

Capitolul 2

Populație și caracteristici. Eșantion

2.1 Definiții și clasificări

Pentru început, vom introduce noțiunile (termenii) de bază care se vor regăsi în toate etapele unui studiu statistic.

Def. Se numește *populație (colectivitate) statistică* orice mulțime nevidă de elemente supusă studiului în legătură cu un anumit fenomen. Elementele sale se mai numesc *indivizi* sau *unități statistice*.

Populația statistică o vom nota cu Ω , iar elementele sale cu $\omega_1, \omega_2, \dots$

Obs. 1 Din punct de vedere matematic, populația este modelată prin intermediul noțiunii de mulțime.

Def. *Caracteristica statistică* reprezintă însușirea, proprietatea sau trăsătura comună tuturor indivizilor unei populații, și care stă la baza studiului acelei populații. Se notează cu litere mari, X, Y, \dots și se mai numește *variabilă statistică* sau *variabilă aleatoare*.

Exemple

- a) Pentru populația alcătuită din studenții unei facultăți, se pot studia caracteristici precum: sex, vârstă, stare civilă, notele la examene, media, etc.

- b) Dacă populația este reprezentată de un anumit bun realizat în serie, atunci caracteristicile pot fi: calitatea (control de calitate), dimensiunea, prețul, etc.
- c) În cazul în care populația este constituită din bunurile desfăcute pe piață de o firmă, caracteristicile considerate vor fi: sortimentele desfăcute, cantitatea din fiecare tip de bun, prețul la desfacere, etc.
- d) Populația dintr-o anumită zonă geografică poate fi studiată din punct de vedere al: sex, vârstă, număr de ani de școlarizare, stare civilă, grupa socio-profesională, stare de sănătate, culoare ochi, culoare păr, etc.

Obs. 2 Pentru fiecare populație se poate lua în studiu un număr mai mare sau mai mic de caracteristici, în funcție de cerințele analizei. Numărul caracteristicilor poate fi limitat din considerente de eficiență și de raționalitate.

Caracteristicile determină descompunerea populației în grupe omogene, care pot fi studiate separat. Fiecare caracteristică are mai multe forme de manifestare numite și *modalități* sau *valori*, care trebuie să satisfacă următoarele principii:

- *principiul completitudinii*: fiecare individ al populației aparține unei clase definite de una dintre modalități;
- *principiul unicității*: un individ aparține unei singure clase;
- *principiul organizării ierarhice a claselor*: clasele pot fi unificate prin mărirea gradului de generalitate al modalităților.

Vom nota modalitățile cu litere mici: x, y, \dots

Obs. 3 Din cele spuse mai sus rezultă că variabila statistică poate fi modelată matematic prin intermediul noțiunii de funcție, având ca domeniu de definiție mulțimea Ω , iar ca domeniu de valori o mulțime nevidă, de orice natură (numerică sau nu).

Clasificarea caracteristicilor se poate face după mai multe criterii. Prezentăm două criterii:

- a) Clasificare după *forma de prezentare*:

- caracte
no
pri

- caracte
cu

Obs
tica de c
cel canti
sexul as
tativ la
studiată

b) Clasi

- caracte
xe

- caracte
ce

Del
Cardin
este alc
valoare

2.2

Așa cu
duri. E
scala d

2.2.1

Valoril
sau sir
mate f

alizat în serie,
ate), dimensi-

- *caracteristici cantitative*: valorile lor sunt exprimate numeric. Exemple: nota, vârsta, greutatea, numărul de aruncări ale unei monede până la prima apariție a stemei, etc.

ăcute pe piață
tele desfăcute,

- *caracteristici calitative*: valorile lor nu se exprimă numeric. Exemple: culoare ochi, profesie, sex.

ă din punct de
civilă, grupa
re păr, etc.

Obs. 4 O caracteristică de un tip se poate transforma într-o caracteristică de celălalt tip. Cea mai frecventă transformare este din tipul calitativ în cel cantitativ, prin numerotarea modalităților. Spre exemplu, putem codifica sexul astfel: 1-masculin, 2-feminin. Invers, pentru a trece de la tipul cantitativ la cel calitativ, se pot defini grupe tipice de valori pentru caracteristica studiată.

in număr mai
izei. Numărul
i de raționali-

- b) Clasificare după *numărul modalităților*:

upe omogene,
ulte forme de
atisfacă urmă-

- *caracteristici discrete*: iau o mulțime cel mult numărabilă de valori. Exemple: sex, nota.

ine unei clase

- *caracteristici de tip continuu*: iau valori într-un interval; pot fi doar de tip cantitativ. Exemplu: dimensiunile unei piese.

Def. O submulțime finită $\Omega_1 \subseteq \Omega$ se numește *eșantion* sau *selecție*. Cardinalul său poartă numele de *volum de selecție*. În general, *eșantionul* este alcătuit din acei indivizi ai populației pentru care s-a observat efectiv valoarea caracteristicii studiate.

unificate prin

2.2 Elemente de teoria scalării

Așa cum s-a văzut mai sus, caracteristicile pot fi măsurate în mai multe moduri. Există patru tipuri principale de scală: scala nominală, scala ordinală, scala de interval și cea de raport.

2.2.1 Scala nominală

Valorile unei caracteristici măsurată pe *scala nominală* pot fi nume, numere sau simboluri, care însă nu se pot ordona. Ca urmare, caracteristicile exprimate în această scală sunt de tip calitativ și trebuie să îndeplinească condițiile:

stică poate fi
d ca domeniu
vidă, de orice

nulte criterii.

1. Valorile scalei nominale sunt exhaustive și mutual exclusive; fiecare observație trebuie să corespundă unei singure valori. Indivizii care au aceeași valoare se consideră echivalenți în raport cu caracteristica studiată.
2. Numele sau simbolurile care desemnează valorile pot fi interschimbate fără a altera informațiile transmise de scală.

Exemple. Grupa sanguină (O, A, B, AB); culoarea ochilor, sex, etc.

2.2.2 Scala ordinală

Valorile unei caracteristici măsurată pe *scala ordinală* pot fi ordonate după un anumit criteriu, însă nu permit efectuarea de operații aritmetice. Caracteristicile exprimate în această scală sunt tot de tip calitativ și verifică condițiile:

1. Diferențele dintre valori nu sunt neapărat egale și pot fi chiar imposibile de măsurat.
2. Simbolurile asociate valorilor nu sunt importante atâta timp cât relația de ordine este păstrată.

