

UNIVERSITATEA TEHNICA DE CONSTRUCTII
BUCURESTI

PROBABILITATI SI STATISTICA

Viorel PETREHUS, Sever-Angel POPESCU

BUCURESTI 2005

Cuprins

Cuvânt înainte	v
I Probabilități	1
1 Definiția probabilității	2
1.1 Definiția clasică	2
1.2 Definiția axiomatică a probabilității	3
1.3 Probabilități condiționate	6
1.4 Rezumat	9
1.5 Exerciții	10
2 Variabile aleatoare simple	14
2.1 Definiție și proprietăți	14
2.2 Spațiul de probabilitate produs	19
2.3 Rezumat	20
2.4 Exerciții	21
3 σ Câmpuri de probabilitate	24
3.1 Variabile aleatoare pe σ câmpuri de probabilitate	25
3.2 Media unei variabile aleatoare oarecare	27
3.3 Funcția de repartiție densitatea de probabilitate	30
3.4 Integrala Stieltjes	32
3.5 Media și funcția de repartiție	33
3.6 Rezumat	37
3.7 Exerciții	39

4	Legi clasice	43
4.1	Repartiția binomială	43
4.2	Repartiția Poisson	44
4.3	Repartiția uniformă	46
4.4	Repartiția Normală	46
4.5	Repartiția exponențială negativă	51
4.6	Repartiția Gamma	51
4.7	Repartiția χ^2 (hi pătrat)	53
4.8	Repartiția Student	54
4.9	Rezumat	55
4.10	Exerciții	56
5	Legi limită	60
5.1	Legea numerelor mari	61
5.2	Teoreme limită centrală	63
5.3	Rezumat	68
5.4	Exerciții	69
6	Dependența între variabilele aleatoare	72
6.1	Coeficientul de corelație	72
6.2	Variabile aleatoare bidimensionale	74
6.3	Funcția de repartiție condiționată	79
6.4	Distribuția sumei și câtului	81
6.4.1	Distribuția sumei	81
6.4.2	Distribuția câtului	82
6.5	Distribuția Student	82
6.6	Distribuția Snedecor-Fisher	84
6.7	Exerciții	85
7	Procese aleatoare	89
7.1	Procese Poisson	90
7.2	Procese Markov discrete	92
7.3	Procese de naștere și moarte	95
7.3.1	Model de așteptare cu o singură stație de deservire și un număr mare de unități ce au nevoie de serviciile stației	96
7.3.2	Model de așteptare cu o singură stație iar numărul de unități care au nevoie de serviciile stației este limitat la o valoare dată.	98

7.3.3	Model de aşteptare cu n stații de deservire și cu N unități ce trebuie deservite ($1 < n < N$)	99
7.4	Procese aleatoare staționare	99
7.5	Exerciții	102
II	Statistică	106
8	Statistica descriptivă	107
8.1	Statistica unei variabile	107
8.2	Statistica a două variabile	113
8.3	Exerciții	116
9	Statistici. Estimarea parametrilor	118
9.1	Principiul verosimilității maxime	125
9.2	Metoda momentelor (K. Pearson)	129
9.3	Exerciții	130
10	Intervale de încredere	132
10.1	Intervale de încredere pentru medie	133
10.1.1	Dispersia este cunoscută	134
10.1.2	Dispersia este necunoscută	135
10.2	Intervale de încredere pentru dispersie	136
10.3	Intervale de încredere pentru cîtul a două dispersii	137
10.4	Intervale de încredere în cazul unor selecții mari	138
10.5	Rezumat	139
10.6	Exerciții rezolvate	141
10.7	Exerciții	142
11	Ipoteze statistice. Teste statistice	144
11.1	Ipoteze și testarea lor	144
11.1.1	Testul Z privind media unei populații normale cu dispersia cunoscută σ^2	146
11.1.2	Testul T privind media unei populații normale cu dispersia estimată prin estimatorul nedeplasat S'^2	152
11.1.3	Test pentru proporția de succese	154
11.1.4	Testul T pentru compararea a două esantioane	155
11.2	Tipuri de erori. Reguli de decizie	156
11.3	Puterea unui test statistic	162
11.4	Rezumat	164
11.5	Exerciții rezolvate (mai dificile)	166
11.6	Exerciții propuse	168

12 Testul neparametric χ^2	174
12.1 Principiul testului χ^2	174
12.1.1 Teste asupra formei unei distributii	179
12.1.2 Teste de independentă	179
12.1.3 Teste de omogenitate	180
12.2 Rezumat	184
12.3 Exerciții rezolvate	185
12.4 Exerciții	187
13 Alte teste neparametrice	190
13.1 Testul de concordantă Kolmogorov-Smirnov	190
13.2 Testul lungimilor (secventelor)	192
13.3 Testul lui Wilcoxon I (cazul observațiilor necuplate)	195
13.4 Testul semnelor	196
13.5 Testul lui Wilcoxon II (cazul observațiilor cuplate)	196
13.6 Exerciții	197
14 Analiza dispersiei și analiza regresiei	200
14.1 Analiza dispersiei	200
14.2 Analiza regresiei	202
14.2.1 Metoda celor mai mici pătrate (C. F. Gauss)	204
14.2.2 Condițiile Gauss–Markov pentru metoda celor mai mici pătrate	205
14.2.3 Măsura deviației la metoda celor mai mici pătrate	207
14.2.4 Intervale de încredere și teste pentru β_0 și β_1	210

Cuvânt înainte

Cursul de față a fost scris în perioada 1996-1997 de către Viorel Petrehuș (partea I, probabilități) și Sever-Angel Popescu (partea a II-a, statistică) pentru studenții anului II din Universitatea Tehnică de Construcții București și a apărut în 1997 multiplicat în atelierele universității. El a fost gândit în 14 lecții, câte una pe săptămână, pe parcursul unui semestru. Fiecare lecție se încheie cu exerciții.

Autorii sunt recunoscători tuturor celor care au contribuit cu observațiile lor la buna organizare a materialului prezentat.

Autorii

Partea I

Probabilități

Lecția 1

Definiția probabilității

1.1 Definiția clasică

Probabilitatea a fost privită fie dintr-un punct de vedere "psihologic" ca măsurând gradul de siguranță al observatorului relativ la producerea sau neproducerea unui fenomen, fie "statistic" ca frecvența de apariție a unui fenomen într-un număr mare de experimente independente. Din punct de vedere clasic, definiția care s-a dovedit cea mai eficientă în calcule a fost aceea care a plecat de la conceptul de egală posibilitate. Acest lucru înseamnă că numărul de posibilități într-un experiment este finit și toate posibilitățile au aceeași șansă.

Probabilitatea unui eveniment care constă din mai multe asemenea posibilități este raportul dintre numărul cazurilor favorabile și numărul cazurilor posibile. Utilizarea acestei definiții presupune că într-un fel sau altul putem număra stările posibile și pe cele favorabile.

Exemplul 1.1 *Se aruncă un zar de două ori. Să se determine probabilitățile evenimentelor:*

- a) *Suma celor două zaruri este 6.*
- b) *Ambele zaruri au același număr.*

Soluție. Cazurilor posibile în cele două situații sunt (1,1), (1,2), ... (6,6), în număr de 36. Pentru punctul a) numai cazurile (1,5), (2,4), (3,3), (4,2), (5,1) sunt favorabile. Probabilitatea este deci

$$p = \frac{\text{nr. cazuri favorabile}}{\text{nr. cazuri posibile}} = \frac{5}{36}$$

Analog, probabilitatea celui de-al doilea eveniment este $p=6/36=1/6$.

Exemplul 1.2 *Dintr-un pachet de 36 cărți se extrag trei la întâmplare. Care este probabilitatea ca cel puțin o carte să fie as?*

Soluție. Numărul cazurilor posibile este $C_{36}^3=7140$. Numărul cazurilor favorabile este

$$\overbrace{4 \cdot C_{32}^2}^{\text{un as}} + \overbrace{C_4^2 \cdot C_{32}^1}^{\text{doi ași}} + \overbrace{C_4^3 \cdot 1}^{\text{trei ași}} = 2180$$

Probabilitatea este deci $p = 2180/7140 \approx 0,30532..$

În exemplul următor vedem cât de dificil este uneori să numărăm cazurile posibile sau cazurile favorabile.

Exemplul 1.3 *O persoană scrie n scrisori către n persoane distincte, le pune în n plicuri și apoi scrie adresele la întâmplare. Care este probabilitatea ca cel puțin o scrisoare să ajungă la destinatarul potrivit?*

Numărul cazurilor posibile de a scrie adresele este $n!$. Enumerarea cazurilor favorabile este mai dificilă. După o mică dezvoltare a calculului probabilităților se poate rezolva elegant această problemă (exercițiul 5).

Exemplul 1.4 *Care este probabilitatea ca luând un punct la întâmplare în pătratul $[0,1] \times [0,1]$ el să fie deasupra bisectoarei $y=x$?*

În acest caz numărul cazurilor posibile și favorabile este infinit și definiția clasică nu se mai aplică.

1.2 Definiția axiomatică a probabilității

În general evenimentele se exprimă prin propoziții. Propozițiile obținute prin operațiile logicii matematice (sau,și,non), între propoziții care exprimă evenimente, exprimă la rândul lor alte evenimente. În cele ce urmează probabilitatea este definită pe o mulțime Ω de evenimente care odată cu evenimentele A și B conține și evenimentele următoare exprimate prin operațiile logice sau, și, non:

$$\begin{array}{ll} A \wedge B & A \text{ și } B \\ A \vee B & A \text{ sau } B \\ \overline{A} & \text{non } A \\ 1 & \text{evenimentul sigur} \\ 0 & \text{evenimentul imposibil} \end{array}$$

De asemenea, presupunem că au loc următoarele relații:

a)Comutativitatea

$$A \wedge B = B \wedge A \quad A \vee B = B \vee A$$

b)Asociativitatea

$$\begin{aligned} A \vee (B \vee C) &= (A \vee B) \vee C \\ A \wedge (B \wedge C) &= (A \wedge B) \wedge C \end{aligned}$$

c)Distributivitatea

$$\begin{aligned} A \wedge (B \vee C) &= (A \wedge B) \vee (A \wedge C) \\ A \vee (B \wedge C) &= (A \vee B) \wedge (A \vee C) \end{aligned}$$

d) Absorbția

$$(A \cap B) \cup A = A \quad (A \cup B) \cap A = A$$

e) Legile lui Morgan

$$\overline{A \vee B} = \overline{A} \wedge \overline{B}$$

$$\overline{A \wedge B} = \overline{A} \vee \overline{B}$$

f) Evenimentele 1 și 0 se caracterizează prin

$$0 \wedge A = 0, \quad 0 \vee A = A$$

$$1 \wedge A = A, \quad 1 \vee A = 1$$

g) Evenimentul non A sau contrarul lui A are proprietățile:

$$A \wedge \overline{A} = 0 \quad A \vee \overline{A} = 1$$

Remarcăm că aceste relații sunt adevărate dacă A, B, \dots sunt propoziții, iar 0 este propoziția totdeauna falsă și 1 este propoziția totdeauna adevărată.

O mulțime de evenimente cu proprietățile de mai sus se numește câmp de evenimente sau algebră de evenimente sau algebră Boole. Relațiile de mai sus nu sunt independente, deci nu formează un set minimal de axiome pentru algebrele Boole. Pe de altă parte ele implică alte relații importante, ca de exemplu $\overline{\overline{A}} = A$. Următoarea teoremă a lui Stone descrie asemenea câmpuri de evenimente prin submulțimi.

Teorema 1.5 (Stone) Fie M o mulțime de evenimente cu proprietățile de mai sus. Atunci există o mulțime X și o submulțime $\Omega \subset P(X)$ cu proprietățile :

$$1) \emptyset \in \Omega, X \in \Omega$$

$$2) A, B \in \Omega \Rightarrow A \cup B \in \Omega$$

$$3) A, B \in \Omega \Rightarrow A - B \in \Omega$$

și o bijectie $\sigma : M \rightarrow \Omega$ astfel ca:

$$a) \sigma(A \vee B) = \sigma(A) \cup \sigma(B)$$

$$b) \sigma(A \wedge B) = \sigma(A) \cap \sigma(B)$$

$$c) \sigma(\overline{A}) = \overline{\sigma(A)}$$

$$\sigma(0) = \emptyset \quad \sigma(1) = X$$

Teorema arată că în esență orice câmp de evenimente poate fi reprezentat prin submulțimi ale aceleiași mulțimi X . Operației 'și' între evenimente îi corespunde intersecția submulțimilor, operației 'sau' îi corespunde reuniunea, negației îi corespunde complementara, evenimentul sigur corespunde mulțimii totale X , iar evenimentul imposibil corespunde cu mulțimea vidă.

Definiția 1.6 O mulțime $\Omega \subset P(X)$ cu proprietățile 1,2,3 din teorema lui Stone se numește clan, algebră de evenimente sau câmp de evenimente. Uneori o vom numi algebră de mulțimi.

Propoziția 1.7 Fie $\Omega \subset P(X)$ un câmp de evenimente.

- a) Dacă $A_1, A_2, \dots, A_n \in \Omega$ atunci $\bigcup_{i=1..n} A_i \in \Omega$ și $\bigcap_{i=1..n} A_i \in \Omega$
 b) Dacă $A \in \Omega$ atunci $\overline{A} \in \Omega$.

Demonstrație. b) $\overline{A} = X - A$ și din proprietățile 1 și 3 ale algebrei de mulțimi rezultă $\overline{A} \in \Omega$. a) Dacă A și B sunt în Ω atunci $A \cap B = \overline{\overline{A} \cup \overline{B}} \in \Omega$ conform cu proprietatea 2 a algebrei de mulțimi și punctul b). Prin inducție rezultă acum ușor și a).

Vedem că în general un număr finit de operații de intersecție, reuniune, diferență sau complementară cu mulțimi din Ω are ca rezultat tot o mulțime din Ω . Vom defini acum probabilitatea.

Definiția 1.8 Fie $\Omega \subset P(X)$ un câmp de evenimente. Se numește probabilitate pe Ω , o funcție $p: \Omega \rightarrow R$ cu proprietățile:

- 1) $p(A) \in [0, 1]$ pentru oricare $A \in \Omega$
- 2) $p(A \cup B) = p(A) + p(B)$ dacă $A \cap B = \emptyset$
- 3) $p(\emptyset) = 0$ $p(X) = 1$

Tripletul (X, Ω, p) îl vom numi câmp de probabilitate.

Teorema 1.9 Fie (X, Ω, p) un câmp de probabilitate. Atunci:

- 1) $p(\overline{A}) = 1 - p(A)$
- 2) Dacă A_1, A_2, \dots, A_n sunt în Ω și $A_i \cap A_j = \emptyset$ pentru orice i și j atunci

$$p\left(\bigcup_{i=1..n} A_i\right) = \sum_{i=1}^n p(A_i) \quad (1.1)$$

- 3) $p(A \cup B) = p(A) + p(B) - p(A \cap B)$ A și B nefiind neapărat disjuncte.
- 4) $A \subset B$ implică $p(A) \leq p(B)$.

Demonstrație. 1) $X = \overline{A} \cup A$, $A \cap \overline{A} = \emptyset$, deci $1 = p(X) = p(A) + p(\overline{A})$

2) Pentru $n=2$ afirmația rezultă din definiția probabilității, pe urmă se procedează prin inducție.

3) Să privim evenimentele ca mulțimi. Atunci $A \cup B = (A \cap \overline{B}) \cup (B \cap \overline{A}) \cup (A \cap B)$, reuniune disjunctă, deci $p(A \cup B) = p(A \cap \overline{B}) + p(B \cap \overline{A}) + p(A \cap B)$. De asemenea $A = (A \cap \overline{B}) \cup (A \cap B)$ deci $p(A) = p(A \cap \overline{B}) + p(A \cap B)$ și prin urmare $p(A \cap \overline{B}) = p(A) - p(A \cap B)$. Analog $p(B \cap \overline{A}) = p(B) - p(A \cap B)$. Punând aceste valori în expresia precedentă pentru $p(A \cup B)$ se obține punctul 3).

4) $A \subset B \Rightarrow B = (B \cap \overline{A}) \cup A$ (disjuncte) $\Rightarrow p(B) = p(B \cap \overline{A}) + p(A) \geq p(A)$.

Observația 1.10 Datorită formulei (1.1) probabilitatea p se mai numește finit aditivă.

