortogonal de axe xOy , cu intervalul $[m, M]$ plasat pe axa Ox (cu originea O eventual în punctul m). Pe intervalul I_1 se construiește dreptunghiul de înălțime n_1 și baza de lungime h , pe intervalul I_2 se construiește dreptunghiul de înălțime n_2 și baza de lungime h , etc. Intervalul I_k pentru care n_k este maxim este moda selecției. Dacă două subgrupe sunt la fel de înalte și detașate de celelalte, se spune că selecția este bimodală (ex: figura 6.2).

10.2.2 Algoritmul pentru construcția histogramelor

Date inițiale: x_1, x_2, \dots, x_n .

Pasul I. Se calculează $m = \min_i x_i$, $M = \max_i x_i$, $r = [1 + 3.322 \lg n]$ și $h = \frac{M - m}{r}$.

Pasul II. Se determină intervalele $I_k = [m + (r - 1)h, m + rh)$ și frecvențele absolute n_k , $1 \leq k \leq r$.

Pasul III. Se constuiesc dreptunghiurile ce formează histograma $D_k = I_k \times [0, n_k]$, $1 \leq k \leq r$.

Reluăm Exemplul 10.2.1.

Pasul I. $M = \max_i x_i = 99$, $m = \min_i x_i = 49$, $n = 83$,

$$r = [1 + 3.322 \lg 83] = 7$$

$$h = \frac{M - m}{r} = \frac{99 - 49}{7} = 7.14286$$

Pasul II. Se determină intervalele I_k și frecvențele n_k , $1 \leq k \leq 7$.

I_k	[49, 56.1429)	[56.1429, 63.2857)	[63.2857, 70.4286)
n_k	3	8	16
I_k	[70.4286, 77.5714)	[77.5714, 84.7143)	[84.7143, 91.8571)
n_k	18	22	8
I_k	[91.8571; 99]		
n_k	7		

Pasul III. Desenăm histograma.

Algoritmul MATLAB pentru realizarea histogramei este următorul:

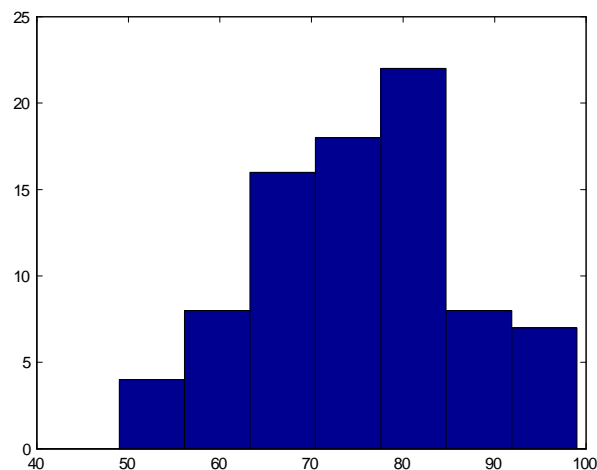
```
x=[64, 62, 76, 82, 66, 76, 72, 71, 74, 72, 71, 73, 70, 75, 77, 84, 92, 86, ...
    62, 58, 78, 80, 79, 84, 83, 82, 66, 68, 68, 82, 84, 78, 76, 69, 77, 58, ...
    62, 82, 85, 58, 78, 84, 94, 88, 77, 78, 88, 91, 70, 71, 78, 58, 65, 53, ...
```

```

60, 49, 68, 74, 71, 66, 68, 71, 73, 70, 85, 78, 65, 54, 51, 78, 89, 66, ...
68, 95, 94, 99, 81, 81, 92, 88, 99, 81, 81];
n=length(x); m=min(x); M=max(x); k=fix(1+3.322*log10(n));
h=(M-m)/k;
grupe=zeros(k,1);
for i=1:n
    if x(i)==M
        grupe(k)=grupe(k)+1;
    else
        t=fix((x(i)-m)/h)+1;
        grupe(t)=grupe(t)+1;
    end
end
grupe
xprim=zeros(k,1);
for i=1:k
    xprim(i)=m+(2*i-1)/2*h;
end
bar(xprim,grupe,1)