Exemplu. Starea unui pacient poate fi clasificată ca înrăutățită, stabilă sau îmbunătățită. Referitor la prima condiție de mai sus, se observă că diferența dintre starea "înrautățită" și "stabilă" nu este neapărat aceeași ca cea dintre "stabilă" și "îmbunătățită".

2.2.3 Scala de interval

Valorile unei caracteristici măsurată pe *scala de interval* sunt numere echidistante. Scala nu are un "punct zero" (o origine care să reprezinte absența caracteristicii măsurate). Cu ajutorul acestei scale se poate determina cât de mult (sau cât de puțin) din caracteristica măsurată reprezintă fiecare valoare. Sunt permise operații de adunare și de scădere între valori. Caracteristicile exprimate în această scală sunt de tip cantitativ.

Exemplu. Temperatura exprimată în grade Celsius (observați că valoarea 0° este cea la care îngheață apa și nu reprezintă absența temperaturii).

2.3. SU

2.2.4

Pe *scale*
sunt ech
exprima

Ex
volum

2.3

Fluctu
binare
pot fi

- var

- var

- etc

2.2.4 Scala de raport

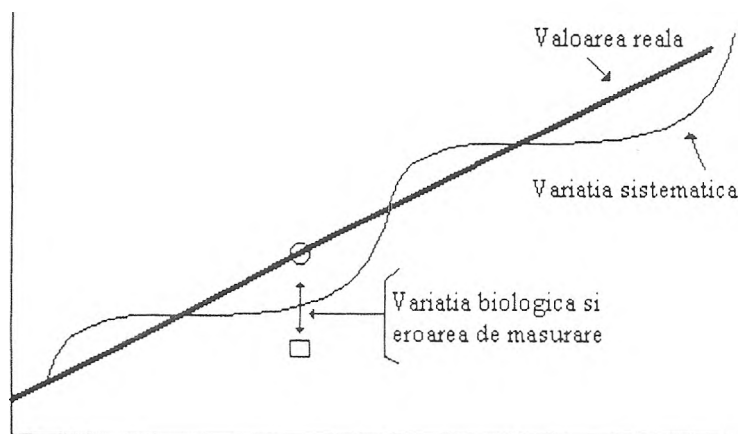
Pe *scala de raport*, măsurătorile încep dintr-un punct zero (origine) și valorile sunt echidistante. Sunt permise toate operațiile aritmetice, iar caracteristicile exprimate în această scală sunt de tip cantitativ.

Exemple. Temperatura exprimată în grade Kelvin, lungimea, greutatea, volumul etc.

2.3 Sursele variațiilor din date

Fluctuațiile ce pot apărea în timpul culegerii datelor reflectă în general combinarea efectelor mai multor fenomene. Spre exemplu, măsurătorile clinice pot fi afectate de:

- *variația biologică* (exemplu: valoarea tensiunii unui același pacient, în condiții identice, poate să difere);
- *variația sistematică* (exemplu: valoarea tensiunii unui pacient diferă în funcție de momentul zilei și de poziția în care stă);
- *eroarea de măsurare*, care poate fi atât aleatoare, cât și sistematică (spre exemplu, datorită calibrării instrumentelor).



- Valoarea reala
- Valoarea masurata, influentata de variatii

I. EȘANTION

tatea specialis-
le.

Capitolul 3

Tabele de distribuție. Gruparea datelor

Având în vedere observația 4 din secțiunea anterioară, în cele ce urmează vom considera numai caracteristici statistice cantitative.

3.1 Serii statistice (date brute)

Forma cea mai simplă de prezentare a datelor statistice provenind dintr-o singură variabilă statistică este prin enumerarea observațiilor efectuate asupra variabilei: x_1, x_2, \dots, x_n . Această enumerare se numește *serie statistică* sau *eșantion*, n fiind volumul său. Unele valori ale eșantionului se pot repeta.

Dacă datele din eșantion se ordonează crescător, rezultă un *șir variațional*, $x_{(1)}, x_{(2)}, \dots, x_{(n)}$, unde

$$\begin{cases} 1) \{x_1, x_2, \dots, x_n\} = \{x_{(1)}, x_{(2)}, \dots, x_{(n)}\} \\ 2) x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)} \end{cases}$$

Exemplu. S-a măsurat înălțimea (în cm) a 20 de copii dintr-o școală, obținându-se următoarea serie statistică:

133	136	120	138	133	131	127	141	127	143
130	131	125	144	128	134	135	137	133	129

Șirul variațional corespunzător este:

120	125	127	127	128	129	130	131	131	133
133	133	134	135	136	137	138	141	143	144

Obs. 1 Deoarece numărul datelor dintr-o serie statistică poate fi foarte mare (de ordinul miilor), aceasta devine greu de prelucrat și de studiat. În urmare, datele se prezintă și sub formă de tabele de distribuție.

Uneori se
numește c
Relații

3.2 Tabele de distribuție. Tipuri de frecvențe

Obs. 2 C
probabilit
a i-a valor

3.2.1 Tabel de distribuție a frecvențelor absolute

Acest tip de tabel se obține dintr-o serie statistică astfel: dacă (x_1, \dots, x_n) este un eșantion de volum n de observații făcute asupra variabilei statistice X , atunci:

3.2.3

1. Se formează șirul variațional: $x_{(1)}, x_{(2)}, \dots, x_{(n)}$.
2. Pe o coloană (linie) se trec doar valorile distincte din șirul variațional, fie ele $x'_{(1)}, x'_{(2)}, \dots, x'_{(k)}$, $k \leq n$.
3. Pe a doua coloană (linie), în dreptul fiecărei valori, se trece numărul său de apariții din eșantionul inițial: n_1, \dots, n_k . Acest număr de apariții al unei valori observate $x'_{(i)}$ se numește *frecvență absolută*, notată n_i .

În scopul
statistice,
prin adun.

Def.
corespunz
variabilei

Tabelul astfel obținut se numește *tabel de distribuție a frecvențelor absolute* sau *distribuție empirică (statistică, observată)* sau *seria statistică a frecvențelor absolute*.