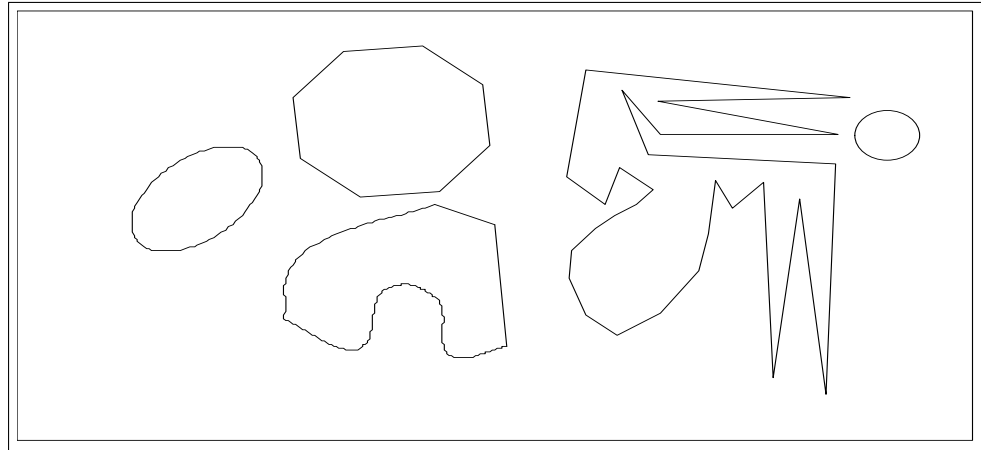
Exemplul 1.11 Se aruncă două zaruri. Pot apare toate perechile de fețe (i,j) cu $1 \leq i \leq 6$ și $1 \leq j \leq 6$, în număr de 36. Fie X mulțimea acestor perechi. Fie $\Omega = P(X)$ și $p : \Omega \rightarrow R$ definită prin $p(A) = \frac{\text{numar de elemente din } A}{\text{numar de elemente din } X}$. Funcția p definită în acest fel este o probabilitate în sensul definiției anterioare. Numărul de elemente din X reprezintă numărul cazurilor posibile, iar numărul de elemente din A reprezintă numărul cazurilor favorabile. În acest fel definiția clasică a probabilităților este cuprinsă în definiția axiomatică de mai sus.

Acum putem să dăm o soluție pentru problema pusă în exemplul 4.

Exemplul 1.12 Fie $X=[0,1] \times [0,1]$. Fie $\Omega \subset P(X)$ formată din submulțimile lui X cu frontiera formată dintr-un număr finit de curbe de clasă C^1 pe porțiuni (adică fiecare curbă poate avea cel mult un număr finit de puncte unde să nu fie cu derivata continuă). Se verifică ușor că Ω este o algebră de mulțimi. Definim acum probabilitatea prin

$$p(A) = \frac{\text{aria}(A)}{\text{aria}(X)} \quad \text{pentru } A \in \Omega$$

Este clar că p astfel definită satisface condițiile din definiția axiomatică a probabilității. În raport cu această probabilitate evenimentul din exemplul 4 este reprezentat prin $A = \{(x, y) \in X \mid y > x\}$ pentru care $p(A) = 1/2$. Asemenea probabilități în care X este o submulțime în R, R^2, R^3 , unde Ω este formată din submulțimi "cu arie" ale lui X , iar $p(A) = \frac{\int_A f(\xi) d\xi}{\int_X f(\xi) d\xi}$ se numesc probabilități geometrice. Funcția f în această definiție se mai numește pondere sau densitate și trebuie să fie pozitivă.



Diverse tipuri de figuri pentru care este definită probabilitatea geometrică

1.3 Probabilități condiționate

Fie (X, Ω, p) un câmp de probabilitate, B în Ω , $p(B) > 0$.

Definiția 1.13 Definim $p_B : \Omega \rightarrow R$ (sau $p(\cdot | B)$) prin

$$p(A|B) = p_B(A) = \frac{p(A \cap B)}{p(B)} \quad (1.2)$$

și o numim probabilitatea lui A condiționată de B.

Teorema 1.14 În condițiile definiției de mai sus p_B este o probabilitate și în plus $p(A \cap B) = p(B) \cdot p_B(A)$.

Demonstrație. Deoarece $\emptyset \subset A \cap B \subset B$ atunci $0 \leq p_B(A) \leq 1$. De asemenea $p_B(\emptyset) = 0$ și $p_B(X) = 1$ sunt evidente. Dacă $A_1 \cap A_2 = \emptyset$ atunci și $(A_1 \cap B) \cap (A_2 \cap B) = \emptyset$, deci

$$\begin{aligned} p_B(A_1 \cup A_2) &= \frac{p((A_1 \cup A_2) \cap B)}{p(B)} = \frac{p((A_1 \cap B) \cup (A_2 \cap B))}{p(B)} \\ &= \frac{p(A_1 \cap B) + p(A_2 \cap B)}{p(B)} = p_B(A_1) + p_B(A_2) \end{aligned}$$

Definiția 1.15 Spunem că două evenimente A și B sunt independente dacă $p(A \cap B) = p(A) \cdot p(B)$, altfel spus dacă $p(A) = p_B(A)$.

Propoziția 1.16 Dacă A și B sunt independente atunci și perechile de evenimente (A, \bar{B}) , (\bar{A}, B) , (\bar{A}, \bar{B}) sunt independente.

Demonstrație. Prin definiție $p(A \cap B) = p(A)p(B)$. Deoarece $A = (A \cap B) \cup (A \cap \bar{B})$ (reuniune disjunctă), avem

$$p(A) = p(A \cap B) + p(A \cap \bar{B}) = p(A)p(B) + p(A \cap \bar{B})$$

de unde rezultă

$$p(A \cap \bar{B}) = p(A) - p(A)p(B) = p(A)(1 - p(B)) = p(A)p(\bar{B})$$

. Analg se demonstrează independența și în celelalte cazuri.

În mod analog se definește independența mai multor evenimente:

Definiția 1.17 A_1, A_2, \dots, A_n sunt independente dacă pentru orice $k \leq n$ și orice evenimente $A_{i_1}, A_{i_2}, \dots, A_{i_k}$ dintre cele A_1, \dots, A_n date avem

$$p(A_{i_1} \cap A_{i_2} \cap \dots \cap A_{i_k}) = p(A_{i_1})p(A_{i_2}) \dots p(A_{i_k}). \quad (1.3)$$

Analog cu propoziția anterioară se poate demonstra și

Propoziția 1.18 Dacă evenimentele A_1, A_2, \dots, A_n sunt independente atunci și evenimentele $A_1, \bar{A}_2, \dots, A_n$ sunt independente (orice combinație de A_i sau complementare).

Demonstrația este lăsată ca exercițiu.

Definiția 1.19 Spunem că evenimentele A_1, A_2, \dots, A_n formează o partiție a lui X dacă

- 1) $\bigcup_{i=1..n} A_i = X$
- 2) $A_i \cap A_j = \emptyset$ pentru orice $i \neq j$.

Avem următoarea

Propoziția 1.20 Fie (X, Ω, p) un câmp de probabilitate și partiția A_1, \dots, A_n . Atunci pentru orice $B \in \Omega$ avem

$$p(B) = \sum_{i=1}^n p(A_i) \cdot p(B|A_i) \quad (1.4)$$

Formula de mai sus se numește formula probabilității totale.

Demonstrație.

$$p(B) = p(B \cap X) = p(B \cap (\bigcup_{i=1..n} A_i)) = p(\bigcup_{i=1..n} (B \cap A_i)) = \sum_{i=1}^n p(B \cap A_i) = \sum_{i=1}^n p(A_i)p(B) \quad (1.5)$$

Q.E.D.

Formula următoare, numită formula lui Bayes, ne dă probabilitatea ca după ce știm că un eveniment ce poate apare din mai multe cauze s-a realizat, acesta să se fi realizat dintr-o cauză anume. Enunțul precis este:

Propoziția 1.21 (Bayes) Fie (X, Ω, p) un câmp de probabilitate și A_1, A_2, \dots, A_n o partiție a lui X . Fie B un eveniment cu probabilitate nenulă. Atunci:

$$p(A_i|B) = \frac{p(A_i)p(B|A_i)}{\sum_{j=1}^n p(A_j)p(B|A_j)} \quad (1.6)$$

Demonstrație.

$$p(A_i|B) = \frac{p(B \cap A_i)}{p(B)} = \frac{p(A_i)p(B|A_i)}{\sum_{j=1}^n p(A_j)p(B|A_j)}$$

conform formulei probabilității totale.

Exemplul 1.22 În trei urne cu bile albe și negre avem compozițiile: 4+6, 2+8, 4+1. Se pun toate bilele la un loc și se extrage o bilă la întâmplare. Se constată că este albă. Care este probabilitatea ca ea să fi provenit din urna a treia?

Soluție. Fie A_i evenimentul care constă în extragerea unei bile provenind din urna i și B evenimentul care constă în extragerea unei bile albe. Avem $p(A_1) = \frac{10}{25}$ $p(A_2) = \frac{10}{25}$ $p(A_3) = \frac{5}{25}$ (probabilitatea evenimentului A_i este proporțională cu numărul bilelor provenite din urna i). Mai departe avem $p(B|A_1) = \frac{4}{10}$ $p(B|A_2) = \frac{2}{10}$ $p(B|A_3) = \frac{4}{5}$. Probabilitatea care ni se cere este $p(A_3|B)$ și conform cu formula lui Bayes

$$p(A_3|B) = \frac{\frac{10}{25} \cdot \frac{4}{5}}{\frac{10}{25} \cdot \frac{4}{10} + \frac{10}{25} \cdot \frac{2}{10} + \frac{5}{25} \cdot \frac{4}{5}} = \frac{2}{5}$$

Observația 1.23 Dacă alegerea unei bile s-ar fi făcut după regulile: i) se alege la întâmplare o urnă; 2) din urna aleasă se extrage la întâmplare o bilă; atunci am avea $p(A_1) = p(A_2) = p(A_3) = \frac{1}{3}$ și tot formula lui Bayes conduce la $p(A_3|B) = \frac{4}{7}$.

1.4 Rezumat

Probabilitatea are o componentă psihologică în sensul de grad de încredere și una practică în sensul de frecvență de apariție a unui fenomen într-un număr mare de experiențe independente, dar pentru a putea introduce numere și a face calcule s-a dovedit utilă o definiție axiomatică, asemănător cum în geometrie dreapta se introduce axiomatic, independent de mulțimea exemplelor practice. Punctele importante în acest demers sunt:

- i) Reprezentarea evenimentelor prin algebre de mulțimi (Stone).
- ii) Definiția probabilității ca funcție pe algebre de evenimente, cu proprietăți asemănătoare cu aria figurilor plane.

În dezvoltarea elementară a calculului cu probabilități următoarele trei puncte sunt esențiale:

- iii) Noțiunea de probabilitate condiționată și probabilitatea de apariție simultană a evenimentelor independente.
- iv) Formula probabilității totale (1.4).
- v) Formula pentru probabilitatea cauzelor (1.6).

FORMULE UTILIZATE FRECVENT :

a) $p = \frac{\text{nr. cazuri favorabile}}{\text{nr. cazuri posibile}}$.

b) Definiția probabilității condiționate $p(A|B) = \frac{p(A \cap B)}{p(B)}$ (1.2);

probabilitatea evenimentelor independente $p(A \cap B) = p(A) \cdot p(B)$ (1.3)

c) Formula probabilității totale $p(B) = \sum_{i=1}^n p(A_i) \cdot p(B|A_i)$ (1.4).

d) Formula lui Bayes $p(A_i|B) = \frac{p(A_i)p(B|A_i)}{\sum_{j=1}^n p(A_j)p(B|A_j)}$ (1.6).

1.5 Exerciții

1. O urnă conține jetoane numerotate de la 1 la 8. Se extrag la întâmplare 3 jetoane. Care este probabilitatea ca suma numerelor extrase să fie superioară sau egală cu suma numerelor rămase?

Indicație. Probabilitatea este raportul dintre numărul cazurilor favorabile și numărul cazurilor posibile. Sunt posibile $C_8^3 = 56$ cazuri. Suma tuturor numerelor este 36, deci suma celor extrase trebuie să depășească 18 pentru a fi caz favorabil. Acum se enumeră pur și simplu toate cazurile favorabile: (8,7,6), (8,7,5), etc.

2. Se aleg la întâmplare a, b, c în intervalul $(0, 1)$. Care este probabilitatea ca ecuația de gradul doi $ax^2 + 2bx + c = 0$ să aibă rădăcini reale? (se consideră $(a, b, c) \in (0, 1) \times (0, 1) \times (0, 1)$ iar probabilitatea este proporțională cu volumul).

Indicație. $(a, b, c) \in (0, 1)^3$. Condiția de rădăcini reale este $4b^2 - 4ac \geq 0$ sau $1 \geq b \geq \sqrt{ac}$. Fie V = volumul lui $(0, 1)^3$ și V_f = volumul cazurilor favorabile. Avem $V_f = \int \int \int_{V_f} da \cdot dc \cdot db = \int_0^1 da \int_0^1 dc \int_{\sqrt{ac}}^1 db$. Ținând seama că $V=1$, se obține imediat probabilitatea $p = V_f/V$.

3. La un teatru sunt $2n$ persoane la coadă la bilete. n persoane au bancnote de 500 lei iar celelalte n persoane au bancnote de 1000 lei. Un bilet costa 500 lei. Care este probabilitatea ca așezându-se întâmplător la coadă să poată toate persoanele să-și cumpere bilet, știind că inițial în casă nu este nici un leu și fiecare persoană primește restul de la casă?

Soluție Ca să numărăm așezările în care toți își pot cumpăra bilete, reprezentăm într-un sistem xOy pe cele $2n$ persoane prin puncte în $1, 2, 3, \dots, 2n$ pe Ox . Fiecărei așezări îi asociem o funcție definită pe $\overline{0, 2n}$ prin $f(0)=0$ și $f(i)=f(i-1)+1$ dacă persoana i de la rînd are 1000 lei și $f(i)=f(i-1)-1$ dacă persoana i are 500 lei. Fiecare traiectorie începe în $(0, 0)$ și se termină în $(2n, 0)$, deoarece de cîte ori se urcă se și coboară. Numărul aranjărilor la coadă este $n!n!C_{2n}^n$ unde C_{2n}^n reprezintă numărul de alegeri ale locurilor de către cei cu 500 lei (sau numărul traiectoriilor) iar $n!$ reprezintă numărul de permutări celor cu 500 lei sau cu 1000 lei între ei. O traiectorie reprezintă o așezare nefavorabilă dacă pentru un i avem $f(i)=1$. În acest caz primul i cu proprietatea de mai sus reprezintă cel cu 1000 lei care nu mai poate primi rest de la casă. Considerînd simetrica acestei traiectorii de la acest i față de $y=1$ obținem o traiectorie ce se termină în $(2n, 2)$ și are $n+1$ urcușuri și $n-1$ coborâșuri. Numărul traiectoriilor de acest fel este C_{2n}^{n+1} , deci numărul cazurilor nefavorabile este $n!n!C_{2n}^{n+1}$. Probabilitatea este deci

$$p = \frac{n!n!C_{2n}^n - n!n!C_{2n}^{n+1}}{n!n!C_{2n}^n} = \frac{1}{n+1}$$

4. Fie (X, Ω, p) un câmp de probabilitate și A_1, A_2, \dots, A_n evenimente din Ω . Să se arate

că

$$\begin{aligned}
 p\left(\bigcup_{i=1..n} A_i\right) &= \sum_{i=1}^n p(A_i) - \sum_{i \neq j} p(A_i \cap A_j) \\
 &\quad + \sum_{i \neq j \neq k} p(A_i \cap A_j \cap A_k) + \\
 &\quad + \dots + (-1)^{n-1} p(A_1 \cap A_2 \cap \dots \cap A_n)
 \end{aligned} \tag{1.7}$$

(formula lui Poincaré).

Indicație. Formula a fost demonstrată pentru $n=2$ în cadrul lecției. Presupunem prin inducție că este adevărată pentru n . Fie $B = \bigcup_{i=1..n} A_i$ și $B_i = A_{n+1} \cap A_i$. Avem

$$\begin{aligned}
 p\left(\bigcup_{i=1..n+1} A_i\right) &= p(B \cup A_{n+1}) \\
 &= p(B) + p(A_{n+1}) - p(B \cap A_{n+1}) \\
 &= (\text{formula 1.7}) + p(A_{n+1}) - p\left(\bigcup_{i=1..n} B_i\right) \\
 &= (\text{formula 1.7}) + p(A_{n+1}) - (\text{formula 1.7 pt. } B_i) \\
 &= (\text{formula 1.7 pt. } n+1)
 \end{aligned}$$

5. O persoană scrie n scrisori către n persoane distincte, le pune în n plicuri și apoi scrie adresele la întâmplare. Care este probabilitatea ca cel puțin o scrisoare să ajungă la destinatarul potrivit? (problema concordanțelor).