```

Histograma obținută este următoarea:



Exemplul 10.2.3 Presupunem că într-un atelier lucrează 225 de muncitori și că s-a analizat într-o perioadă de timp modul în care își realizează sarcinile. S-a constatat că 10 dintre ei realizează 101%, etc. conform tabelului de mai jos

10	15	40	23	17	18	36	45	8	5	8
101%	89%	85%	106%	127%	95%	105%	112%	117%	96%	80%

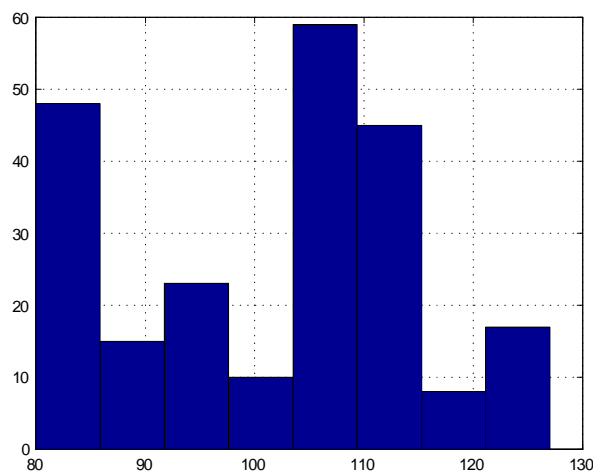
Construim histograma.

În acest caz $m = 80$, $M = 127$, $r = [1 + 3.322 \log_{10} 225] = 8$, $h = \frac{127 - 80}{8} = 5.875$.
 Vom împărți intervalul în 8 clase de lungime 5.875 fiecare

t	[80, 85.875)	[85.875, 91.75)	[91.75, 97.625)
n_j	48	15	23
t	[97.625, 103.5)	[103.5, 109.375)	[109.375, 115.25)
n_j	10	59	45
t	[115.25, 121.125)	[121.125, 127]	
n_j	8	17	

Avem $I_1 = [80, 85.4)$, $n_1 = 48$, $I_2 = [85.4, 90.8)$, $n_1 = 15$, etc. Desenăm histograma

```
x=[10 15 40 23 17 18 36 45 8 5 8];
y=[101 89 85 106 127 95 105 112 117 96 80];
n=length(y); m=min(y); M=max(y); s=sum(x);
k=fix(1+3.322*log10(s));
h=(M-m)/k;
grupe=zeros(k,1);
for i=1:n
    if y(i)==M
        grupe(k)=grupe(k)+x(i);
    else
        t=fix((y(i)-m)/h)+1
        grupe(t)=grupe(t)+x(i)
    end
end
grupe
xprim=zeros(k,1);
for i=1:k
    xprim(i)=m+(2*i-1)/2*h;
end
bar(xprim,grupe,1)
grid
```



Ce concluzii se pot trage din analiza histogramei? În primele patru grupe sunt persoanele care nu-și realizează sarcinile. Caracteristica numerică numită modă este, pentru atelierul respectiv, "să se realizeze sarcinile între 103.5% și 109.375%". 59 din cei 225 de muncitori sunt în modă.

10.3 Caracteristici statistice ale datelor experimentale

Definiția 10.3.1 Dacă avem o selecție x_1, x_2, \dots, x_n de volum n , definim media de selecție (empirică) ca fiind $\bar{x} = \frac{1}{n}(x_1 + x_2 + \dots + x_n)$, adică media aritmetică a celor n valori.

Dacă datele sunt organizate sub forma unui tabel (10.1) atunci media de selecție este $\bar{x} = \frac{1}{n} \sum_{i=1}^k n_i x_i$, unde n_i reprezintă frecvența absolută de apariție a valorii x_i a variabilei în selecția considerată și $\sum_{i=1}^r n_i = n$.