Este ușor de observat că, dacă avem în total k valori distincte în șirul variațional, frecvențele absolute verifică relația

$$n_1 + \dots + n_k = n. \quad (3.1)$$

Def.
corespunz
ator valori
nume

3.2.2 Tabel de distribuție a frecvențelor relative

Frecvențele se pot exprima și în valori relative la numărul total de observații. Aceste frecvențe se numesc *relative* și sunt definite prin

Obs. 3 Se
frecvența

Exem
mătorul

$$f_n(X = x'_{(i)}) = n'_i = \frac{n_i}{n}$$

133

144

Uneori se exprimă în procente, sub forma $100 \frac{n_i}{n} \%$. Tabelul corespunzător se numește *al frecvențelor relative* sau *seria statistică a frecvențelor relative*.

Relația (3.1) devine în acest caz

$$n'_1 + \dots + n'_k = 1.$$

Obs. 2 Când n este foarte mare, $f_n(X = x'_{(i)}) = n'_i \approx P(X = x'_{(i)})$, adică probabilitatea ca variabila statistică X rezultată în urma observațiilor să ia a i -a valoare distinctă din șirul variațional.

3.2.3 Frecvențe cumulate crescător și descrescător

În scopul reprezentării cât mai compacte, inclusiv în formă grafică, a datelor statistice, în statistică sunt utilizate frecvențele cumulate. Acestea se obțin prin adunarea din aproape în aproape a frecvențelor absolute sau relative.

Def. Se numește *frecvență absolută cumulată crescător (descrescător)* corespunzătoare unei valori x , suma frecvențelor absolute ale tuturor valorilor variabilei mai mici sau egale cu x (mai mari sau egale cu x), anume

$$\sum_{y \leq x} n_y \left(\sum_{y \geq x} n_y \right).$$

Def. Se numește *frecvență relativă cumulată crescător (descrescător)* corespunzătoare unei valori x a variabilei, suma frecvențelor relative ale tuturor valorilor variabilei mai mici sau egale cu x (mai mari sau egale cu x), anume

$$\sum_{y \leq x} n'_y \left(\sum_{y \geq x} n'_y \right).$$

Obs. 3 Se observă ușor că frecvența relativă cumulată este raportul dintre frecvența absolută cumulată și volumul populației.

Exemplu. Reluând datele de mai sus (înălțimile celor 20 de copii), avem următorul tabel cu frecvențele absolute, relative și cumulate:

Înăl- țimea	Frecv. abs.	Frecv. relat.	Frec. abs. cumul. cresc.	Frec. abs. cumul. descresc.	Frec. rel. cumul. cresc.	Frecv. cumul. descresc.
120	1	0.05	1	20	0.05	0.05
125	1	0.05	2	19	0.1	0.05
127	2	0.1	4	18	0.2	0.05
128	1	0.05	5	16	0.25	0.05
129	1	0.05	6	15	0.3	0.05
130	1	0.05	7	14	0.35	0.05
131	2	0.1	9	13	0.45	0.05
133	3	0.15	12	11	0.6	0.05
134	1	0.05	13	8	0.65	0.05
135	1	0.05	14	7	0.7	0.05
136	1	0.05	15	6	0.75	0.05
137	1	0.05	16	5	0.8	0.05
138	1	0.05	17	4	0.85	0.05
141	1	0.05	18	3	0.9	0.05
143	1	0.05	19	2	0.95	0.05
144	1	0.05	20	1	1	0.05

3.2.5 Ta

Dacă analiza
se trec într-c
Exempl
parului:

C

3.2.4 Tabele de distribuție grupate

Pentru un volum n foarte mare de date distincte, se recomandă gruparea lor în clase (intervale de valori), sub forma:

Din ace
fiecărei car
numeric, a

Intervalul	Frecv. abs.	Frecv. relat.
$[y_1, y_2)$	n_1	n'_1
$[y_2, y_3)$	n_2	n'_2
...
$[y_s, y_{s+1})$	n_s	n'_s
\sum	n	1

Culo
neg
cast
blo

Valoarea centrală (mijlocul) unei clase este, pentru clasa $[y_i, y_{i+1})$,

$$c_i = \frac{y_i + y_{i+1}}{2}.$$

Exemplu. Să grupăm datele de mai sus în intervale de lungime 5:

Frecv. rel.	Frecv. cumulat.	Frecv. descresc.
0.05	1	
0.1	0.95	
0.2	0.9	
0.25	0.8	
0.3	0.75	
0.35	0.7	
0.45	0.65	
0.6	0.55	
0.65	0.4	
0.7	0.35	
0.75	0.3	
0.8	0.25	
0.85	0.2	
0.9	0.15	
0.95	0.1	
1	0.05	

Intervalul	Frecv. abs.	Frecv. relat.
(120, 125)	1	0.05
(125, 130)	5	0.25
(130, 135)	7	0.35
(135, 140)	4	0.2
(140, 145)	3	0.15
Σ	20	1

3.2.5 Tabele de distribuție pentru două caracteristici

Dacă analiza statistică se face după două caracteristici, rezultatele obținute se trec într-o tabelă cu două intrări.

Exemplu. Clasificarea a 1500 de persoane după culoarea ochilor și a părului:

Culoare păr	Culoare ochi				Σ
	negri	căprui	verzi	albaștri	
negru	145	285	30	11	471
castaniu	62	431	87	67	647
blond	33	36	185	128	382
Σ	240	752	302	206	1500

Din acest tip de tabel sunt ușor de obținut tabelele de distribuție ale fiecărei caracteristici în parte, denumite și *distribuții marginale*. Pe exemplul numeric, avem:

Culoare păr	Frecv. abs.
negru	471
castaniu	647
blond	382
Σ	1500

Culoare ochi	Frecv. abs.
negri	240
căprui	752
verzi	302
albaștri	206
Σ	1500

, y_{i+1}),

Capitolul 4

Reprezentări grafice

4.1 Diagramele frecvențelor necumulate

Se pleacă de la seria frecvențelor absolute sau relative și se pot trasa trei tipuri de diagrame: diagrama cu bastonașe, histograma (diagrama cu dreptunghiuri) și poligonul frecvențelor. Pentru toate cele trei tipuri de diagrame, pe axa OX se trec valorile observate, iar dacă se lucrează cu date grupate, se trec limitele intervalelor sau centrele lor.

4.1.1 Diagrama cu bastonașe

Acest tip de grafic este foarte bun pentru datele negrupate. Astfel, pentru fiecare valoare observată distinctă se trasează un segment vertical de lungime egală cu frecvența corespunzătoare valorii respective.