Indicație. Fie A_k evenimentul ce constă în faptul că persoana k primește scrisoarea potrivită. $p(A_k) = \frac{(n-1)!}{n!} = 1/n$. Există C_n^2 evenimente de tipul $A_i \cap A_j$ și fiecare are probabilitatea $\frac{(n-2)!}{n!}$ (deoarece două persoane primesc scrisorile potrivite iar celelalte $n-2$ persoane primesc scrisorile întâmplător în $(n-2)!$ feluri. Analog, există C_n^3 evenimente de tipul $A_i \cap A_j \cap A_k$ și fiecare are probabilitatea $\frac{(n-3)!}{n!}$, etc. Aplicând formula de mai sus rezultă:

$$\begin{aligned}
 p &= n \cdot \frac{1}{n} - C_n^2 \cdot \frac{(n-2)!}{n!} + C_n^3 \cdot \frac{(n-3)!}{n!} + \dots \\
 &= \frac{1}{1!} - \frac{1}{2!} + \frac{1}{3!} + \dots
 \end{aligned}$$

6. Două persoane A și B s-au înțeles să se întâlnească într-un anumit loc între orele 12 și 13. Persoana care vine prima așteaptă 20 de minute și apoi pleacă. Dacă fiecare persoană vine la întâmplare în acel loc, care este probabilitatea ca persoanele să se întâlnească?

Indicație. Prima problemă este cum reprezentăm toate posibilitățile de întâlnire. Fie pe axa x ora de sosire a primei persoane iar pe axa y ora de sosire a celei de a doua persoane.

Toate variantele de sosire se reprezintă deci prin produsul $[12, 13] \times [12, 13]$. Condiția de întâlnire este $|x - y| \leq \frac{1}{3}$ ore. Admițând că probabilitatea este proporțională cu aria, găsim că aria zonei de întâlnire este $5/9$ iar aria zonei de sosire este 1. Deci probabilitatea este $5/9$.

7. Un magazin se aprovizionează de la trei fabrici A, B, C , cu un anumit produs, în proporție de 30%, 60%, 10%. Proporția de articole defecte dintre cele achiziționate este de 2%, 1%, 5%, pentru A , respectiv B și C . Dacă un cumpărător găsește că produsul achiziționat la magazin este cu defecțiuni, care este probabilitatea ca el să fi provenit de la fabrica B ?

Indicație. Se aplică formula lui Bayes.

8. Într-un circuit sunt legate în serie rezistențele R_1, R_2 și grupul de rezistențe în paralel: R_3, R_4, R_5 . Probabilitățile de defecțiune ale celor cinci elemente independente sunt: 0,1; 0,01; 0,05; 0,04; 0,1. Care este probabilitatea de întrerupere a circuitului?

Indicație. Fie A_k evenimentul care constă în întreruperea rezistenței R_k . Avem $p(A_1) = 0,1$, $p(A_2) = 0,01$ etc. Toate aceste evenimente sunt independente și întreruperea circuitului este: $A_1 \cup A_2 \cup (A_3 \cap A_4 \cap A_5)$. Se aplică formula lui Poincaré și se ține seama că evenimentele A_k sunt independente.

9. Trei țintași nimeresc o țintă cu probabilitățile 0,7; 0,8; 0,9. Fiecare trage câte o lovitură. Care este probabilitatea ca:

- Toți trei să nimerescă ținta?
- Cel puțin unul să nimerescă ținta?
- Unul singur să nimerescă ținta?

Indicație. Fie A_k evenimentul care constă în faptul că trăgătorul k nimereste. Avem $p(A_1) = 0,7$; $p(A_2) = 0,8$; $p(A_3) = 0,9$ iar evenimentele A_1, A_2, A_3 sunt independente. La punctul a) se cere probabilitatea lui $A_1 \cap A_2 \cap A_3$ care este produsul probabilităților $p(A_1) \cdot p(A_2) \cdot p(A_3)$. Analog se consideră și celelalte cazuri.

10. Fie $X = \{1, 2, 3, 4\}$ și $\Omega = P(X)$. Se definește o probabilitate pe X astfel ca $p(\{1, 2, 3\}) = 7/8$, $p(\{2, 3, 4\}) = 1/2$, și $p(\{1, 4\}) = 5/8$. Să se determine $p(\{1\})$ și $p(\{1, 3, 4\})$.

11. Într-o urnă se află bile numerotate 0, 1, 2, ..., 9. Se extrag 3 bile la întâmplare, fără a pune bila înapoi. Se scrie numărul format din cele trei cifre în ordinea apariției. Care este probabilitatea ca numărul să se dividă la 12? Dar dacă extragerea se face cu punerea bilei înapoi?

12. Un grup format din $2n$ băieți și $2n$ fete este despărțit în două subgrupuri de același efectiv. Care este probabilitatea ca în fiecare subgrup numărul de băieți să fie egal cu numărul de fete.

13. Fie A_1, A_2, \dots, A_n evenimente independente. Care este numărul maxim de evenimente

distincte care se poate obține din A_1, \dots, A_n prin aplicarea repetată a operațiilor și, sau, non?

14. Pe un segment AB se iau la întâmplare două puncte C și D , astfel ca $|AC| < |AD|$. Care este probabilitatea ca să se poată forma un triunghi cu segmentele AC , CD , DB ?

15. O urnă conține 5 bile albe, 3 bile negre și 2 bile roșii iar altă urnă conține 3 bile negre, 2 bile albe și 5 bile roșii. Se extrage câte o bilă din fiecare urnă. Care este probabilitatea ca să se extragă bile de aceeași culoare?

16. Fie evenimentele A_1, A_2, \dots, A_n . Să se demonstreze că

$$p(A_1 \cap A_2 \cap \dots \cap A_n) \geq p(A_1) + p(A_2) + \dots + p(A_n) - (n - 1)$$

(inegalitatea lui Boole).

17. Un submarin lansează n torpile asupra unui vas. Dacă fiecare torpilă are probabilitatea $\frac{1}{2}$ de a lovi vasul, independent de celelalte torpile, care este numărul minim de torpile ce trebuie lansate ca probabilitatea de a fi lovit vasul de cel puțin una să depășească 0,9?

Lecția 2

Variabile aleatoare simple

2.1 Definiție și proprietăți

De fiecare dată când aplicăm teoria probabilităților subînțelegem că există o algebră de evenimente realizată ca o mulțime de submulțimi ale unei mulțimi X , și, o probabilitate definită pe acele evenimente, chiar dacă X și algebra de evenimente nu sunt exprimate explicit. În general informația din X se extrage prin funcții. În cele ce urmează, X este o mulțime pe care avem definită o algebră de evenimente Ω și o probabilitate p .

Definiția 2.1 O funcție $f : X \rightarrow R$ se numește variabilă aleatoare (pe scurt v.a.) dacă pentru orice $(a, b) \subset R$ avem $f^{-1}(a, b) \in \Omega$. O funcție $f : X \rightarrow C$ se numește variabilă aleatoare dacă $Re(f)$ și $Im(f)$ sunt variabile aleatoare reale.

În cele ce urmează variabilele aleatoare considerate iau un număr finit de valori. Asemenea variabile aleatoare le vom numi simple. Fie v_1, v_2, \dots, v_n valorile distincte ale lui f și fie $A_i = f^{-1}(v_i)$. În aceste condiții variabila aleatoare se scrie

$$f(x) = \begin{cases} v_1 \text{ dacă } x \in A_1 \\ v_2 \text{ dacă } x \in A_2 \\ \dots \\ v_n \text{ dacă } x \in A_n \end{cases} \quad (2.1)$$

Este clar că definiția de mai sus este echivalentă cu a spune că pentru orice valoare v rezultă că $f^{-1}(v) \in \Omega$. Vom mai scrie $A_i = \{f = v_i\}$, $p_i = p(A_i) = p(f^{-1}(v_i)) = p(f = v_i)$. Este mai puțin important din punctul de vedere al calculului cu probabilități, care este mulțimea A_i , cât care sunt probabilitățile $p(A_i)$, $p(A_i \cap A_j)$ etc. Fiecărei variabile aleatoare f îi vom asocia o diagramă (notată tot f):

$$f = \begin{pmatrix} v_1 & v_2 & v_3 & \dots & v_n \\ p_1 & p_2 & p_3 & \dots & p_n \end{pmatrix}$$

Valorile v_1, v_2, \dots, v_n sunt în general distincte iar evenimentele $A_i, i = 1, 2, \dots, n$ formează o partiție a lui X , adică A_i și A_j sunt disjuncte dacă $i \neq j$ și $\bigcup_{i=1..n} A_i = X$. Este clar că $p_1 + p_2 + \dots + p_n = 1$. Putem considera și diagrame în care unele valori v coincid între ele ceea ce ar corespunde definirii lui f prin (2.1) și unde am avea de exemplu $v_1 = v_2$ dar neapărat $A_1 \cap A_2 = \emptyset$. În acest caz A_i nu mai este neapărat $f^{-1}(v_i)$ dar A_1, \dots, A_n formează încă o partiție, $p_i = p(A_i)$ și $p_1 + \dots + p_n = 1$. De exemplu funcția f^k ia valorile $(v_i)^k$ pe mulțimile A_i , deci diagrama asociată este:

$$f^k = \begin{pmatrix} v_1^k & v_2^k & \dots & v_n^k \\ p_1 & p_2 & \dots & p_n \end{pmatrix}$$

iar dacă $v_1 = -v_2$, $k=2$ atunci $v_1^2 = v_2^2$ și avem o diagramă unde unele valori din rândul de sus coincid. Diagramele în care $v_i \neq v_j$ pentru $i \neq j$ le vom numi standard. Dată o variabilă aleatoare simplă, ca mai sus, definim:

Definiția 2.2 1)media lui f : $M(f) = \sum_{i=1}^n v_i p_i = \sum_{i=1}^n v_i \cdot p(f = v_i)$
 2)momentul de ordin k : $M_k(f)$ (sau m_k) = $\sum_{i=1}^n v_i^k p_i = M(f^k)$
 3)momentul centrat de ordin k : $\mu_k(f) = \sum_{i=1}^n (v_i - M(f))^k p_i = M((f - M(f))^k)$
 4)dispersia: $D(f) = \sigma^2(f) = \sum_{i=1}^n (v_i - M(f))^2 p_i = M((f - M(f))^2) = \mu_2(f)$
 5)funcția caracteristică $f_c : R \rightarrow C$ sau $\varphi_f : R \rightarrow C$, $f_c(t) = \varphi_f(t) = \sum_{i=1}^n e^{\sqrt{-1}v_i t} p_i$
 6)funcția de repartiție $F : R \rightarrow R$, $F(t) = p(\{x \in X | f(x) < t\})$.

Observația 2.3 Dacă pentru un i avem $p_i = 0$ atunci valoarea corespunzătoare v_i o putem exclude din diagramă deoarece nu are nici o contribuție la caracteristicile numerice ale lui f . Despre două v.a. f și g , pentru care mulțimea pentru care $f \neq g$ are probabilitatea 0, spunem că ele coincid aproape peste tot (pe scurt a.p.t.). Ele au aceleași diagrame standard și aceleași caracteristici numerice: medie, dispersie,.... Dacă unei v.a. f i se asociază două diagrame distincte, atunci definițiile de mai sus pentru medie, dispersie, momente, funcție caracteristică, dau același rezultat, indiferent de diagrama folosită. Se poate demonstra acest lucru ușor comparând valorile obținute dintr-o diagramă cu cele date de diagrama standard (exercițiu).

Vom nota în general cu $\{a < f < b\}$ evenimentul $f^{-1}(a, b) = \{x \in X | a < f(x) < b\} \in \Omega$ iar probabilitatea lui cu $p(a < f < b)$. Analog vom nota evenimentele ce se obțin folosind inegalități de tipul \leq, \geq , etc.

Exemplul 2.4 Într-o clasă s-au obținut următoarele note: 5 de către 8 elevi, 7 de către 3 elevi, 8 de către 5 elevi, 10 de către 4 elevi. Funcția nota ia valorile 5,7,8,10, iar probabilitățile corespunzătoare sunt: $\frac{8}{20}$ $\frac{3}{20}$ $\frac{5}{20}$ $\frac{4}{20}$. Diagrama asociată funcției este:

$$\text{nota} \left(\begin{array}{cccc} 5 & 7 & 8 & 10 \\ \frac{8}{20} & \frac{3}{20} & \frac{5}{20} & \frac{4}{20} \end{array} \right)$$

media notelor este $M_1 = 5 \cdot \frac{8}{20} + 7 \cdot \frac{3}{20} + 8 \cdot \frac{5}{20} + 10 \cdot \frac{4}{20} = 7,05$

momentul de ordinul doi este $M_2 = 5^2 \cdot \frac{8}{20} + 7^2 \cdot \frac{3}{20} + 8^2 \cdot \frac{5}{20} + 10^2 \cdot \frac{4}{20} = 53,35$

dispersia este $D = (5 - 7,05)^2 \cdot \frac{8}{20} + (7 - 7,05)^2 \cdot \frac{3}{20} + (8 - 7,05)^2 \cdot \frac{5}{20} + (10 - 7,05)^2 \cdot \frac{4}{20} = 3,6475$

funcția caracteristică este $f_c(t) = e^{i5t} \cdot \frac{8}{20} + e^{i7t} \cdot \frac{3}{20} + e^{i8t} \cdot \frac{5}{20} + e^{i10t} \cdot \frac{4}{20}$ unde $i = \sqrt{-1}$

Observația 2.5 Media se poate numi centrul valorilor luate de o v.a. iar dispersia este o măsură a împrăstierii valorilor acelei v.a. în jurul mediei. In afară de medie, celelalte caracteristici se utilizează doar pentru v.a. reale.

Definiția 2.6 Două variabile aleatoare reale, f și g definite pe același câmp de probabilitate X se numesc independente dacă oricare ar fi intervalele (a,b) și (c,d) avem

$$p(f^{-1}(a,b) \cap g^{-1}(c,d)) = p(f^{-1}(a,b)) \cdot p(g^{-1}(c,d))$$

Definiția independenței a două v.a. simple se poate da în mai multe feluri așa cum rezultă din propoziția următoare:

Propoziția 2.7 Fie (X, Ω, p) un câmp finit de probabilitate și f, g două variabile aleatoare simple cu valori reale. Următoarele afirmații sunt echivalente:

- 1) f și g sunt independente
- 2) evenimentele $\{a < f < b\}$ și $\{c < g < d\}$ sunt independente pentru $a, b, c, d \in R$
- 3) pentru orice $v, u \in R$ evenimentele $\{f = v\}$ și $\{g = u\}$ sunt independente
- 4) pentru orice submulțimi A, B din R , evenimentele $\{f \in A\}$ și $\{g \in B\}$ sunt independente

Demonstrația este lăsată ca exercițiu.

Observația 2.8 Sub forma (3) cu $u, v \in C$ sau sub forma (4) cu $A, B \in C$, din propoziția anterioară, definiția independenței a două v.a. simple se poate enunța și pentru v.a. cu valori complexe.

Propoziția 2.9 Dacă $f, g : X \rightarrow C$ sunt v.a. simple independente iar $F, G : C \rightarrow C$ atunci $F \circ f$ și $G \circ g$ sunt independente.

Demonstrație. Avem:

$$\begin{aligned}
 & p((F \circ f = v) \cap (G \circ g = u)) \\
 &= p((f \in F^{-1}(v)) \cap (g \in G^{-1}(u))) \\
 &= p(f \in F^{-1}(v)) \cdot p(g \in G^{-1}(u)) \\
 &= p(F \circ f = v)p(G \circ g = u)
 \end{aligned}$$

Q.E.D.

De exemplu $(f - a)^2$ și $(g - b)^2$ sunt independente pentru orice a și b .

Teorema 2.10 (proprietăți ale mediei) Fie (X, Ω, p) un câmp de probabilitate finit și f, f_1, f_2, \dots variabile aleatoare pe X . Atunci:

1) $f \equiv C$ (constant) $\Rightarrow M(f) = C$

2) $\alpha = \text{const} \Rightarrow M(\alpha f) = \alpha M(f)$

3) $M(f_1 + f_2) = M(f_1) + M(f_2)$ și de aici $M(f_1 + f_2 + \dots + f_n) = M(f_1) + M(f_2) + \dots + M(f_n)$, $M(af + bg) = aM(f) + bM(g)$, și $M(f - M(f)) = 0$, (a și b fiind constante)

4) Dacă f_1, f_2 sunt independente atunci $M(f_1 f_2) = M(f_1) \cdot M(f_2)$.