În cazul datelor grupate x_i se înlocuiește cu x_i^* care este valoarea din mijloc a intervalului $[x_i, x_{i+1})$ deci $\bar{x} = \frac{1}{n} \sum_{i=1}^r n_i x_i^*$, unde r este numărul intervalelor.

Teorema 10.3.2 Au loc relațiile: $\overline{x+c} = \bar{x} + c$, $\overline{kx} = k\bar{x}$, $\overline{x+y} = \bar{x} + \bar{y}$.

Mediana este valoarea $\alpha \in \mathbb{R}$ astfel ca numărul de valori $x_j \leq \alpha$ este egal cu numărul de valori $x_j \geq \alpha$. Deci mediana împarte șirul ordonat de date în două părți egale. Dacă șirul are $2k+1$ unități, atunci mediana coincide cu unitatea de ordin $k+1$, dacă șirul are $2k$ unități, mediana este media aritmetică a unităților de ordin k și $k+1$. Dacă există mai multe asemenea valori pentru α , atunci ele formează un interval și mediana este prin definiție mijlocul acestui interval.

Exemplul 10.3.3 Pentru șirul de date: 2.5, 3.7, 1.4, 0.2, 5.4, 8.9, 4.2 șirul ordonat este 0.2, 1.4, 2.5, 3.7, 4.2, 5.4, 8.9. Avem $2k+1 = 7$, apoi $k = 3$ și unitatea de ordin $k+1 = 4$ este 3.7 care este mediana. \diamond

Dacă $x_{me} = \bar{x}$ atunci repartiția este simetrică.

Modul selecției. Pentru datele negrupate este valoarea observată care are frecvența maximă. Pentru date grupate este o valoare din clasa cu cel mai mare număr de observații.

Dispersia de selecție (empirică) este $\overline{\sigma^2} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$.

Dispersia de selecție modificată este $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$.

Observăm că dacă s^2 este mic atunci datele diferă unele de altele foarte puțin, dacă s^2 este mare, datele diferă foarte mult.

Pentru date negrupate avem

$$\begin{aligned} s^2 &= \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} = \frac{\sum_{i=1}^n (x_i^2 - 2\bar{x}x_i + \bar{x}^2)}{n-1} = \frac{\sum_{i=1}^n x_i^2 - 2\bar{x} \sum_{i=1}^n x_i + n\bar{x}^2}{n-1} = \\ &= \frac{\sum_{i=1}^n x_i^2 - n\bar{x}^2}{n-1} = \frac{\sum_{i=1}^n x_i^2 - \frac{n^2\bar{x}^2}{n}}{n-1} = \frac{\sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2}{n-1} \\ &= \frac{\sum_{i=1}^n x_i^2}{n-1} - \frac{n}{n-1} \bar{x}^2, \end{aligned}$$

iar pentru date grupate

$$s^2 = \frac{\sum_{i=1}^r n_i x_i^{*2} - \frac{1}{n} \left(\sum_{i=1}^r n_i x_i^* \right)^2}{n-1} = \sum_{i=1}^r \frac{n_i}{n-1} (x_i^* - \bar{x})^2.$$

Motivația introducerii acestei dispersii o vom vedea ulterior în legătură cu conceptul de estimator nedeplasat.

\bar{x} se mai numește *speranța matematică* a selecției, dispersia de selecție se mai numește *varianța* selecției.

Amplitudinea unei serii statistice este diferența dintre cea mai mare și cea mai mică valoare a variabilei $\mathcal{A} = x_{\max} - x_{\min}$.

În cazul v. a. continue, amplitudinea este diferența dintre limita superioară a ultimului interval și limita inferioară a primului.

Abaterea medie pătratică de selecție este $\bar{\sigma} = \sqrt{\sigma^2}$.

Abaterea medie pătratică de selecție modificată este $s = \sqrt{s^2}$.

Coefficientul de variație este $v = \frac{\sqrt{s^2}}{\bar{x}}$. O selecție de valori pozitive se consideră omogenă dacă $v < \frac{1}{3}$.