4.1.2 Histograma

Se aplică pentru datele grupate în intervale. Este alcătuită din dreptunghiuri având ca bază intervalele de valori, iar aria lor este proporțională cu frecvențele. În funcție de lungimea intervalelor, avem două cazuri:

1. *Histograme pentru intervale de lungimi egale*, caz în care înălțimea fiecărui dreptunghi coincide cu frecvența intervalului care stă la baza sa.
2. *Histograme pentru intervale de lungimi inegale*; în acest caz, pentru calcularea înălțimii dreptunghiurilor se aplică următorul algoritm:

- pe o coloană separată se trec lungimile tuturor intervalelor;
- se alege o lungime standard (cea care apare mai frecvent sau care este mai mică);
- pentru un interval oarecare, dacă lungimea sa este egală cu $(\alpha \times \text{lungime standard})$, atunci înălțimea dreptunghiului corespunzător va fi $\left(\frac{1}{\alpha} \times \text{frecvența corespunzătoare aceluia interval}\right)$.

4.1.3 P

Se obține u
din diagra
tunghiurilo

Exemplul 1. Considerăm datele anterioare (înălțimile celor 20 de copii) grupate în clase de lungime 5. Histograma este cea din fig. 1.

4.2 I

Exemplul 2. Considerăm masele a 35 de obiecte, măsurate în kg. Datele sunt grupate ca în tabelul de mai jos, iar în fig. 2 este trasată histograma corespunzătoare.

Frecvențele
sau prin p

Masa	Frecvență	Lungime interval	Înălțime dreptunghi
[6, 9)	4	3 standard	4
[9, 12)	6	3 standard	6
[12, 18)	10	6 $2 \times$ standard	$10/2 = 5$
[18, 21)	3	3 standard	3
[21, 30)	12	9 $3 \times$ standard	$12/3 = 4$

4.3 I

Ariile sect
frecvențele

Exem
tate A. B.

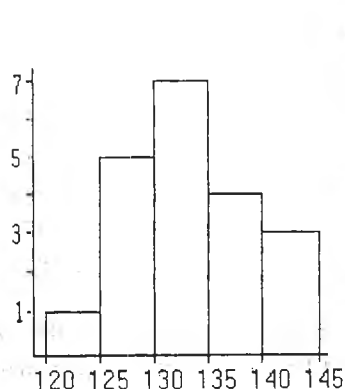


Fig. 1

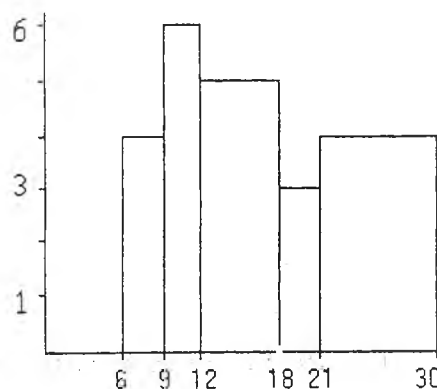


Fig. 2

tervalelor;

frecvent sau care

gală cu $(\alpha \times \text{lungime})$
 răspunzător va fi

4.1.3 Poligonul frecvențelor

Se obține unind prin linii (frânte) extremitățile superioare ale bastonașelor din diagrama cu bastonașe sau mijloacele extremităților superioare ale dreptunghiurilor din histogramă.

4.2 Diagramele frecvențelor cumulate

Frecvențele cumulate se reprezintă ca mai înainte, dar numai prin histograme sau prin poligoanele frecvențelor.

4.3 Diagrame circulare ("pie" sau "plăcintă")

Ariile sectoarelor de cerc dintr-o asemenea diagramă sunt proporționale cu frecvențele.

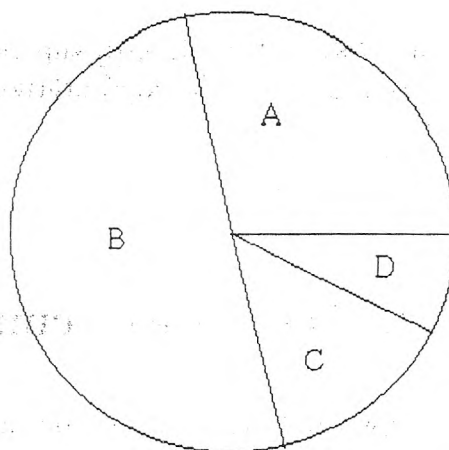
Exemplu. Considerăm vânzările de benzină de la patru benzinării, notate A, B, C, D.

Benzi- nărie	Vânzări (sute litri)	Unghiul sectorului
A	90	$360^\circ \cdot 90/280 = 115.7^\circ$
B	140	$360^\circ \cdot 140/280 = 180^\circ$
C	30	$360^\circ \cdot 30/280 = 38.6^\circ$
D	20	$360^\circ \cdot 20/280 = 25.7^\circ$
Σ	280	

dreptunghi

$$2 = 5$$

$$3 = 4$$



Aceste diagrame se pot utiliza pentru a compara două sau mai multe seturi de date similare. Spre exemplu, dacă se reprezintă structura culturilor agricole pentru un teren într-o diagramă circulară, putem compara diagramele a două zone diferite.

Tipul culturii	Suprafață (ha)	Procent (%)
Grâu	120	40
Porumb	80	26,7
Soia	60	20
Alte culturi	40	13,3

Cap
Val

5.1

Parametre
valorile
parametre

5.1.1

Este valoarea
modul
num:

pentru

pentru
r

Capitolul 5

Valori caracteristice

5.1 Parametri de poziție

ia sau mai multe
ită structura cul-
, putem compara

Parametrii de poziție sunt valori de referință la care se raportează toate valorile unui eșantion. În cele ce urmează vom prezenta principalii asemenea parametri, anume media, mediana și moda (modul).

5.1.1 Media (de selecție), \bar{x}

Este vorba despre media aritmetică obișnuită. Calculul său depinde însă de modul de grupare al datelor. Astfel, cu notațiile din secțiunile anterioare, avem:

• pentru un eșantion (date brute) x_1, x_2, \dots, x_n , media este

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i ;$$

• pentru datele dintr-un tabel de distribuție a frecvențelor absolute sau relative, media se calculează prin

$$\bar{x} = \frac{1}{n} \sum_{i=1}^k x'_{(i)} n_i = \sum_{i=1}^k x'_{(i)} n'_i ;$$

- pentru datele grupate în intervale,

$$\bar{x} = \frac{1}{n} \sum_{i=1}^k c_i n_i.$$

Reamintim și definiția *mediei ponderate* a eșantionului x_1, x_2, \dots, x_n , cu ponderile $f_i > 0, i = \overline{1, n}$,

$$\bar{x}_f = \frac{\sum_{i=1}^n f_i x_i}{\sum_{i=1}^n f_i}.$$

Interpretare. \bar{x} este valoarea în jurul căreia se grupează valorile caracteristicii studiate.