5) $f \geq 0 \Rightarrow M(f) \geq 0$.

6) $f \geq g \Rightarrow M(f) \geq M(g)$.

7) $|M(f)| \leq \max |f|$.

Demonstrație. 1) și 2) sunt evidente.

3) Fie v_1, v_2, \dots, v_n valorile lui f_1 și $A_i = f_1^{-1}(v_i)$ iar u_1, \dots, u_m valorile lui f_2 și $B_j = f_2^{-1}(u_j)$. Variabila aleatoare $f_1 + f_2$ ia valorile $v_i + u_j$ pe evenimentele $A_i \cap B_j$ cu probabilitățile $p_{ij} = p(A_i \cap B_j) = p(f_1 = v_i \cap f_2 = u_j)$. Avem acum

$$\begin{aligned}
 M(f_1 + f_2) &= \sum_{i=1..n, j=1..m} (v_i + u_j) p_{ij} = \sum_{i=1}^n v_i \sum_{j=1}^m p_{ij} + \sum_{j=1}^m u_j \sum_{i=1}^n p_{ij} = \\
 &= \sum v_i p(f_1 = v_i) + \sum u_j p(f_2 = u_j) = M(f_1) + M(f_2)
 \end{aligned}$$

Am folosit

$$\begin{aligned}
 \sum_{j=1}^m p_{ij} &= \sum_{j=1}^m p(A_i \cap B_j) = (\text{fiind disjuncte}) \\
 &= p(\bigcup_{j=1..m} (A_i \cap B_j)) = p(A_i \cap (\bigcup_{j=1..m} B_j)) = \\
 &= p(A_i \cap X) = p(A_i) = p(f_1 = v_i)
 \end{aligned}$$

și în mod analog cealaltă sumă este $\sum_{i=1}^n p_{ij} = p(f_2 = u_j)$.

4) $M(f_1 f_2) = \sum_{i,j} v_i u_j p(f_1 = v_i \cap f_2 = u_j) =$ (din independența variabilelor aleatoare)
 $= \sum_{i,j} v_i u_j p(f_1 = v_i) p(f_2 = u_j) = \sum_i v_i p(f_1 = v_i) \sum_j u_j p(f_2 = u_j) = M(f_1) M(f_2)$

5), 6), 7) sunt evidente.

Q.E.D.

Teorema 2.11 (proprietățile dispersiei). 1) $f \equiv C$ a.p.t.(constant) $\Leftrightarrow D(f) = 0$

2) $D(af) = a^2 D(f)$ a fiind o constantă.

3) Dacă f_1 și f_2 sunt independente atunci $D(f_1 + f_2) = D(f_1) + D(f_2)$

Demonstrație. 1) $D(f) = \sum (v_i - M(f))^2 p(f = v_i)$ și suma este zero numai dacă $v_i = M(f)$ pentru orice i cu $p_i \neq 0$, adică $f \equiv M(f)$ a.p.t.

2) Avem

$$D(af) = M((af - M(af))^2) = M(a^2(f - M(f))^2) = a^2 M((f - M(f))^2) = a^2 D(f)$$

3) În formulele de mai jos utilizăm faptul că f_1 și f_2 independente implică $(f_1 - a)^2$ și $(f_2 - b)^2$ sunt independente, precum și proprietățile mediei:

$$\begin{aligned} D(f_1 + f_2) &= M(((f_1 + f_2) - M(f_1 + f_2))^2) = \\ &= M(((f_1 - M(f_1)) + (f_2 - M(f_2)))^2) = \\ &= M((f_1 - M(f_1))^2 + (f_2 - M(f_2))^2 + 2M((f_1 - M(f_1))(f_2 - M(f_2)))) = \\ &= M((f_1 - M(f_1))^2) + M((f_2 - M(f_2))^2) + \underbrace{2M(f_1 - M(f_1)) \cdot M(f_2 - M(f_2))}_{=0} = \\ &= D(f_1) + D(f_2) \end{aligned}$$

Q.E.D.

O generalizare a punctului 3) este:

Propoziția 2.12 Fie f_1, f_2, \dots, f_n variabile aleatoare independente două câte două. Atunci

$$D(f_1 + f_2 + \dots + f_n) = D(f_1) + D(f_2) + \dots + D(f_n)$$

Demonstrația este analoagă celei pentru două funcții și e lăsată ca exercițiu.

Definiția 2.13 Dacă f este o variabilă aleatoare, expresia $\frac{f - M(f)}{\sqrt{D(f)}}$ se numește deviația standard a lui f .

Se vede imediat că $M(\frac{f - M(f)}{\sqrt{D(f)}}) = 0$ și $D(\frac{f - M(f)}{\sqrt{D(f)}}) = 1$.

Teorema 2.14 (proprietăți ale funcției caracteristice) Fie f_c funcția caracteristică a variabilei aleatoare f . Atunci:

- 1) f_c este continuă (chiar uniform continuă) pe \mathbb{R}
- 2) $f_c(0) = 1$, $|f_c(t)| \leq 1$ pentru $-\infty < t < \infty$
- 3) Dacă $g = a \cdot f + b$, cu a și b constante, atunci $g_c(t) = e^{\sqrt{-1}bt} \cdot f_c(at)$
- 4) Dacă f_1 și f_2 sunt independente, atunci $(f_1 + f_2)_c(t) = f_{1c}(t) \cdot f_{2c}(t)$
- 5) $M_k(f) = \frac{1}{(\sqrt{-1})^k} f_c^{(k)}(0)$
- 6) Dacă $f_{1c} \equiv f_{2c}$ atunci f_1 și f_2 au aceleași diagrame standard, adică $f_1 = f_2$ a.p.t.

Demonstrație. 1) f_c este o sumă finită de funcții continue $f_c(t) = \sum p_k e^{\sqrt{-1}v_k t}$ deci este o funcție continuă. Având derivata mărginită rezultă că este uniform continuă.

2) $g_c(t) = \sum p_k e^{\sqrt{-1}t(\alpha v_k + \beta)} = e^{\sqrt{-1}t\beta} \sum p_k e^{\sqrt{-1}t\alpha v_k} = e^{\sqrt{-1}t\beta} f_c(\alpha t)$.

$$3) f_c(0) = \sum p_k e^{\sqrt{-1}0v_k} = \sum p_k = 1 \text{ și } |f_c(t)| \leq \sum p_k |e^{\sqrt{-1}v_k t}| = \sum p_k = 1$$

4) Dacă f_1 și f_2 sunt independente atunci e^{if_1} și e^{if_2} sunt independente.

Acum observăm că

$$\begin{aligned} (f_1 + f_2)_c(t) &= M(e^{it(f_1+f_2)}) = M(e^{itf_1} \cdot e^{itf_2}) = (\text{find indep.}) \\ &= M(e^{itf_1})M(e^{itf_2}) = f_{1c}(t) \cdot f_{2c}(t) \end{aligned}$$

5) $f'_c(t) = i \sum_k p_k v_k e^{itv_k}$ și deci $f'_c(0) = i \sum_k p_k v_k = iM_1(f)$. Analog prin derivare succesivă se obține $f^{(n)}(0) = i^n M_n(f)$.

6) Avem, folosind diagrame reduse, $\sum p_k e^{iv_k t} \equiv \sum q_l e^{iu_l t}$, unde v_k sunt diferite între ele și u_l sunt diferite între ele. Dacă nu se reduc toți termenii, adică pentru fiecare k să existe l astfel ca $p_k = q_l, v_k = u_l$ atunci după toate reducerile posibile egalitatea devine: $\sum r_s e^{iw_s t} \equiv 0$, ($1 \leq s \leq m$) unde r_s sunt nenule, iar w_s sunt diferite între ele. Derivând de mai multe ori expresia de mai sus și punând $t=0$ în derivate găsim:

$$\begin{cases} r_1 + r_2 + \dots + r_m = 0 \\ r_1 w_1 + r_2 w_2 + \dots + r_m w_m = 0 \\ \dots\dots\dots \\ r_1 w_1^{m-1} + r_2 w_2^{m-1} + \dots + r_m w_m^{m-1} = 0 \end{cases}$$

ceea ce nu se poate decât dacă $r_1 = r_2 = \dots = r_m = 0$. Contradicția obținută arată că reducerea are loc, deci f_1 și f_2 iau aceleași valori, cu aceleași probabilități.

QED.

Definiția 2.15 Spunem că n variabile aleatoare f_1, \dots, f_n sunt independente dacă pentru orice intervale I_1, \dots, I_k , și orice alegere f_{i_1}, \dots, f_{i_k} , $k \leq n$, rezultă că evenimentele $f_{i_1}^{-1}(I_1), f_{i_2}^{-1}(I_2), \dots, f_{i_k}^{-1}(I_k)$ sunt independente.

Exercițiul 2.16

Dacă variabilele aleatoare f_1, \dots, f_n sunt independente atunci

$$(f_1 + f_2 + \dots + f_n)_c(t) = f_{1c}(t) \cdot f_{2c}(t) \cdot \dots \cdot f_{nc}(t)$$

2.2 Spațiul de probabilitate produs

Fie (X_1, Ω_1, p_1) și (X_2, Ω_2, p_2) două câmpuri de probabilitate finite. Atunci pe mulțimea $X = X_1 \times X_2$ se poate defini o algebră de mulțimi astfel:

$$\begin{aligned} \Omega = \{A \subset X | A \text{ este reuniune finită de elemente de forma } A_1 \times A_2, \\ \text{cu } A_1 \in \Omega_1 \text{ și } A_2 \in \Omega_2\} \end{aligned}$$

Ținând seama că Ω_1 și Ω_2 sunt algebre de evenimente, rezultă că și Ω este. Definim $p : \Omega \rightarrow R$ prin $p(A_1 \times A_2) = p_1(A_1) \cdot p_2(A_2)$. Dacă A este o reuniune disjunctă de mulțimi

de forma $\cup A_{i_i} \times A_{2_i}$ atunci $p(A) = \sum p_1(A_{1_i}) \cdot p_2(A_{2_i})$. Se demonstrează că p este o probabilitate pe Ω iar câmpul (X, Ω, p) se numește câmpul produs al câmpurilor (X_1, Ω_1, p_1) și (X_2, Ω_2, p_2) . În mod asemănător se definește produsul unui număr finit de câmpuri de probabilitate $(X_1, \Omega_1, p_1), \dots, (X_n, \Omega_n, p_n)$. Spațiul total este $X = X_1 \times X_2 \times \dots \times X_n$ iar mulțimea de evenimente se definește prin $A \in \Omega \Leftrightarrow A$ este o reuniune finită de elemente de forma $A_1 \times A_2 \times \dots \times A_n$ cu $A_i \in \Omega_i$ pentru orice i . Se arată că Ω este o algebră de evenimente pe care se poate defini o probabilitate care pe evenimentele de forma $A_1 \times A_2 \times \dots \times A_n$ are valoarea $p(A_1 \times A_2 \times \dots \times A_n) = p_1(A_1) \cdot p_2(A_2) \cdot \dots \cdot p_n(A_n)$. Este foarte ușor de arătat că evenimentele

$$\begin{aligned} B_1 &= A_1 \times X_2 \times \dots \times X_n \\ B_2 &= X_1 \times A_2 \times \dots \times X_n \\ &\dots\dots\dots \\ B_n &= X_1 \times X_2 \times \dots \times A_n \end{aligned}$$

sunt independente(exercițiu). De asemenea următoarele variabile aleatoare

$$\begin{aligned} pr_1(x_1, x_2, \dots, x_n) &= x_1 \\ pr_2(x_1, x_2, \dots, x_n) &= x_2 \\ &\dots\dots\dots \\ pr_n(x_1, x_2, \dots, x_n) &= x_n \end{aligned}$$

numite proiecțiile canonice ale lui X , sunt independente(exercițiu). De asemenea dacă $f_k : X_k \rightarrow R$ sunt variabile aleatoare atunci variabilele aleatoare $F_k : X \rightarrow R$ definite prin $F_k(x_1, x_2, \dots, x_n) = f_k(x_k)$ sunt independente(exercițiu).

2.3 Rezumat

Informațiile dintr-un spațiu probabilizat se extrag prin funcții, numite variabile aleatoare. In cazul variabilelor aleatoare ce iau un număr finit de valori (numite și variabile simple), media este definită într-un mod analog cu media notelor la școală. Dispersia e definită ca o măsură a împrăstierii valorilor în jurul mediei. Funcția caracteristică este o noțiune auxiliară importantă pentru calculul mediei, dispersiei, etc. Ea definește unic diagrama variabilei aleatoare. Proprietățile acestor mărimi sunt listate în cele trei teoreme anterioare. E de remarcat comportarea acestor mărimi la adunarea variabilelor aleatoare independente. Un exemplu important de spațiu probabilizat și variabile aleatoare independente este spațiul produs și proiecțiile pr_k .

FORMULE UTILIZATE FRECVENT:

- Definițiile mediei, momentelor, dispersiei, etc. pentru calculul direct al lor pag(15)
- Proprietățile mediei, dispersiei, funcției caracteristice. De remarcat:

- i) $M(af + bg) = aM(f) + bM(g)$;
- ii) $f \geq 0 \Rightarrow M(f) \geq 0$;
- iii) $|M(f)| \leq \max |f|$;
- iv) f_1, f_2 independente $\Rightarrow M(f_1 f_2) = M(f_1) M(f_2)$, $D(f_1 + f_2) = D(f_1) + D(f_2)$ și $(f_1 + f_2)_c(t) = f_{1c}(t) f_{2c}(t)$;
- v) $D(f) = M_2(f) - M^2(f)$
- vi) $M_k(f) = \frac{1}{(\sqrt{-1})^k} f_c^{(k)}(0)$.

2.4 Exerciții

1. Fie variabila aleatoare

$$f = \begin{pmatrix} -1 & 2 & 4 & 10 \\ \frac{1}{2} & \frac{1}{10} & \frac{3}{10} & \frac{1}{10} \end{pmatrix}$$

Să se calculeze media, dispersia și momentul de ordinul 3.

Indicație. Se utilizează definițiile.

2 Fie v.a.

$$f = \begin{pmatrix} -2 & -1 & 2 & 4 & 10 \\ \frac{1}{5} & \frac{1}{10} & \frac{3}{10} & \frac{1}{10} & \alpha \end{pmatrix}$$

Sa se determine α astfel ca f sa fie o v.a. discreta. Să se calculeze media și dispersia ei.

3. Stiind ca pentru v.a. $f = \begin{pmatrix} -1 & 0 & 1 \\ p_1 & p_2 & p_3 \end{pmatrix}$ avem $M(f) = -\frac{1}{5}$ și $M(f^2) = \frac{33}{5}$ să se determine p_1, p_2, p_3 .

4. Fie o v.a. f cu media 10 . Se cere $b \in R$ astfel ca $M(f + b) = 0$.

5. În N urne se găsesc câte n bile numerotate de la 1 la n , în fiecare. Se extrage câte o bilă din fiecare urnă și se notează cu f valoarea celui mai mare număr obținut. Se cere explicitarea probabilităților $p(f=k)$, valoarea medie m_n și limita expresiei m_n/n când $n \rightarrow \infty$.

Indicație. $p(f \leq k)$ =probabilitatea ca din fiecare urnă să se extragă un număr mai mic sau egal cu k . Numărul cazurilor favorabile este k^N iar numărul celor posibile este n^N , deci probabilitatea este $p(f \leq k) = \frac{k^N}{n^N}$. De aici deducem $p(f = k) = \frac{k^N}{n^N} - \frac{(k-1)^N}{n^N}$. Aplicând definiția mediei găsim $m_n = n - \sum_{k=1}^n \frac{(k-1)^N}{n^N}$ deci $\frac{m_n}{n} \rightarrow 1 - \int_0^1 x^N dx$.

6. O întreprindere recrutează un nou angajat. Se prezintă n candidați într-o ordine aleatoare și primul candidat care trece testul proous de comisie este angajat. Probabilitatea de a trece testul este p pentru fiecare candidat. Fie variabila aleatoare f ce ia valoarea j dacă al j -ulea candidat este admis, și valoarea 0 dacă nu este nimeni admis.

Să se determine $p(f=j)$, valoarea medie $M(f)$ și dispersia $D(f)$.

Indicație. Ca să fie admis candidatul j trebuie ca toți candidații $1, 2, \dots, j-1$ să fie respinși iar j să fie admis. Probabilitatea acestui eveniment este $p(f = j) = q^{j-1}p$, unde $q = 1 - p$.

Probabilitatea ca toți să fie respinși este q^n . Prin urmare f are diagrama:

$$\begin{pmatrix} 1 & 2 & \dots & k & \dots & n \\ p & qp & \dots & q^{k-1}p & \dots & q^n \end{pmatrix}$$

Este foarte simplu acum de a calcula funcția caracteristică, media și dispersia.