Proprietate.

a) Media verifică relația

$$x_{(1)} \leq \bar{x} \leq x_{(n)}.$$

b) Dacă valorile x_i din eșantionul inițial sunt transformate în $y_i = a + bx_i$, atunci $\bar{y} = a + b\bar{x}$.

c) Dacă se combină (concatenează, comasează) două eșantioane de volume m, n și de medii \bar{x} , respectiv \bar{y} , atunci media eșantionului rezultat este

$$\bar{z} = \frac{m\bar{x} + n\bar{y}}{m + n}.$$

5.1.2 Mediana

Mediana este acea valoare care are proprietatea că jumătate dintre observațiile din eșantion sunt mai mici sau egale cu ea, cealaltă jumătate dintre observații fiind mai mari sau egale decât ea. Cu alte cuvinte, mediana este acea valoare numerică ce împarte șirul variațional corespunzător eșantionului dat în două părți egale. Calculul său se face tot în funcție de modul de grupare al datelor:

- pentru un eșantion x_1, x_2, \dots, x_n , mediana se calculează astfel:

1. se ordonează crescător eșantionul;

2. dacă n este impar, atunci mediana este observația de rang $\frac{n+1}{2}$, iar dacă n este par, mediana este $(x_{\frac{n}{2}} + x_{\frac{n}{2}+1}) / 2$.

Exer
36, 41, 27

pentru
Se ci put
Mediane

1. se

2. se

3. se

la

fi

es

O a
mare. (

variație

pentru

1. s

2. p

c

Pr

5.1.3

Moda

Secve

- pe

Exemplu. Pentru eșantionul 7, 7, 2, 3, 4, mediana este 4, iar pentru 16, 41, 27, 32, este $(32 + 36) / 2 = 34$.

pentru datele dintr-un tabel de distribuție: presupunem că se cunoaște n , deci putem afla frecvențele absolute (inclusiv pe baza frecvențelor relative). Mediana rezultă astfel:

x_1, x_2, \dots, x_n , cu

1. se calculează frecvențele absolute cumulate crescător;
2. se calculează $(n + 1) / 2$;
3. se determină prima linie din tabel pentru care frecvența absolută cumulată crescător este $\geq (n + 1) / 2$, valoarea similară de pe linia precedentă fiind $< (n + 1) / 2$. Valoarea caracteristicii corespunzătoare acestei linii este mediana.

valorile caracte-

O altă metodă care se poate aplica în acest caz, dacă n nu este prea mare, constă în scrierea valorilor caracteristicii în ordine crescătoare (șirul variațional) până la cea de-a $(n + 1) / 2$ -a valoare, care va fi chiar mediana.

în $y_i = a + bx_i$

ioane de volume
rezultat este

pentru datele grupate în intervale, aflarea mediane se face astfel:

1. se trasează grafic curba frecvențelor absolute cumulate crescător;
2. pentru $y = (n + 1) / 2$, se află prin interpolare x -ul corespunzător de pe curba trasată. Acesta este mediana căutată.

Proprietate. Fie $\alpha \in \{x'_{(1)}, x'_{(2)}, \dots, x'_{(k)}\}$. Atunci

$$\min_{\alpha} \left(\sum_{i=1}^k |x'_{(i)} - \alpha| n_i \right) = \sum_{i=1}^k |x'_{(i)} - \text{mediana}| n_i.$$

ate dintre obser
i jumătate dintre
ite, mediana este
nzător eșantion
cție de modul de

5.1.3 Moda (modul)

Moda (modul) este valoarea din eșantion căreia îi corespunde cea mai mare frecvență absolută (relativă).

- pentru un eșantion (date brute) x_1, x_2, \dots, x_n , se construiește un tabel de distribuție și se procedează ca mai jos;

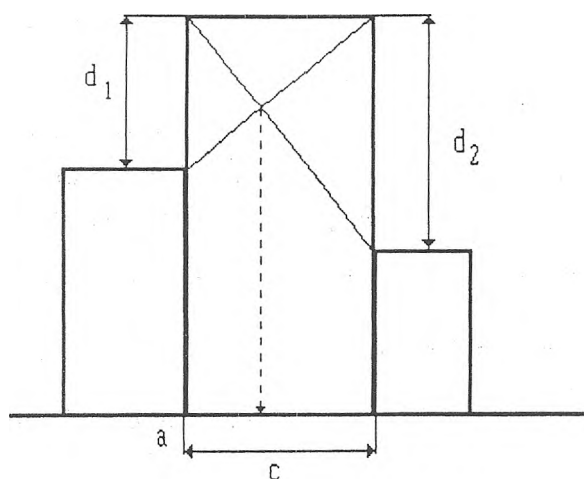
e rang $\frac{n+1}{2}$

- pentru datele dintr-un tabel de distribuție a frecvențelor absolute sau relative, se determină $n_x = \max_y n_y$, x fiind chiar moda;

- pentru datele grupate în intervale, se parcurg pașii:

1. se trasează histograma;
2. se alege cel mai înalt dreptunghi din histogramă;

3. cu notațiile din figură, $\text{moda} = a + \frac{d_1}{d_1 + d_2} \cdot c$.



5.1.4 Cuartile și percentile

Se numesc *cuartile* cele trei valori care împart șirul variațional în patru părți egale. Cele 99 de valori care împart șirul variațional într-o 100 de părți egale se numesc *percentile*. Astfel, considerând n elemente aranjate crescător, avem:

Cuartila inferioară	Q_1	valoarea a $\frac{1}{4}(n+1)$ -a
Mediana	Q_2	valoarea a $\frac{1}{2}(n+1)$ -a
Cuartila superioară	Q_3	valoarea a $\frac{3}{4}(n+1)$ -a
A 10-a percentilă	P_{10}	valoarea a $\frac{10}{100}(n+1)$ -a
A 90-a percentilă	P_{90}	valoarea a $\frac{90}{100}(n+1)$ -a

etc.

lor absolute sau

Definim și *distanța intercuartile* $= Q_3 - Q_1$, iar *semi-distanța intercuartile* $= (Q_3 - Q_1) / 2$. Avantajul acestor distanțe este acela că ele depind exclusiv de valorile din mijlocul șirului variațional, nefiind afectate de valorile extreme.

Aflarea cuartilelor și percentilelor se face similar cu aflarea medianei.

Exemplu. Să se afle cuartilele Q_1, Q_3 și semi-distanța intercuartilelor pentru următorul eșantion:

3, 2, 3, 9, 6, 6, 12, 11, 8, 2, 3, 5, 7, 5, 4, 4, 5, 12, 9

Soluție. Avem 19 numere, pe care le ordonăm crescător:

2, 2, 3, 3, 3, 4, 4, 5, 5, 5, 6, 6, 7, 8, 9, 9, 11, 12, 12

Cuartila inferioară este valoarea a 5-a, deci $Q_1 = 3$. Quartila superioară este valoarea a 15-a, deci $Q_3 = 9$. Rezultă că semi-distanța intercuartilelor $= (Q_3 - Q_1) / 2 = 3$. □

5.2 Parametri de împrăștiere

Acești parametri măsoară gradul de împrăștiere al valorilor observate față de parametrul de poziție dat.

5.2.1 Amplitudinea

Amplitudinea absolută este diferența dintre valoarea cea mai mare și valoarea cea mai mică din eșantion, deci

$$A = x_{\max} - x_{\min} = x_{(n)} - x_{(1)}.$$

Amplitudinea relativă este raportul dintre amplitudinea absolută și medie, adică A/\bar{x} .

Amplitudinea nu se definește pentru datele grupate în intervale.

Obs. Spre deosebire de amplitudinea absolută, amplitudinea relativă poate fi utilizată pentru compararea gradului de împrăștiere a două eșantioane diferite, chiar și atunci când valorile acestora au unități de măsură diferite.

ional în patru părți
ntr-o 100 de părți
aranjate crescător.

-a	etc.
-a	
-a	
1)-a	
1)-a	

5.2.2 Abateri medii

Acestea pot fi de mai multe tipuri, în funcție de valoarea de referință:

- *abaterea medie liniară absolută față de medie* este, pentru eșantionul x_1, \dots, x_n ,

$$\frac{1}{n} \sum_{i=1}^n |x_i - \bar{x}|;$$

- *abaterea medie liniară absolută față de mediană* este

$$\frac{1}{n} \sum_{i=1}^n |x_i - \text{mediana}|;$$

- *abaterea medie liniară absolută față de modă* este

$$\frac{1}{n} \sum_{i=1}^n |x_i - \text{moda}|.$$

5.2.3 Dispersia și abaterea medie pătratică

Dispersia unui eșantion se notează cu s^2 (uneori cu σ^2) și se calculează astfel:

- pentru un eșantion (date brute) x_1, x_2, \dots, x_n ,

$$s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2;$$

- pentru datele dintr-un tabel de distribuție a frecvențelor absolute sau relative,

$$s^2 = \frac{1}{n} \sum_{i=1}^k [x'_{(i)} - \bar{x}]^2 n_i = \sum_{i=1}^k [x'_{(i)} - \bar{x}]^2 n'_i;$$

- pentru datele grupate în intervale,

$$s^2 = \frac{1}{n} \sum_{i=1}^k (c_i - \bar{x})^2 n_i.$$

Interpretare. Dispersia măsoară gradul de împrăștiere al valorilor observate în jurul mediei. Cu cât dispersia este mai mare, cu atât împrăștierea valorilor din eșantion față de medie este mai mare, și reciproc.

Abaterea medie pătratică se definește ca fiind radical din dispersie, deci

$$s = \sqrt{s^2}.$$

Interpretare. Abaterea medie pătratică este o măsură foarte utilă a gradului de împrăștiere. Pentru cele mai multe distribuții, majoritatea observațiilor se situează în intervalul $(\bar{x} - 3s, \bar{x} + 3s)$ (regula celor "3σ").

Proprietăți.

a. Dacă $x_1 = x_2 = \dots = x_n$ atunci $s^2 = 0$.

b. O formulă alternativă pentru calculul dispersiei este

$$s^2 = \overline{x^2} - (\bar{x})^2.$$

c. Dacă valorile x_i din eșantionul inițial sunt transformate în $y_i = a + bx_i$, atunci

$$s_y^2 = b^2 s_x^2.$$

d. Considerând funcția

$$f(a) = \frac{1}{n} \sum_{i=1}^n (x_i - a)^2,$$

atunci

$$\min_{a \in \mathbb{R}} f(a) = f(\bar{x}) = s^2.$$

e. Dacă se combină două eșantioane de volume n, m , de medii \bar{x}', \bar{x}'' , și de dispersii s_1^2, s_2^2 , atunci media eșantionului rezultat o notăm cu \bar{x} (formula sa este dată la media de selecție), iar dispersia este

$$s^2 = \frac{ns_1^2 + ms_2^2}{n+m} + \frac{n(\bar{x}' - \bar{x})^2 + m(\bar{x}'' - \bar{x})^2}{n+m}.$$

Primul termen al sumei reprezintă dispersia în interiorul "grupelor" (eșantioanelor), iar al doilea termen reprezintă dispersia dintre "grupe" (eșantioane).

Capitolul 7

Eșantionul. Repartiții de selecție

În acest capitol, vom prezenta mai întâi câteva elemente despre selecție și eșantion, în special selecția aleatoare simplă. Prin studiul principalelor caracteristici de selecție, vom introduce și conceptul de repartiție de selecție.

7.1 Noțiuni de selecție

7.1.1 Populația și eșantionul

Am văzut anterior că prin *populație* se înțelege în general mulțimea de indivizi care se studiază, iar *eșantionul* este o submulțime a acestei mulțimi, alcătuită din indivizii care s-au observat efectiv în timpul unui studiu neexhaustiv.

Ca exemple de populații amintim populația umană dintr-o anumită regiune geografică, o mulțime de bunuri etc. În astfel de cazuri, conceptul de populație este ușor de înțeles, fiind vorba despre *populații finite*, ai căror indivizi pot fi numărați fără dificultate. Noțiunea de populație nu este însă întotdeauna atât de simplă. Există numeroase populații, adesea ipotetice, care sunt *infinite*. Spre exemplu, insectele ce trăiesc într-o anumită regiune reprezintă o populație greu de numărat; mulțimea tuturor experimentelor realizabile de un anumit tip sau mulțimea tuturor variantelor unui joc sunt două asemenea exemple.