7. O persoană se deplasează aleator plecând din punctul 0. Cu probabilitatea p face un pas în față și cu probabilitatea $q=1-p$ face un pas înapoi. Valoarea unui pas este de 1m iar ritmul este de un pas pe minut. Fie f variabila aleatoare egală cu abscisa la care se ajunge după o oră (un pas în față este +1m iar un pas în spate este -1m). Se cer probabilitățile $p(f=k)$, valoarea medie și dispersia lui f .

Indicație. Dacă se fac în față k pași, se ajunge în poziția $k-(60-k)=2k-60$. Probabilitatea de a face k pași în față este $C_{60}^k p^k q^{60-k}$ (vezi în lecția 4 despre schema Bernoulli). Deci $p(f = 2k - 60) = C_{60}^k p^k q^{60-k}$. Deci funcția caracteristică este $f_c(t) = \sum_{k=0}^{60} C_{60}^k p^k q^{60-k} e^{i(2k-60)t} = e^{-60ti} (pe^{2ti} + q)^{60}$. Cu formulele uzuale se deduc acum media, momentele, dispersia.

8. Probabilitatea ca un bec să reziste la suprasarcină este $p=0,90$. Se fac 5 încercări de becuri. fie f variabila aleatoare care reprezintă numărul de becuri care au rezistea. Se cere media și dispersia lui f .

9. Se aruncă două zaruri și se notează cu f suma punctelor obținute. Se cere media și dispersia lui f .

Indicație. Se poate obține suma 2 din varianta (1,1), adică un caz din 36. Se poate obține suma 3 din variantele (1,2) sau (2,1), deci 2 cazuri din 36, etc. De aici rezultă imediat legea lui f etc.

10. O urnă conține n bile albe și n bile negre. Se extrag una câte una bilele până se epuizează urna. Fie f variabila aleatoare care pentru o variantă de extragere succesivă a tuturor bilelor ia valoarea k dacă prima bilă albă apare la extragerea k .

a) Să se determine legea lui f

b) Se cere media lui f .

Indicație. k ia valori între 1 și $n+1$. Ca f să ia valoarea k trebuie să iasă primele $k-1$ bile negre. Probabilitatea ca să iasă prima bilă neagră este $\frac{n}{2n}$. Condiționat de aceasta, probabilitatea ca a doua bilă să iasă neagră este $\frac{n-1}{2n-1}$. Deci probabilitatea ca primele două bile să iasă negre este $\frac{n(n-1)}{2n(2n-1)}$. Dacă primele două bile au ieșit negre, probabilitatea ca a treia să iasă neagră este $\frac{n-2}{2n-2}$. Deci probabilitatea ca primele trei bile să iasă negre este $\frac{n(n-1)(n-2)}{2n(2n-1)(2n-2)}$. Analog găsim probabilitatea ca primele $k-1$ bile să iasă negre $\frac{n(n-1)(n-2)\dots(n-k+2)}{2n(2n-1)\dots(2n-k+2)} = \frac{C_n^{k-1}}{C_{2n}^{k-1}}$. Știind aceasta, probabilitatea ca la extragerea k să iasă bila albă este $\frac{n}{2n-k+1}$. Deci $p(f = k) = \frac{n}{2n-k+1} \frac{C_n^{k-1}}{C_{2n}^{k-1}}$.

11. Două variabile aleatoare f și g iau aceleași valori v_1, v_2, \dots, v_n . Se mai știe că $M(f) = M(g)$, $M_2(f) = M_2(g), \dots, M_{n-1}(f) = M_{n-1}(g)$. Să se arate că $p(f = v_k) = p(g = v_k)$ pentru orice $k \in \overline{1..n}$.

Indicație. Fie $p_k = p(f = v_k)$ și $p'_k = p(g = v_k)$. Avem:

$$M_k(f) = p_1 v_1^k + p_2 v_2^k + \dots + p_n v_n^k = M_k(g) = p'_1 v_1^k + p'_2 v_2^k + \dots + p'_n v_n^k$$

pentru $k=0,1,2,\dots,n$. Prin urmare p_1, p_2, \dots, p_n și p'_1, p'_2, \dots, p'_n satisfac același sistem nesingular.

12. Să se arate că $D(f + C) = D(f)$, C fiind o constantă.

Indicație. Se aplică definiția dispersiei și se ține seama că $M(f + C) = M(f) + C$.

13. Fie X o variabilă aleatoare finită cu media m și dispersia σ^2 . Cum trebuie să fie a și b (constante) ca variabila aleatoare $af+b$ să aibă media zero și dispersia unu?

Indicație. $0 = M(af+b) = aM(f) + b = am + b$ și $1 = D(af+b) = D(af) = a^2 D(f) = a^2 \sigma^2$

14. Două variabile aleatoare independente ξ, η au distribuțiile

$$\xi \left(\begin{array}{cccc} 0 & 1 & 2 & 3 \\ \frac{1}{8} & \frac{1}{4} & \frac{a}{4} & \frac{a}{2} \end{array} \right) \quad \eta \left(\begin{array}{cccc} 0 & 1 & 2 & 3 \\ \frac{1}{2} & b & \frac{1}{8} & \frac{1}{4} \end{array} \right)$$

a) Să se determine a și b

b) Care este distribuția variabilei $\xi - \eta$?

c) Care este distribuția variabilei $2\xi + 3\eta$?

Indicație. Pentru ξ avem $\frac{1}{8} + \frac{1}{4} + \frac{a}{4} + \frac{a}{4} = 1$ și analog pentru η . Pentru $\xi - \eta$ se ține seama că ξ și η sunt independente, deci, de exemplu $p(\xi = 1 \text{ și } \eta = 2) = p(\xi = 1) \cdot p(\eta = 2)$, etc.

15. Doi jucători, A și B, convin asupra următoarei reguli de joc: la o aruncare cu zarul, dacă apar fețele 1 sau 2 atunci primul jucător dă celui de al doilea 3 lei, iar în caz contrar al doilea jucător dă primului 1 leu. Care este câștigul mediu al fiecărui jucător după n aruncări?

Indicație. Probabilitatea de câștig a lui A este $1/3$ iar a lui B este de $2/3$. Se scriu variabilele aleatoare care reprezintă câștigul pentru A și B, ca la problema 7.

16. Trei jucători A, B, C pun jos câte o sumă de a , b , respectiv c lei, apoi aruncă o monedă, pe rând, în ordinea A, B, C, A, B, C, ... până apare fața. Primul jucător la care apare fața ia toți banii, iar dacă nu a apărut după ce fiecare a aruncat de 3 ori, se anulează jocul.

a) Care sunt câștigurile medii ale fiecărui jucător după 3 jocuri?

b) Cum ar trebui să fie sumele a , b , c ca jocul să fie echitabil?

Lecția 3

σ Câmpuri de probabilitate

Fie (X, Ω, p) un câmp de probabilitate. Dacă algebra de evenimente verifică în plus axioma:

Oricare ar fi șirul de evenimente $(A_n)_{n \in \mathbb{N}}$, cu $A_n \in \Omega$ atunci $\cup A_n \in \Omega$,

spunem că Ω este σ -algebră sau un σ -câmp de evenimente. Dacă p este o probabilitate pe Ω care în plus verifică relația

$$p\left(\bigcup_{n=1, \infty} A_n\right) = \sum_{n=1}^{\infty} p(A_n)$$

pentru orice șir A_n de evenimente disjuncte, atunci spunem că p este o σ probabilitate. În aceste condiții tripletul (X, Ω, p) se numește σ -câmp de probabilitate.

Teorema 3.1 Fie (X, Ω, p) un σ -câmp de probabilitate. Atunci:

- 1) Dacă $A_n \in \Omega$ pentru orice n , avem și $\bigcap_{n=1, \infty} A_n \in \Omega$.
- 2) Dacă $A_1 \supseteq A_2 \supseteq \dots \supseteq A_n \supseteq \dots$ atunci $p\left(\bigcap_{n=1, \infty} A_n\right) = \lim_{n \rightarrow \infty} p(A_n)$
- 4) Dacă $A_1 \subseteq A_2 \subseteq \dots \subseteq A_n \subseteq \dots$ atunci $p\left(\bigcup_{n=1, \infty} A_n\right) = \lim_{n \rightarrow \infty} p(A_n)$
- 3) $p\left(\bigcup_{n=1, \infty} A_n\right) \leq \sum_{n=1}^{\infty} p(A_n)$

Demonstrație. p fiind o probabilitate, toate proprietățile pentru probabilitățile finit aditive rămân adevărate.

1) $\bigcap_{n=1, \infty} A_n = \overline{\bigcup_{n=1, \infty} \bar{A}_n}$, și deoarece $\bar{A}_n \in \Omega$, rezultă că reuniunea mulțimilor \bar{A}_n este în Ω , deci și complementara acestei reuniuni, adică intersecția mulțimilor A_n este în Ω .

2) Mulțimile $A_k - A_{k+1}$ sunt disjuncte între ele și sunt disjuncte de $\bigcap_{k=1 \dots \infty} A_k$ iar $A_1 =$

$\left(\bigcap_{k=1.. \infty} A_k\right) \cup (A_1 - A_2) \cup (A_2 - A_3) \cup \dots$ este o reuniune disjunctă, de unde:

$$p(A_1) = p\left(\bigcap_{k=1.. \infty} A_k\right) + \sum_{k=1}^{\infty} p(A_k - A_{k+1})$$

de unde:

$$\begin{aligned} p\left(\bigcap_{k=1.. \infty} A_k\right) &= p(A_1) - \lim_{n \rightarrow \infty} \sum_{k=1}^{n-1} p(A_k - A_{k+1}) = \\ &= \lim_{n \rightarrow \infty} (p(A_1) - \sum_{k=1}^{n-1} p(A_k - A_{k+1})) = \lim_{n \rightarrow \infty} p(A_n) \end{aligned}$$

3) Demonstrația este asemănătoare cu cea de la punctul 2).

4) $p(A_1 \cup A_2) = p(A_1) + p(A_2) - p(A_1 \cap A_2) \leq p(A_1) + p(A_2)$ și prin inducție se demonstrează că $p(A_1 \cup A_2 \cup \dots \cup A_n) \leq p(A_1) + p(A_2) + \dots + p(A_n)$. Acum având $A_1 \subseteq A_1 \cup A_2 \subseteq \dots \subseteq (A_1 \cup A_2 \cup \dots \cup A_n) \subseteq \dots$, rezultă:

$$\begin{aligned} p\left(\bigcup_{k=1.. \infty} A_k\right) &= \lim_{n \rightarrow \infty} p(A_1 \cup A_2 \cup \dots \cup A_n) \leq \\ &\leq \lim_{n \rightarrow \infty} (p(A_1) + p(A_2) + \dots + p(A_n)) = \sum_{k=1}^{\infty} p(A_k) \end{aligned}$$

QED.

Exemplul 3.2 Fie $X = [0, 1] \times [0, 1]$ și Ω mulțimea submulțimilor lui X cu frontiera formată dintr-un număr finit de segmente. Este clar că Ω este o mulțime închisă la reuniune, intersecție și complementară în X , deci este o algebră de mulțimi. O probabilitate pe Ω se poate defini prin formula $p(A) = \frac{\text{aria}(A)}{\text{aria}(X)} = \text{aria}(A)$. Pe de altă parte Ω nu este o σ algebră deoarece spre exemplu un cerc plin C nu face parte din Ω deși este reuniunea unui șir crescător A_n de mulțimi din Ω , cu A_n egal cu poligonul regulat cu 2^n laturi înscris în C . Se demonstrează că există o cea mai mică σ algebră ce conține pe Ω , fie ea B , iar probabilitatea p (sau aria p) se poate extinde la o σ probabilitate pe B . Probabilitatea (sau aria) astfel definită se numește probabilitatea Lebesgue (aria sau măsura Lebesgue).

În mod asemănător se definește măsura Lebesgue pe un segment sau pe un paralelipiped în spațiu, ca o extindere a noțiunii de lungime a unui segment sau de volum al unui corp mărginit de un număr finit de fețe plane.

3.1 Variabile aleatoare pe σ câmpuri de probabilitate

Definiția 3.3 Fie (X, Ω, p) un σ câmp de probabilitate. O funcție $f : X \rightarrow \mathbb{R}$ se numește variabilă aleatoare dacă pentru orice număr $c \in \mathbb{R}$, mulțimea $\{\omega \in \Omega \mid f(\omega) < c\}$ este un eveniment din Ω . Adesea vom nota mulțimea de mai sus prin $\{f < c\}$ sau $f^{-1}(-\infty, c)$.

Propoziția 3.4 $f: X \rightarrow R$ este variabilă aleatoare dacă și numai dacă este îndeplinită una din condițiile următoare:

- a) $\forall c \in R \Rightarrow \{f < c\} \in \Omega$
- b) $\forall c \in R \Rightarrow \{f \leq c\} \in \Omega$
- c) $\forall c \in R \Rightarrow \{f > c\} \in \Omega$
- d) $\forall c \in R \Rightarrow \{f \geq c\} \in \Omega$
- e) $\forall \alpha, \beta \in R \Rightarrow \{\alpha \leq f < \beta\} \in \Omega$

Demonstrație. a) este definiția curentă a variabilei aleatoare. $a) \rightarrow b) \{f \leq c\} = \bigcap_n \{f < c + \frac{1}{n}\} \in \Omega$. $b) \rightarrow c) \{f > c\} = X - \{f \leq c\} \in \Omega$ $c) \rightarrow d) \{f \geq c\} = \bigcap_n \{f > c - \frac{1}{n}\} \in \Omega$ $d) \rightarrow e) \{\alpha \leq f < \beta\} = \{f \geq \alpha\} - \{f \geq \beta\} \in \Omega$ $e) \rightarrow a) \{f < c\} = \bigcup_n \{-n \leq f < c\} \in \Omega$

QED.

Prin urmare o variabilă aleatoare este caracterizată prin faptul că pentru orice interval $I \subset R$, închis, deschis, semiînchis, finit sau nu, avem relația $f^{-1}(I) \in \Omega$. Spre deosebire de cazul variabilelor aleatoare simple studiate în lecția trecută, condiția $f^{-1}(v) \in \Omega$ pentru $v \in R$ nu este suficientă ca f să fie variabilă aleatoare.

Teorema 3.5 Fie $f: X \rightarrow R$ o variabilă aleatoare și $\alpha \in R$. Atunci următoarele funcții sunt variabile aleatoare: a) $\alpha + f$; b) $\alpha \cdot f$; c) f^2 ; d) $\frac{1}{f}$ dacă $f \neq 0$; e) $|f|$.

Demonstrație. a) $\{\alpha + f < c\} = \{f < c - \alpha\} \in \Omega$. b) $\{\alpha f < c\} = \{f < \frac{c}{\alpha}\}$ dacă $\alpha > 0$ și $\{\alpha f < c\} = \{f > c/\alpha\}$ dacă $\alpha < 0$. În ambele cazuri se obțin mulțimi în Ω . c) $\{f^2 < c\} = \{-\sqrt{c} < f < \sqrt{c}\} \in \Omega$. d) $\{\frac{1}{f} < c\} = (\{f > 0\} \cap \{cf > 1\}) \cup (\{f < 0\} \cap \{cf < 1\}) \in \Omega$. e) exercițiu.

QED.

Teorema 3.6 Fie $(f_i)_{i \in N}$ un șir de variabile aleatoare. Atunci:

- a) $h(\omega) = \sup_{1 \leq i < \infty} f_i(\omega)$ este o variabilă aleatoare.
- b) $g(\omega) = \inf_{1 \leq i < \infty} f_i(\omega)$ este o variabilă aleatoare.
- c) Dacă există $\lim_{i \rightarrow \infty} f_i(\omega) = f(\omega)$ atunci f este o variabilă aleatoare.

Demonstrație. a) $\{h > c\} = \bigcup_{i=1}^{\infty} \{f_i > c\} \in \Omega$ b) $\{g < c\} = \bigcup_{i=1}^{\infty} \{f_i < c\} \in \Omega$ c) Fie $\bar{f}_i = \max_{j \geq i} f_j$. Conform cu a) \bar{f}_i este o variabilă aleatoare. $f = \lim_{i \rightarrow \infty} f_i = \min_i \bar{f}_i$ care este conform cu b) o variabilă aleatoare.

QED.

Teorema 3.7 Dacă f și g sunt variabile aleatoare atunci $f+g$, fg , $\frac{f}{g}$ sunt variabile aleatoare.