De altfel, în general, termenii de populație și de eșantion sunt utilizați referitor la mulțimi de măsuri și nu la mulțimile de indivizi asupra cărora aceste

măsurile au fost aplicate. Se consideră, spre exemplu, populația înălțimilor sau vârstelor locuitorilor unei anumite regiuni, precizând deci care este caracteristica particulară care prezintă interes în legătură cu indivizii studiate.

7.1.2 Selecția aleatoare simplă

Pentru a cunoaște în totalitate caracteristicile unei anumite populații, ar trebui ca fiecare element al populației respective să poată fi studiat. O asemenea operație se mai numește și *recensământ*. Din păcate, o asemenea situație este foarte rară, devenind imposibilă în cazul populațiilor numeroase (din motive precum costul și rapiditatea culegerii datelor) sau infinite. Ca urmare, se recurge la studiul unui eșantion reprezentativ pentru populație.

Def. Prin *selecție* sau *sondaj* vom înțelege mulțimea operațiilor de alegere sau selectare, dintr-o populație, a indivizilor care vor constitui eșantionul.

Există mai multe metode de construire a unui eșantion, dintre care ne oprim asupra selecției aleatoare simple.

Def. Un eșantion se numește *aleator* când probabilitatea ca un individ al populației să facă parte din eșantion este aceeași, indiferent de individ. Eșantionul se numește *aleator și simplu* dacă este aleator și, în plus, selecțiile indivizilor ce vor constitui eșantionul se fac independent una de alta.

Unul dintre procedeele cele mai utilizate pentru obținerea unor asemenea eșantioane din populații finite sau infinite, este *metoda numerelor aleatoare*.

Pentru a asigura caracterul aleator și simplu al unui eșantion cu ajutorul numerelor aleatoare, se asociază mai întâi un grup de cifre (un număr) fiecărui individ din populație. Se generează apoi un șir de numere aleatoare sau se ia din tabelele de numere aleatoare, pornind de la o anumită valoare, și se includ în eșantion indivizii ai căror numere apar în șir. În cazul populațiilor finite trebuie procedat cu grijă, astfel încât să nu se selecteze un același individ de mai multe ori (selecțiile trebuie efectuate "fără revenire").

7.2 Caracteristici de selecție

Pentru a studia valorile unui eșantion sunt necesare cel puțin două caracteristici: o valoare centrală și o măsură a împrăștiilor valorilor eșantionului în jurul acestei valori centrale.

populația înălțimilor
d deci care este ca-
u indivizii studiați.

te populații, ar tre-
tudiat. O asemenea
emenea situație este
eroase (din motive
ite. Ca urmare, se
ulație.

mea operațiilor de
vor constitui eșan-

ion, dintre care ne

tatea ca un individ
diferent de individ.
și, în plus, selecțiile
una de alta.

erea unor asemenea
numerelor aleatoare.
eșantion cu ajutorul
(un număr) fiecărui
e aleatoare sau se ia
valoare, și se includ
d populațiilor finite
n același individ de
e").

puțin două caracte-
rilor eșantionului în

Abaterea unei valori centrale este o soluție a problemei următoare: să se în-
valorile observate (x_1, \dots, x_n) cu o unică valoare c , cât mai apropiată
valorilor ale eșantionului. Aceasta presupune și alegerea unei măsuri
dintre c și x_i , deci a unei măsuri a împrăstierii. Altfel spus, con-
spațiul eșantioanelor de volum n , \mathbb{R}^n , se caută o distanță pe acest
pentru a măsura abaterea dintre (x_1, \dots, x_n) și eșantionul constant

se folosește distanța $\frac{1}{n} \sum_{i=1}^n (x_i - c)^2$, minimul acesteia va fi atins

media $c = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$, iar măsura gradului de împrăstiere devine

dispersia $s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$. Se poate folosi și o altă distanță, pre-

$\sum_{i=1}^n |x_i - c|$, caz în care minimul va fi atins pentru valoarea mediană a

eșantionului, iar măsura gradului de împrăstiere va fi abaterea medie față de

Cuplul (\bar{x}, s^2) este însă cel mai utilizat.

ce urmează, vom presupune pentru simplificare că populația stu-

este infinită, iar selecția este aleatoare și simplă.

7.2.1 Media de selecție

Considerăm un prim eșantion de n observații, (x_1, \dots, x_n) , pentru care s-a

media $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$. Dacă se selectează în condiții similare un al

eșantion de același volum, (x'_1, \dots, x'_n) , media corespunzătoare

$\bar{x}' = \frac{1}{n} \sum_{i=1}^n x'_i$, va fi în general diferită de prima medie observată. Acest lucru

valabil și pentru mediile altor eșantioane prelevate în condiții similare:

(x''_1, \dots, x''_n) etc.

Se considera șirul infinit al observațiilor de rang i din fiecare eșantion

x_i, x'_i, x''_i, \dots ca fiind observații efectuate asupra unei aceleiași variabile

aleatoare $X_i, i = \overline{1, n}$. Valorile medii observate $\bar{x}, \bar{x}', \bar{x}'', \dots$ devin astfel și

ale unei variabile aleatoare \bar{X} , care depinde de X_1, \dots, X_n astfel

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i.$$

Def. \bar{X} se numește *medie de selecție*. Fiind variabile aleatoare, X_1, \dots, X_n și

\bar{X} au repartiții denumite *repartiții de selecție*.

Selecția fiind aleatoare și simplă, variabilele X_1, \dots, X_n sunt independente și identic repartizate. Să notăm cu m media lor și cu σ^2 dispersia.