Demonstrație. $\{f+g < c\} = \{f < c-g\} = \bigcup_{r \in Q} (\{f < r\} \cap \{c-g > r\}) \in \Omega$, știind că toate numerele $r \in Q$ se pot pune într-un șir. Deci $f+g$ este o variabilă aleatoare. Analog se

procedează cu $f - g$. Putem scrie $f \cdot g = \frac{1}{4} [(f + g)^2 - (f - g)^2]$ și conform punctului anterior f este o variabilă aleatoare. Deoarece $\frac{f}{g} = f \cdot \frac{1}{g}$ rezultă că și $\frac{f}{g}$ este o variabilă aleatoare.

QED.

Din cele de mai sus vedem că dacă a și b sunt constante atunci $af + bg$ este variabilă aleatoare odată cu f și g . Deci mulțimea variabilelor aleatoare pe o aceeași σ algebră formează un spațiu vectorial. Se poate arăta că dacă f este o variabilă aleatoare și $F : R \rightarrow R$ este o funcție continuă atunci $F \circ f : X \rightarrow R$ este o variabilă aleatoare (demonstrația este lăsată ca exercițiu)

Dacă f este o v.a. simplă atunci există o partiție a spațiului $X = \bigcup_{i=1}^n E_i$, $E_i \cap E_j = \emptyset$ astfel că f este constantă pe E_i . Dacă f și g sunt variabile aleatoare simple atunci $f \pm g$, fg , f^2 , f^α , $|f|$ sunt simple. Orice variabilă aleatoare este limita unui șir de variabile aleatoare simple.

Teorema 3.8 Fie (X, Ω) o σ algebră și f o variabilă aleatoare pe X . Atunci există un șir de variabile aleatoare simple $f_n \rightarrow f$. Convergența este înțeleasă în sensul că pentru orice $x \in X$ rezultă $\lim_{n \rightarrow \infty} f_n(x) = f(x)$ în R (convergență simplă).

Demonstrație. Fie $E_{n,j} = \{\omega \in X \mid \frac{j-1}{2^n} \leq f < \frac{j}{2^n}\}$ $-n2^n \leq j \leq n2^n$. Fie acum f_n definite astfel:

$$f_n(\omega) = \begin{cases} \frac{j-1}{2^n}, & \omega \in E_{n,j} \\ n, & f(\omega) \geq n \\ -n, & f(\omega) < -n \end{cases}$$

f_n este o variabilă aleatoare simplă și se vede că $f_n(\omega) \rightarrow f(\omega)$. De asemenea se vede că f_n tinde uniform la f dacă f este mărginită. În plus dacă $f > 0$, atunci f_n tinde crescător la f .

QED.

3.2 Media unei variabile aleatoare oarecare

Dacă variabila aleatoare este simplă atunci media ei a fost definită în lecția anterioară. Fie $f(\omega) = v_i$ dacă $\omega \in E_i$, pentru $i=1,2,\dots,n$. Atunci media lui f se definește prin $M(f) = \sum_i v_i p(E_i)$. Vom mai nota această expresie prin $\int_X f \cdot dp$ sau $\int_X f(\omega) \cdot dp(\omega)$. Dacă variabila aleatoare nu este simplă atunci se aproximează f prin variabile aleatoare simple iar media lui f se aproximează prin media acestora. În continuare vom folosi noțiunea de convergență aproape peste tot (pe scurt a.p.t.).

Definiția 3.9 Un șir $(f_n)_{n \in \mathbb{N}}$ de v.a. se spune că este convergent a.p.t. la variabila aleatoare f dacă și numai dacă mulțimea $\left\{x \in X \mid \lim_{n \rightarrow \infty} f_n(x) \neq f(x)\right\}$ are probabilitatea 0.

Putem descrie acum felul în care se definește media în general:

a) Considerăm un șir de variabile aleatoare simple $(f_n)_{n \in N}$ care converge la f a.p.t. (un asemenea șir există conform teoremei precedente, și care converge peste tot la f).

b) Stim că $M(f_n) = \int_X f_n dp$ există. Dacă șirul e Cauchy în sensul că $\forall \epsilon > 0, \exists n(\epsilon) \in N, a.î \forall n, m \in N, n \geq n(\epsilon), m \geq n(\epsilon) \Rightarrow \int_X |f_n - f_m| dp < \epsilon$, atunci se arată că există $\lim_{n \rightarrow \infty} \int_X f_n dp$ și această limită se numește media lui f . Ea se notează cu $\int_X f dp$ sau $M(f)$ și nu depinde de șirul de variabile aleatoare simple care tind Cauchy la f . În cazul particular când f este simplă, se regăsește vechea definiție.

Se demonstrează că proprietățile mediei demonstrate anterior pentru variabile aleatoare simple se păstrează prin trecere la limită, în general. Avem deci (scriind $\int_X f dp$ în loc de $M(f)$):

1)

$$f = c \Rightarrow \int_X f dp = c$$

2)

$$\int_X (\alpha f + \beta g) dp = \alpha \int_X f dp + \beta \int_X g dp$$

3)

$$\int_X (f_1 + f_2 + \dots + f_n) dp = \int_X f_1 dp + \int_X f_2 dp + \dots + \int_X f_n dp$$

Spunem că variabilele aleatoare f_1, f_2 sunt independente, la fel ca în cazul variabilelor aleatoare simple, dacă oricare ar fi intervalele I_1, I_2 avem $p(\{f_1 \in I_1\} \cap \{f_2 \in I_2\}) = p(f_1 \in I_1) \cdot p(f_2 \in I_2)$. Cu aceasta putem enunța următoarea proprietate:

4) Dacă f_1 și f_2 sunt variabile aleatoare independente atunci:

$$\int_X (f_1 f_2) dp = \left(\int_X f_1 dp \right) \cdot \left(\int_X f_2 dp \right)$$

5)

$$f \geq 0 \Rightarrow \int_X f dp \geq 0$$

6)

$$f \geq g \Rightarrow \int_X f dp \geq \int_X g dp$$

7)

$$\left| \int_X f dp \right| \leq \max |f|$$

Proprietatea următoare este mai specială:

8) Dacă $\lim_{n \rightarrow \infty} f_n = f$ a.p.t. și dacă există o variabilă aleatoare $g \geq 0$ astfel ca $\int_X g dp < \infty$ și $|f_n| \leq g$ pentru orice n , atunci

$$\lim_{n \rightarrow \infty} \int_X f_n dp = \int_X f dp$$

Deci dacă un șir de variabile aleatoare este convergent și dominat de o v.a. pozitivă cu medie finită, atunci media limitei este egală cu limita mediilor.

O variabilă aleatoare complexă este o funcție $f : X \rightarrow \mathbb{C}$ $f = u + iv$ astfel ca u și v să fie variabile aleatoare reale. Se definește media lui f prin $\int_X f dp = \int_X u dp + i \int_X v dp$. Se demonstrează că proprietățile 1)-8) se păstrează.

Momentul de ordin k al unei variabile aleatoare oarecare se definește ca pentru variabilele simple

$$M_k(f) = M(f^k) = \int_X f^k dp$$

Momentul centrat de ordinul k se definește analog:

$$\mu_k(f) = M\left((f - M(f))^k\right) = \int_X (f - M(f))^k dp$$

Dispersia se definește de asemenea ca pentru variabile simple:

$$D(f) = M((f - M(f))^2) = \int_X (f - M(f))^2 dp = M_2(f) - (M(f))^2$$

Proprietățile demonstrate pentru dispersie în cazul variabilelor aleatoare simple rămân adevărate în general, ele deducându-se din proprietățile mediei, care rămân adevărate în general.

Funcția caracteristică a unei variabile aleatoare generale se definește ca pentru variabile aleatoare simple, cu ajutorul mediei:

$$f_c(t) = M(e^{\sqrt{-1}tf}) = \int_X e^{itf} dp$$

Si în acest caz proprietățile demonstrate pentru variabile aleatoare simple rămân adevărate în general.

Unei variabile aleatoare generale nu-i putem atașa o diagramă ca în cazul variabilelor aleatoare simple, ele luînd în general o infinitate de valori. În cazul particular când valorile se pot enumera $v_1, v_2, \dots, v_n, \dots$ iar evenimentele $E_i = \{f = v_i\}$ au probabilitățile p_i putem asocia lui f diagrama:

$$\begin{pmatrix} v_1 & v_2 & v_3 \dots & v_n \dots \\ p_1 & p_2 & p_3 \dots & p_n \dots \end{pmatrix}$$

În acest caz variabilele simple

$$f_n(\omega) = \begin{cases} v_i & \omega \in E_i, i \leq n \\ 0 & \omega \in E_i, i > n \end{cases}$$

tind spre f , și deci $\int_X f dp = \lim_{n \rightarrow \infty} \int_X f_n dp = \lim_{n \rightarrow \infty} (v_1 p_1 + v_2 p_2 + v_3 p_3 + \dots v_n p_n) = \sum_{i=1}^{\infty} v_i p_i = M(f)$. Analog, dispersia este

$$D(f) = \sum_{i=1}^{\infty} (v_i - M(f))^2 p_i$$

Exemplul 3.10 Fie o variabilă aleatoare ce ia valorile $0, 1, 2, \dots, n, \dots$ cu probabilitățile $p_0, p_1, p_2, \dots, p_n, \dots$. Fie $G_f(t) = \sum_{i=0}^{\infty} p_i t^i$. Să se arate că:

a) $M(f) = G'_f(1)$;

b) $M_2(f) = G''_f(1) + G'_f(1)$;

c) $D(f) = G''_f(1) + G'_f(1) - (G'_f(1))^2$. Funcția G_f se numește funcția generatoare a v.a. f .

Soluție. a) În acest caz avem $v_i = i$. $G'_f(t) = \sum i p_i t^{i-1}$ și deci $G'_f(1) = \sum i p_i = \sum v_i p_i = M(f)$. b) $G''_f(t) = \sum i(i-1) p_i t^{i-2}$ deci $G''_f(1) = \sum i^2 p_i - \sum i p_i$ adică $G''_f(1) = M_2(f) - M(f)$ de unde rezultă b). c) Avem $D(f) = M_2(f) - (M(f))^2$ ceea ce dă, conform cu a) și b), exact c).

3.3 Funcția de repartiție densitatea de probabilitate

Fie $f : X \rightarrow R$ o variabilă aleatoare definită pe spațiul probabilizat X . Asociem lui f o funcție $F : R \rightarrow R$ prin formula $F(t) = p(f < t)$. Funcția F nu mai apare ca depinzând explicit de X . Ea conține informațiile "probabilistice" despre f , independent de natura elementelor din X . Este posibil ca aceeași funcție F să corespundă la variabile aleatoare diferite, definite pe același spațiu X sau pe spații diferite.

Definiția 3.11 Funcția $F : R \rightarrow R$ cu $F(t) = p(f < t)$, f fiind o variabilă aleatoare, se numește funcția de repartiție a acestei variabile aleatoare.

Definiția 3.12 Fie o variabilă aleatoare f cu funcția de repartiție F . O funcție $\rho : R \rightarrow [0, \infty)$, integrabilă, cu proprietatea că $F(t) = \int_{-\infty}^t \rho(x) dx$ se numește densitate de repartiție a variabilei f .

Teorema 3.13 Funcția de repartiție are următoarele proprietăți:

1) F este monoton crescătoare

$$2) F(-\infty) = \lim_{t \rightarrow -\infty} F(t) = 0 \quad F(\infty) = \lim_{t \rightarrow \infty} F(t) = 1$$

3) F este continuă la stînga.

$$4) p(a \leq f < b) = F(b) - F(a).$$

5) Reciproc, dacă o funcție $F : R \rightarrow R$ are proprietățile de mai sus, există un câmp de probabilitate (X, Ω, p) și o variabilă aleatoare pe X , care are pe F ca funcție de repartiție.

Demonstrație. 1) $t_1 < t_2 \Rightarrow \{f < t_1\} \subseteq \{f < t_2\} \Rightarrow p(f < t_1) \leq p(f < t_2)$ adică $F(t_1) \leq F(t_2)$. 3) Fie $t_n \rightarrow t$ monoton crescător. În aceste condiții mulțimile $\{f < t_n\}$ formează un șir crescător spre $\{f < t\}$ deci $p(f < t) = \lim_{n \rightarrow \infty} p(f < t_n)$, adică F este continuă

în t la stînga. 2) Șirul de mulțimi $A_n = \{f < t_n\}$ este monoton descrescător cu intersecția vidă pentru orice șir (t_n) monoton descrescător spre $-\infty$. Din continuitatea probabilității, rezultă $F(-\infty) = \lim_{n \rightarrow \infty} F(t_n) = \lim_{n \rightarrow \infty} p(A_n) = p(\emptyset) = 0$. Analog se demonstrează că $F(\infty) = 1$.

4) Avem $\{f < b\} = \{f < a\} \cup \{a \leq f < b\}$, reuniune disjunctă. Luând probabilitățile în ambii membri găsim formula din punctul 4).

5). Se poate lua $X = R$, $\Omega =$ cea mai mică σ -algebră ce conține intervalele de forma $[a, b)$, iar probabilitatea este definită prin $p([a, b)) = F(b) - F(a)$. Detaliile nu fac obiectul prezentului curs.

Teorema 3.14 Densitatea de probabilitate a unei variabile aleatoare are proprietățile:

$$1) \rho(t) \geq 0$$

$$2) \int_{-\infty}^{\infty} \rho(t) dt = 1$$

$$3) p(a \leq f < b) = \int_a^b \rho(x) dx.$$

4) Reciproc, dacă o funcție integrabilă $\rho : R \rightarrow R$ are proprietățile 1)-3) atunci există un câmp de probabilitate (X, Ω, p) și o variabilă aleatoare pe X , care admite pe ρ ca densitate de probabilitate.

Demonstrație. 1) și 2) sunt evidente iar 3) rezultă din punctul 4) al teoremei precedente. 5) Se poate lua, ca în teorema precedentă, $X = R$, $\Omega =$ cea mai mică σ -algebră ce conține intervalele de forma $[a, b)$, iar probabilitatea este definită prin $p([a, b)) = F(b) - F(a)$, unde $F(x) = \int_{-\infty}^x \rho(t) dt$.

QED.

Funcția de repartiție a unei variabile aleatoare există întotdeauna pe când densitatea de repartiție nu. Dacă există densitate de repartiție ρ , atunci F este continuă iar în punctele de continuitate ale lui ρ funcția F este în plus derivabilă și există relația $F'(t) = \rho(t) = \lim_{\Delta t \rightarrow 0} \frac{p(t \leq f < t + \Delta t)}{\Delta t}$.

Pentru a vedea cum putem calcula efectiv media, dispersia, etc. în general, studiem pe scurt un concept numit integrala Stieltjes.

3.4 Integrala Stieltjes

Fie $F : [a, b] \rightarrow R$ o funcție monoton crescătoare. Vrem să definim $\int_a^b f(x)dF$ pentru o funcție $f : [a, b] \rightarrow R$. Fie

$$\Delta : x_0 < x_1 < \dots < x_{n-1} < x_n = b$$

o diviziune a intervalului $[a, b]$ și fie câte un element $\xi_i \in [x_{i-1}, x_i]$. Definim suma Stieltjes analog cu suma Riemann:

$$S_\Delta = \sum_{i=1}^n f(\xi_i) (F(x_i) - F(x_{i-1}))$$

Spunem că f este integrabilă Stieltjes în raport cu F dacă pentru orice șir de diviziuni $(\Delta_n)_{n \in \mathbb{N}}$ cu norma $|\Delta_n| = \max_i |x_i - x_{i-1}|$ tinzând la zero, șirul de sume Stieltjes are o aceeași limită, independent de șirul de diviziuni sau de punctele ξ_i alese. Această limită se notează cu $\int_a^b f(x)dF(x)$ sau $\int_a^b f dF$. Se demonstrează că pentru orice funcție f continuă integrala Stieltjes există. Următoarele proprietăți ale integralei Stieltjes seamănă cu cele ale integralei Riemann:

a) Dacă f_1 și f_2 sunt integrabile Stieltjes atunci $\alpha f_1 + \beta f_2$ este integrabilă și

$$\int_a^b (\alpha f_1 + \beta f_2) dF = \alpha \int_a^b f_1 dF + \beta \int_a^b f_2 dF.$$

b) Dacă $a < c < b$ atunci $\int_a^b f dF = \int_a^c f dF + \int_c^b f dF$.

c) $f \geq 0$ implică $\int_a^b f dF \geq 0$.

d) $\int_a^b f dF \leq \max_{x \in [a, b]} |f(x)| \cdot (F(b) - F(a))$.

e) Dacă $F = F_1 + F_2$ sumă de două funcții monotone, atunci $\int_a^b f dF = \int_a^b f dF_1 + \int_a^b f dF_2$.

Demonstrația acestor proprietăți nu o facem aici. Să calculăm câteva integrale Stieltjes.

Exemplul 3.15 Fie $F : [a, b] \rightarrow R$, $F(x) = x$. Atunci $F(x_i) - F(x_{i-1}) = x_i - x_{i-1}$ și integrala Stieltjes devine integrala Riemann.

Exemplul 3.16 $F : [a, b] \rightarrow R$ este monoton crescătoare, derivabilă, cu derivata $\rho(x) = F'(x)$ integrabilă Riemann. Atunci $F(x_i) - F(x_{i-1}) = \rho(\xi'_i)(x_i - x_{i-1})$ (conform teoremei lui Lagrange) și dacă luăm în suma Stieltjes $\xi_i = \xi'_i$ obținem suma Riemann pentru funcția $f\rho$. Deci

$$\int_a^b f dF = \int_a^b f(x)\rho(x)dx \quad (3.1)$$

Exemplul 3.17 Fie $c \in [a, b]$ și $F : [a, b] \rightarrow R$ dată prin

$$F(x) = \begin{cases} \alpha & x \in [a, c] \\ \beta & x \in (c, b] \end{cases}$$

Atunci în suma Stieltjes numai termenul $f(\xi_i)(x_i - x_{i-1})$ pentru care $x_{i-1} \leq c \leq x_i$ este nenul și suma este $(\beta - \alpha)f(\xi_i)$ care tinde spre $(\beta - \alpha)f(c)$ dacă f este continuă în c . Deci

$$\int_a^b f dF = (\beta - \alpha)f(c) = (F(c+0) - F(c-0)) \cdot f(c) \quad (3.2)$$

Exemplul 3.18 Fie $F : [a, b] \rightarrow R$ o funcție constantă pe porțiuni, cu salturi în punctele c_1, c_2, \dots, c_n unde avem $F(c_i + 0) - F(c_i - 0) > 0$. Atunci, la fel ca în exemplul anterior, dacă f este continuă în punctele c_i avem

$$\int_a^b f dF = \sum_{i=1}^n f(c_i) (F(c_i + 0) - F(c_i - 0))$$

Exemplul 3.19 Fie F continuă pe porțiuni, cu discontinuități în c_1, c_2, \dots, c_n iar între aceste puncte derivabilă, $F'(x) = \rho(x)$, integrabilă. Atunci în $\int_a^b f dF$ avem două părți: una de tipul (3.1) iar alta de tipul (3.2), deci

$$\int_a^b f dF = \int_a^b f(x) \rho(x) dx + \sum_{i=1}^n f(c_i) (F(c_i + 0) - F(c_i - 0))$$

Se poate defini $\int_{-\infty}^{+\infty} f dF$ prin $\lim_{a \rightarrow -\infty, b \rightarrow \infty} \int_a^b f dF$. Proprietățile a)..e) se păstrează la fel ca și metodele de calcul descrise în exemplele anterioare.

3.5 Media și funcția de repartiție

Fie o variabilă aleatoare f cu valori într-un interval $[a, b)$ și cu funcția de repartiție F . Fie

$$\Delta : x_0 = a < x_1 < \dots < x_{n-1} < x_n = b$$

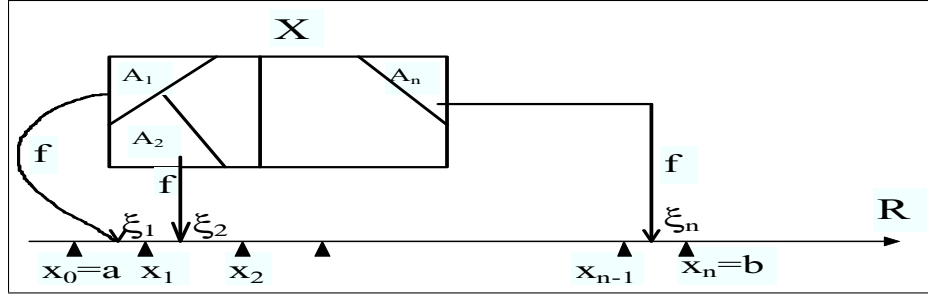
o diviziune a intervalului $[a, b]$. Luăm cîte un punct $\xi_i \in [x_{i-1}, x_i]$. Fie variabila aleatoare

$$f_{\Delta}(\omega) = \begin{cases} 0 & \text{dacă } f(\omega) \notin [a, b) \\ \xi_i & \text{dacă } f(\omega) \in [x_{i-1}, x_i) \end{cases}$$

Altfel spus, dac  $A_i = \{\omega \in X \mid f(\omega) \in [x_{i-1}, x_i)\}$, atunci

$$f_{\Delta}(\omega) = \begin{cases} \xi_1, & \text{dacă } \omega \in A_1 \\ \xi_2, & \text{dacă } \omega \in A_2 \\ \\ \xi_n, & \text{dacă } \omega \in A_n \end{cases}$$

Variabila f_Δ este deci simplă, cu media $M(f_\Delta) = \sum_i \xi_i \cdot p(x_{i-1} \leq f < x_i) = \sum_i \xi_i (F(x_i) - F(x_{i-1}))$. Atunci când norma diviziunii $|\Delta| \rightarrow 0$, variabilele f_Δ tind la f (a se vedea figura)



Exemplu de diviziune pentru calculul mediei

deci prin definiție

$$M(f) = \lim_{|\Delta| \rightarrow 0} M(f_\Delta) = \lim_{|\Delta| \rightarrow 0} \sum_i \xi_i (F(x_i) - F(x_{i-1})) = \int_a^b x dF(x) \quad (3.3)$$

Am obținut astfel formula simplă de mai sus care exprimă media unei variabile aleatoare printr-o integrală Stieltjes. Dacă variabila aleatoare nu este mărginită de a și b finite, atunci integrala de mai sus trebuie luată între $-\infty$ și ∞ . Variabila aleatoare f^k este aproximată de variabilele aleatoare f_Δ^k care au mediile

$$\sum_i \xi_i^k (F(x_i) - F(x_{i-1})) \rightarrow \int_a^b x^k dF(x)$$

Dispersia lui f este media variabilei $(f - M(f))^2$. Această variabilă aleatoare este aproximată prin variabilele simple $(f_\Delta - M(f))^2$ care au mediile

$$\sum_i (\xi_i - M(f))^2 (F(x_i) - F(x_{i-1})) \rightarrow \int_a^b (x - M(f))^2 dF(x)$$

Funcția caracteristică a lui f este media variabilei aleatoare e^{itf} care este aproximată prin variabilele $e^{\sqrt{-1}tf_\Delta}$ cu mediile

$$\sum_i e^{\sqrt{-1}t\xi_i} (F(x_i) - F(x_{i-1})) \rightarrow \int_a^b e^{\sqrt{-1}tx} dF(x)$$

Dacă f nu este mărginită atunci integralele trebuie luate de la $-\infty$ la ∞ . Avem deci

formulele ($\rho = F'$ este densitatea de probabilitate):

$$\begin{aligned}
 M(f) &= \int_X f dp = \int_{-\infty}^{\infty} x dF(x) = \int_{-\infty}^{\infty} x \rho(x) dx \\
 M_k(f) &= M(f^k) = \int_X f^k dp = \int_{-\infty}^{\infty} x^k dF(x) = \int_{-\infty}^{\infty} x^k \rho(x) dx \\
 \mu_k(f) &= M\left((f - M(f))^k\right) = \int_X (f - M(f))^k dp = \\
 &= \int_{-\infty}^{\infty} (x - M(f))^k dF(x) = \int_{-\infty}^{\infty} (x - M(f))^k \rho(x) dx \\
 D(f) &= M\left((f - M(f))^2\right) = \int_X (f - M(f))^2 dp = \\
 &= \int_{-\infty}^{\infty} (x - M(f))^2 dF(x) = \int_{-\infty}^{\infty} (x - M(f))^2 \rho(x) dx \\
 f_c(t) &= M(e^{\sqrt{-1}tf}) = \int_X e^{\sqrt{-1}tf} dp = \int_{-\infty}^{\infty} e^{\sqrt{-1}tx} dF(x) = \int_{-\infty}^{\infty} e^{\sqrt{-1}tx} \rho(x) dx
 \end{aligned} \tag{3.4}$$

Prin urmare cunoașterea funcției de repartiție este suficientă pentru determinarea caracteristicilor numerice ale unei variabile aleatoare. Așa cum s-a menționat mai înainte proprietățile acestor mărimi, demonstrate pentru variabile aleatoare simple, rămân adevărate în general. Se vede în ultima formulă că funcția caracteristică este până la un factor transformata Fourier inversă a densității de probabilitate, deci, densitatea de probabilitate este determinată de funcția caracteristică (prin transformata Fourier). Mai precis se demonstrează teorema următoare:

Teorema 3.20 *Fie două variabile aleatoare f și g , cu funcțiile de repartiție F, G și funcțiile caracteristice f_c și g_c . Dacă funcțiile caracteristice sunt egale, atunci $F = G$.*

Nu demonstrăm această teoremă dar o vom folosi mai departe pentru a pune în evidență egalitatea unor funcții de repartiție.

Exemplul 3.21 *Fie $X = [0, 1] \times [0, 1]$, $\Omega =$ mulțimea figurilor cu arie, $p(A) =$ aria lui A . Fie $f : X \rightarrow R$ $f(x, y) = x + y$. Este clar că f este o variabilă aleatoare pentru că $\{f < t\} = \{(x, y) \in X \mid x + y < t\}$ este o figură cu arie. Se cere:*

- 1) *Să se determine un șir mărginit de v.a. simple ce tind la f .*
- 2) *Să se determine $\int_X f dp$ utilizând șirul de v.a. de la punctul precedent.*
- 3) *Să se determine funcția de repartiție și densitatea de repartiție (dacă există).*
- 4) *Să se determine media și dispersia folosind funcția de repartiție.*

Soluție. 1,2) Putem în mai multe moduri determina șiruri de v.a. simple ce tind la f . De exemplu pentru un n dat fie

$$\begin{aligned}
 0 &= x_0 < x_1 < \dots < x_n = 1 & x_i &= i/n \\
 0 &= y_0 < y_1 < \dots < y_n = 1 & y_j &= j/n
 \end{aligned}$$

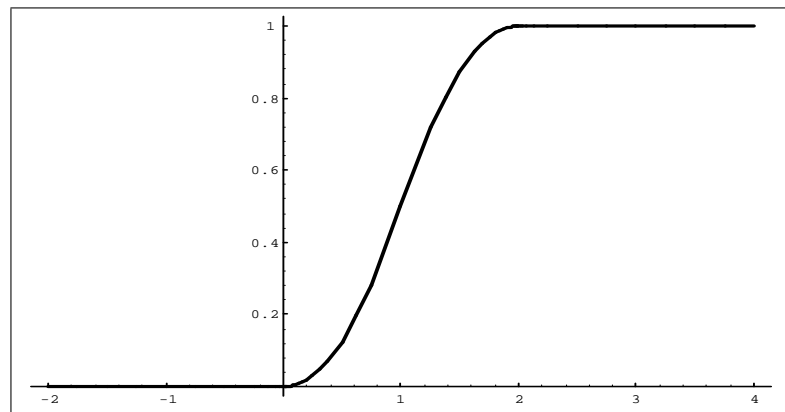
și fie $f_n(x, y) = (x_{i-1} + x_i) / 2 + (y_{j-1} + y_j) / 2$ dacă $(x, y) \in [x_{i-1}, x_i] \times [y_{j-1}, y_j]$, adică valoarea lui $f(x, y)$ în centrul P_{ij} dreptunghiului $[x_{i-1}, x_i] \times [y_{j-1}, y_j]$ pe care o notăm și $f(P_{ij})$. Dacă (x, y) nu se află într-un dreptunghi $[x_{i-1}, x_i] \times [y_{j-1}, y_j]$ atunci definim $f_n(x, y) = 0$. Este clar că $f_n \rightarrow f$ a.p.t. și avem:

$$\begin{aligned} \int_X f dp &= \lim_{n \rightarrow \infty} \int_X f_n dp = \lim_{n \rightarrow \infty} \sum_{i,j=1..n} f(P_{ij}) p([x_{i-1}, x_i] \times [y_{j-1}, y_j]) = \\ &= \lim_{n \rightarrow \infty} \sum_{i,j=1..n} f(P_{ij}) (x_i - x_{i-1}) (y_j - y_{j-1}) = \int \int_X f(x, y) dx dy = \\ &= \int_0^1 \int_0^1 (x + y) dx dy = 1 \end{aligned}$$

3) Funcția de repartiție este

$$\begin{aligned} F(t) &= \text{aria}(\{(x, y) \in X | x + y < t\}) = \\ &= \begin{cases} 0 & t \leq 0 \\ \frac{t^2}{2} & 0 \leq t \leq 1 \\ 1 - \frac{(2-t)^2}{2} & 1 \leq t \leq 2 \\ 1 & t \geq 2 \end{cases} \end{aligned}$$

Graficul funcției de repartiție este:

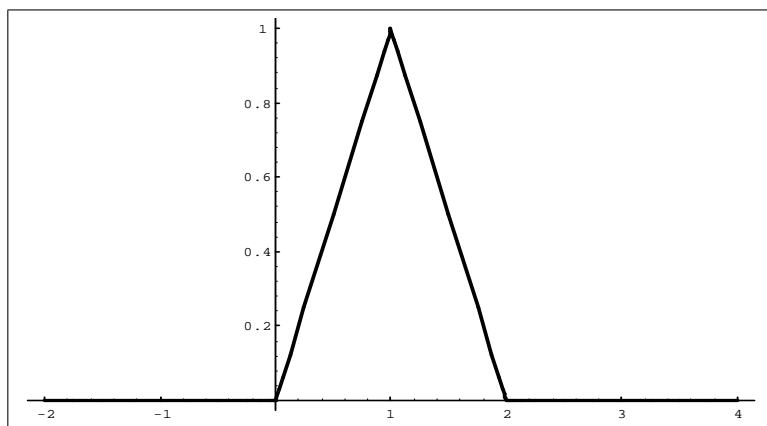


Graficul funcției de repartiție al variabilei aleatoare f

Deoarece F este derivabilă avem densitatea de probabilitate:

$$\rho(t) = \begin{cases} 0 & t \leq 0 \\ 2t & 0 \leq t \leq 1 \\ 2(2-t) & 1 \leq t \leq 2 \\ 0 & t \geq 2 \end{cases}$$

Graficul densității este:

Graficul densității de probabilitate a variabilei aleatoare f

4)

$$\begin{aligned}
 M(f) &= \int_{-\infty}^{\infty} t dF(t) = \int_{-\infty}^{\infty} t \rho(t) dt = \int_0^2 t \rho(t) dt = \\
 &= \int_0^1 2t^2 dt + \int_1^2 t \cdot 2(2-t) dt = 1
 \end{aligned}$$

Analog se poate calcula dispersia.

3.6 Rezumat

Pentru a putea lucra și cu variabile aleatoare cu un număr infinit de valori trebuie lărgit conceptul de algebră de evenimente la cel de σ algebră iar conceptul de probabilitate la cel de σ probabilitate. Concret:

i) O σ algebră de evenimente (mulțimi) este o algebră de evenimente (mulțimi) cu proprietatea că dacă A_n sunt în algebră pentru orice n , atunci $\bigcup_{n=1.. \infty} A_n$ este în algebră.

ii) O probabilitate este σ probabilitate dacă pentru orice șir de evenimente (mulțimi) disjuncte A_n avem $p(\bigcup_{n=1.. \infty} A_n) = \sum_{n=1}^{\infty} p(A_n)$.

iii) σ probabilitatea are o remarcabilă proprietate de continuitate: dacă A este intersecția unui șir monoton descrescător sau reuniunea unui șir monoton crescător de evenimente A_n atunci $p(A) = \lim_{n \rightarrow \infty} p(A_n)$.

iv) O variabilă aleatoare generală este o funcție reală care întoarce intervalele în evenimente, deci au sens expresii ca $p(f = a)$ sau $p(a \leq f < b)$ etc. Operațiile uzuale cu funcții, inclusiv trecerea la limită, aplicate variabilelor aleatoare duc tot la variabile aleatoare. Orice variabilă aleatoare este limită de variabile aleatoare simple (ce iau un număr finit de valori).

v) Caracteristicile numerice ale unei variabile aleatoare se definesc prin limita caracteristicilor corespunzătoare ale variabilelor aleatoare simple ce tind la variabila generală. Această

limită nu există întotdeauna. Proprietățile acestor caracteristici: medie, dispersie, momente, funcție caracteristică, demonstrate anterior pentru variabile aleatoare simple se păstrează prin trecere la limită și pentru variabile aleatoare generale.

vi) Pentru calculul concret al caracteristicilor numerice se introduce funcția de repartiție definită prin $F(t) = p(f < t)$, F fiind funcția de repartiție a variabilei aleatoare f . Toate informațiile numerice legate de f sunt conținute în funcția de repartiție F , independent de spațiul probabilitat pe care e definită f .

vii) Densitatea de probabilitate ρ , este derivata funcției de repartiție (dacă există), adică $\rho(x) = \lim_{\Delta t \rightarrow 0} \frac{p(t \leq f < t + \Delta t)}{\Delta t}$, sau, mai exact se definește prin: $F(x) = \int_{-\infty}^x \rho(t) dt$. Proprietățile funcției de repartiție și ale densității de probabilitate sunt studiate mai sus.

viii) Pentru a ajunge la scopul propus, exprimarea caracteristicilor numerice ale unei variabile aleatoare prin intermediul funcției de repartiție, avem nevoie de un concept auxiliar numit integrala Stieltjes

$$\int_a^b f(x) dg(x) = \lim_{|\Delta| \rightarrow \infty} \sum f(\xi_i) (g(x_i) - g(x_{i-1}))$$

, o generalizare a integralei Riemann. Proprietățile integralei Stieltjes seamănă în mare măsură cu cele ale integralei Riemann.

ix) În final, prin formulele (3.4), media, dispersia, etc. sunt exprimate direct prin funcția de repartiție ori prin densitatea de probabilitate, independent de spațiul pe care e definită variabila aleatoare.