Proprietate.

$$a) \mathbb{E}\bar{X} = m, \text{Var}(\bar{X}) = \frac{\sigma^2}{n}.$$

b) Din legea numerelor mari, $\bar{X} \xrightarrow[n \rightarrow \infty]{a.s.} m$.

c) Din teorema limitei centrale,

$$\frac{\bar{X} - m}{\frac{\sigma}{\sqrt{n}}} \xrightarrow[n \rightarrow \infty]{repart.} o \text{ v.a. repartizată } N(0, 1).$$

Dem.

$$a) \mathbb{E}\bar{X} = \frac{1}{n} \sum_{i=1}^n \mathbb{E}X_i = m, \text{Var}(\bar{X}) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) = \frac{\sigma^2}{n}.$$

7.2.2 Dispersia de selecție

Procedând ca în cazul mediei de selecție, putem calcula dispersiile pentru diferitele eșantioane:

$$s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2, s'^2 = \frac{1}{n} \sum_{i=1}^n (x'_i - \bar{x}')^2, s''^2 = \frac{1}{n} \sum_{i=1}^n (x''_i - \bar{x}'')^2, \dots,$$

aceste dispersii fiind considerate ca valori observate asupra variabilei aleatoare

$$S^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2.$$

Def. S^2 se numește *dispersie de selecție*.

Proprietate.

$$a) \mathbb{E}S^2 = \frac{n-1}{n} \sigma^2, \text{Var}(S^2) = \frac{n-1}{n^3} [(n-1) \mu_4 - (n-3) \sigma^4]$$

unde μ_4 este momentul centrat de ordin 4 al v.a. X_i .

$$b) S^2 = \overline{X^2} - \bar{X}^2.$$

X_i sunt independente
dispersia.

$$c) S^2 \xrightarrow[n \rightarrow \infty]{a.s.} \sigma^2. (LNM)$$

$$d) \frac{(S^2 - \sigma^2) \sqrt{n}}{\sqrt{\mu_4 - \sigma^4}} \xrightarrow[n \rightarrow \infty]{repart.} o \text{ v.a. repartizată } N(0, 1). (TLC)$$

Dem.

$$S^2 = \frac{1}{n} \sum_{i=1}^n (X_i - m + m - \bar{X})^2 = \frac{1}{n} \sum_{i=1}^n (X_i - m)^2 + \frac{2(m - \bar{X})}{n} \sum_{i=1}^n (X_i -$$

$$+ (m - \bar{X})^2 = \frac{1}{n} \sum_{i=1}^n (X_i - m)^2 - 2(m - \bar{X})^2 + (m - \bar{X})^2 \Rightarrow$$

$$ES^2 = \frac{1}{n} \sum_{i=1}^n Var(X_i) - Var(\bar{X}) = \frac{1}{n} n \sigma^2 - \frac{1}{n} \sigma^2 = \frac{n-1}{n} \sigma^2. \square$$

$$Var(X_i) = \frac{\sigma^2}{n}. \square$$

Obs. Considerând variabila transformată

$$S^{*2} = \frac{n}{n-1} S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2,$$

$$ES^{*2} = \sigma^2.$$

7.2.3 Funcția de repartiție empirică

Notăm cu F funcția de repartiție comună variabilelor X_1, \dots, X_n .

Def. Se numește *funcție de repartiție empirică* (sau *funcție de repartiție empirică*) $F_n^* : \mathbb{R} \times \Omega \rightarrow [0, 1]$ definită prin

$$F_n^*(x, \omega) = \frac{|\{i \in \overline{1, n} \mid X_i(\omega) < x\}|}{n} = \frac{\text{nr. de observații} < x}{n}.$$

Proprietate. Funcția de repartiție empirică are următoarele proprietăți:

1. Pentru $\forall \omega \in \Omega$ fixat, $F_n^*(\cdot, \omega)$ este o funcție în scară.
2. Pentru $\forall x \in \mathbb{R}$ fixat, $F_n^*(x, \cdot)$ este o v.a. simplă cu proprietatea că

$$nF_n^*(x, \cdot) \sim \text{Binomial}(n, F(x)).$$

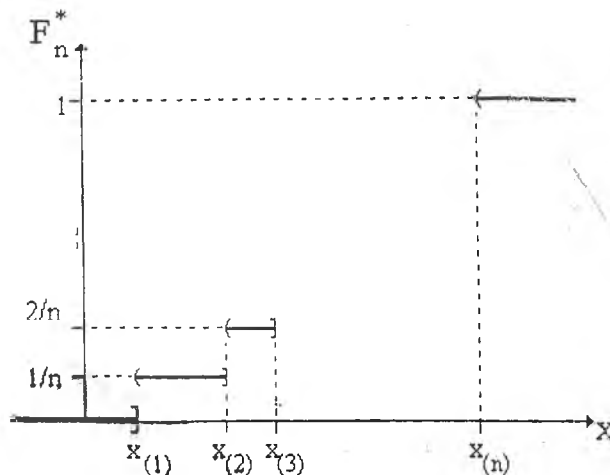
c) $\mathbb{E}(F_n^*(x, \cdot)) = F(x)$, $\text{Var}(F_n^*(x, \cdot)) = \frac{1}{n} F(x) [1 - F(x)]$.

d) Pentru $\forall x \in \mathbb{R}$ fixat, $F_n^*(x, \cdot) \xrightarrow[n \rightarrow \infty]{a.s.} F(x)$.

e) **Teorema Glivenko-Cantelli:** $\sup_{x \in \mathbb{R}} |F_n^*(x, \cdot) - F(x)| \xrightarrow[n \rightarrow \infty]{a.s.} 0 \Leftrightarrow$

$\Leftrightarrow P\left(\limsup_{n \rightarrow \infty} \sup_{x \in \mathbb{R}} |F_n^*(x, \cdot) - F(x)| = 0\right) = 1$ (convergența lui F_n^* către F este uniform a.s.).

Dem. a)



b) Se consideră v.a. auxiliare $Z_k(\omega) = \begin{cases} 1, & X_k(\omega) \leq x \\ 0, & \text{altfel} \end{cases}$, $k = \overline{1, n} \Rightarrow Z_k \sim \text{Bernoulli}(F(x))$. Deci

$$nF_n^*(x, \cdot) = \sum_{k=1}^n Z_k \sim \text{Binomial}(n, F(x)).$$

c) Rezultă imediat din proprietățile repartiției Binomiale.

d) $F_n^*(x, \cdot) = \frac{1}{n} \sum_{k=1}^n Z_k \xrightarrow[n \rightarrow \infty]{a.s.} \mathbb{E} Z_k = F(x)$ (LNM). \square

7.2.4 Eșantre

Valorile unui eșantion crescătoare corespunde și mai numesc și există, cu f densitatea $X_{(i)}$, avem

Proprietăți

Repartiția

Densitatea

Dem

$$F_1(x) =$$

$$= 1 - P$$

Construcții

7.3 Re

Pe baza exs
câteva