FORMULE UTILIZATE FRECVENT:

a) Pentru o probabilitate σ -aditivă: Dacă $A_1 \supseteq A_2 \supseteq \dots \supseteq A_n \supseteq \dots$ atunci $p(\bigcap_{n=1, \infty} A_n) = \lim_{n \rightarrow \infty} p(A_n)$ iar dacă $A_1 \subseteq A_2 \subseteq \dots \subseteq A_n \subseteq \dots$ atunci $p(\bigcup_{n=1, \infty} A_n) = \lim_{n \rightarrow \infty} p(A_n)$.

b) $\int_a^b f dF = \int_a^b f(x) \rho(x) dx + \sum_{i=1}^n f(c_i) (F(c_i + 0) - F(c_i - 0))$ unde F e funcția de repartiție, ρ e densitatea de probabilitate iar c_i sunt punctele de discontinuitate ale lui F

c) Formulele (3.4) pentru medie, momente, etc.

d) $p(a \leq f < b) = F(b) - F(a) = \int_a^b \rho(x) dx$.

3.7 Exerciții

1. Fie o variabilă aleatoare :

$$f = \begin{pmatrix} 0 & 1 & 2 & \dots & n & \dots \\ \alpha p^0 & \alpha p & \alpha p^2 & \dots & \alpha p^n & \dots \end{pmatrix}$$

a) Să se determine α astfel ca diagrama de mai sus să corespundă unei variabile aleatoare.

b) Să se calculeze media, dispersia și funcția caracteristică.

Indicație. Trebuie ca $\sum_{k=0}^{\infty} \alpha p^k = 1$. De aici rezultă α . O problemă analoagă este rezolvată în cadrul lecției.

2. Fie variabila aleatoare uniformă

$$\begin{pmatrix} 1 & 2 & \dots & k & \dots & n \\ \frac{1}{n} & \frac{1}{n} & \dots & \frac{1}{n} & \dots & \frac{1}{n} \end{pmatrix}$$

Să se determine funcția caracteristică, media și dispersia.

Indicație. Se utilizează definițiile.

3. O urnă conține s bile din care $s-1$ negre și una albă. Se fac extrageri din urnă cu punerea bilei înapoi pînă se extrage bila albă. Fie variabila aleatoare f care ia valoarea k dacă bila albă apare la extragerea k . Se cere valoarea medie a lui X .

Indicație. Probabilitatea de a extrage o bilă albă este $p=1/s$ iar probabilitatea de a extrage o bilă neagră este $q=(s-1)/s$. Probabilitatea ca bila albă să apară abia la extragerea k este $p_k = \left(\frac{s-1}{s}\right)^{k-1} \cdot \frac{1}{s}$. Mai departe problema seamănă cu 1).

4. Un jucător trage la o țintă; probabilitatea de a o nimeri este $p \in (0, 1)$ iar probabilitatea de a rata lovitură este $q = 1 - p$. Partida se termină când fie numărul lovirilor este cu 2 mai mare ca al ratărilor, caz în care câștigă partida, fie când numărul ratărilor este cu 2 mai mare ca numărul lovirilor, caz în care pierde partida. Fie A_n evenimentul care constă în a câștiga la tragerea n iar B_n evenimentul care constă în a pierde la tragerea n .

a) $p(A_{2n+1}) = ?$ $p(B_{2n+1}) = ?$.

b) Care e probabilitatea ca jucătorul să câștige?

c) Care e probabilitatea ca jucătorul să piardă?

d) Care e probabilitatea ca partida să nu se termine niciodată?

e) Fie X variabila aleatoare care ia valoarea k dacă partida se termină la tragerea k . Se ce probabilitățile $p(X=k)$ și valoarea medie a lui X .

Indicație. Dacă jucătorul câștigă (sau pierde) la mutarea k , atunci nimerește (sau ratează) cu două lovituri în plus față de situația contrară. Prin urmare numărul lovirilor este par, deci $p(A_{2n+1}) = p(B_{2n+1}) = 0$. Fie p_{2l} probabilitatea ca partida să fie câștigată la tragerea $2l$, fie q_{2l} probabilitatea ca partida să fie pierdută la tragere $2l$ și fie r_{2l} probabilitatea de "remiză" la tragerea $2l$. Avem $p_{2l} = r_{2l-2} \cdot p^2$ $q_{2l} = r_{2l-2} \cdot q^2$ $r_{2l} = r_{2l-2} \cdot (1 - p^2 - q^2)$.

5. a) Pentru care λ funcția

$$\rho(x) = \begin{cases} \lambda(-x^2 + 1) & x \in [-1, 1] \\ 0 & x \notin [-1, 1] \end{cases}$$

este o densitate de probabilitate?

b) Se cere media și dispersia variabilei cu densitatea ρ .

c) Care este probabilitatea ca o v.a. cu densitatea ρ să ia valori între 0,2 și 0,5?

Indicație. λ rezultă din $\int_{-\infty}^{\infty} \rho(x) dx = 1$. În rest se folosesc formulele (3.4).

6. a) O variabilă aleatoare f are densitatea

$$\rho(x) = \begin{cases} 0, & x < 0 \\ e^{-x}, & x \geq 0 \end{cases}$$

Care este funcția de repartiție și densitatea de probabilitate a variabilelor f^2 , $2f + 1$, e^f .

b) Să se arate că dacă o v.a. f are densitatea $\rho(x)$ iar $g: R \rightarrow R$ este bijectie derivabilă cu $g' > 0$, atunci variabila aleatoare $g \circ f$ are densitatea $\rho(g^{-1}(t)) \frac{1}{g'(t)}$.

Indicație. Fie G funcția de repartiție a variabilei f^2 . Avem

$$\begin{aligned} G(t) &= p(f^2 < t) = p(-\sqrt{t} < f < \sqrt{t}) = \int_{-\sqrt{t}}^{\sqrt{t}} \rho(x) dx = \\ &= \int_0^{\sqrt{t}} e^{-x} dx = 1 - e^{-\sqrt{t}} \quad \text{pentru } t \geq 0 \end{aligned}$$

În mod evident $G(t)=0$ pentru $t \leq 0$. Densitatea lui f^2 se obține prin derivarea lui G . Analog se procedează și cu celelalte expresii de f .

7. Se dă funcția

$$f(x) = \begin{cases} 0 & x < 0 \\ x-1 & 0 \leq x \leq 1 \\ x^2 & 1 < x < 2 \\ 5-x & x \geq 2 \end{cases}$$

Să se calculeze $\int_{-1}^5 (2+x) df(x)$.

Indicație. Se aplică tehnicile de calcul ale integralei Stieltjes din lecție.

8. Funcția de repartiție a unei v.a. continue este

$$F(x) = \begin{cases} a, & x \leq 0 \\ bx^2, & x \in (0, 1) \\ c, & x \geq 1 \end{cases}$$

Se cer a, b, c , și $P(1/4 \leq f \leq 1/2)$.

9. Să admitem că timpul de așteptare într-o stație de metrou este o v.a. cu funcția de repartiție:

$$F(t) = \begin{cases} 0, & t \leq 0 \\ t/4, & 0 < t \leq 2 \\ 1/2, & 2 < t \leq 4 \\ t/8, & 4 < t \leq 8 \\ 1, & t > 8 \end{cases}$$

- a) Se cere densitatea de probabilitate.
 b) Care este probabilitatea ca un călător să aștepte mai mult de 3 minute?
 c) Care este probabilitatea ca un călător să aștepte mai mult de 3 minute știind că a așteptat mai mult de 1 minut?
 d) Care este probabilitatea ca un călător să aștepte mai puțin de 3 minute știind că a așteptat mai mult de 1 minut?

10. Funcția de repartiție a unei v.a. X este

$$F_X(t) = \begin{cases} 0, & t < 0 \\ x^2, & t \in (0, 1] \\ 1, & t > 1 \end{cases}$$

Se cer $M(X)$, $D(X)$, $\phi_X(t)$, $F_{3X+2}(t)$, $F_{X^2}(t)$.

11. Fie o v.a. cu funcția de repartiție F care admite medie și dispersie finite. Să se arate că

$$\int_{-\infty}^{\infty} (x-a)^2 dF(x)$$

este minimă când $a=M(f)$.

Indicație. Fie $i(a) = \int_{-\infty}^{\infty} (x-a)^2 dF(x)$. La extrem trebuie să avem $i'(a) = 0$.

12. O variabilă aleatoare f ia valori în intervalul $[a, b]$. Să se arate că $M(f) \in [a, b]$ și $D(f) \in [0, (b-a)^2/4]$.

Indicație. Se exprimă media și dispersia prin integrale Stieltjes apoi se majorează corespunzător mărimile de sub integrală. Se poate utiliza problema precedentă.

13. Intr-un sistem de axe rectangulare xOy se dau $A(2,0)$, $B(0,1)$, $O(0,0)$. Se iau la întâmplare $A' \in [O, A]$ și $B' \in [O, B]$. A' și B' au repartiții uniforme. Se cer:

- a) Valoarea medie a ariei triunghiului $OA'B'$.
 b) Valoarea medie a perimetrului triunghiului $OA'B'$.
 c) Valoarea medie a lungimii segmentului $A'B'$.

Indicație. Se procedează ca în exemplul 7).

14. Se iau la întâmplare două puncte în intervalul $[0, 1]$. Care este funcția de repartiție a distanței dintre ele? Care e valoarea medie a acestei distanțe?

Indicație. Se procedează ca la problema 7).

15. O variabilă aleatoare f are densitatea $p(x) = \frac{1}{\pi} \frac{1}{1+x^2}$. Se cere repartiția variabilei $g = \min(|f|, 1)$ precum și media lui g .

Indicație. Fie G funcția de repartiție a lui g . Dacă $t \leq -1$ atunci

$$G(t) = p(g < t) = p(\min(|f|, 1) < t) = p(\emptyset) = 0$$

În mod analog $G(t) = 1$ pentru $t > 1$. Dacă $t \in (0, 1)$ atunci

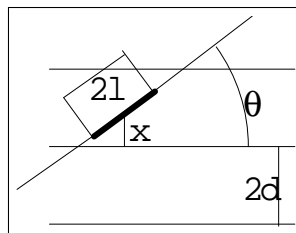
$$\begin{aligned} G(t) &= p(g < t) = p(\min(|f|, 1) < t) = p(|f| < t) \\ &= p(-t < f < t) = \int_{-t}^t \frac{1}{\pi} \frac{1}{1+x^2} dx = \frac{2}{\pi} \arctan(t) \end{aligned}$$

Asemănător se poate calcula $G(1)$ și apoi media cu formulele **3.4**

16. Doi jucători A și B , având fiecare un capital de a respectiv b lei, joacă până la ruina unuia din ei. La un joc șansele de câștig sunt p pentru jucătorul A și $q=1-p$ pentru jucătorul B . La fiecare joc, cel ce pierde plătește 1 leu câștigătorului. Care sunt șansele de câștig ale fiecăruia?

17. Problema lui Buffon. Pe un plan orizontal sunt trasate drepte paralele între ele la aceeași distanță $2d$. Se aruncă pe acest plan un ac de lungime $2l < 2d$. Care este probabilitatea ca acul să intersecteze o paralelă oarecare?

Indicație. Starea de intersecție a acului cu rețeaua de linii depinde de orientarea acului, θ , și de distanța x a mijlocului acului față de cea mai apropiată linie (vezi figura).



$\theta \in [0, \pi)$ iar $x \in [0, d]$. Pentru intersecție trebuie ca $x \leq d \cos \theta$. Spațiul pozițiilor posibile este $X = [0, \pi) \times [0, d]$, iar mulțimea pozițiilor favorabile intersecției este $A = \{(\theta, x) \in X \mid x \leq d \cos \theta\}$. În cazul când toate pozițiile sunt egal probabile, probabilitatea unui eveniment este proporțională cu aria mulțimii punctelor prin care se reprezintă acel eveniment.

Lecția 4

Legi clasice

Am văzut că valorile caracteristice ale unei variabile aleatoare sunt determinate de funcția de repartiție (sau de densitatea de probabilitate, dacă există). În cazul variabilelor simple, caracteristicile numerice sunt determinate de diagrama asociată lor. În cele ce urmează sunt studiate câteva legi frecvent întâlnite în aplicații. Când discutăm o asemenea lege subînțelegem că există un câmp de probabilitate (X, Ω, p) și o funcție $f : X \rightarrow R$ care are ca funcție de repartiție pe F și ca densitate pe ρ . Pentru calcule nu este nevoie să avem explicit date X, Ω, p și f , ci numai pe F sau ρ .

4.1 Repartiția binomială

Fie variabila aleatoare simplă

$$f = \left(\begin{array}{ccccccc} 0 & 1 & \dots & k & \dots & n \\ C_n^0 p^0 q^n & C_n^1 p^1 q^{n-1} & \dots & C_n^k p^k q^{n-k} & \dots & C_n^n p^n q^0 \end{array} \right)$$

unde $p \in [0, 1]$ $p+q=1$ $n \in N$. O asemenea variabilă aleatoare se numește binomială. Funcția caracteristică este

$$f_c(t) = \sum_{k=0}^n e^{itk} C_n^k p^k q^{n-k} = (pe^{it} + q)^n$$

unde $i = \sqrt{-1} \in C$. De aici obținem:

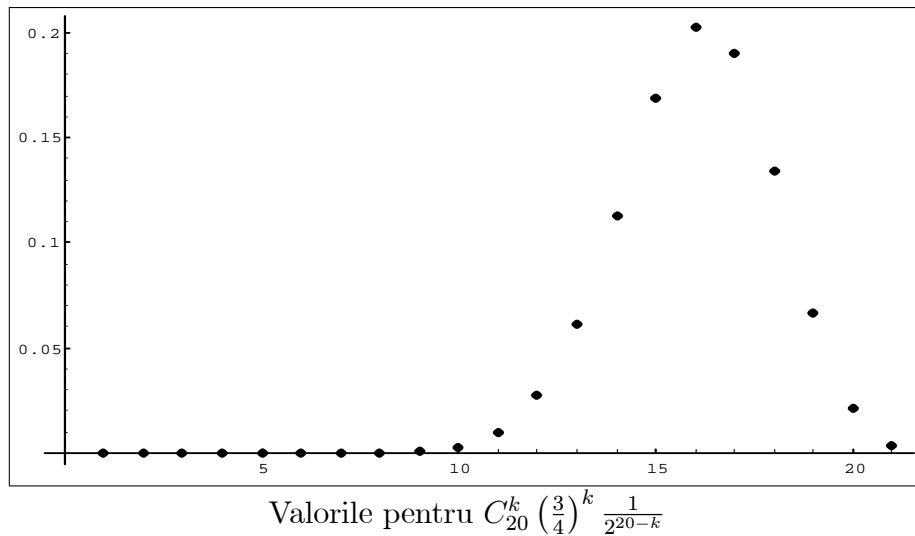
a) $M(f) = \frac{1}{i} f'_c(0) = \frac{1}{i} n (pe^{it} + q)^{n-1} pie \Big|_{t=0} = np$

b) $M_2(f) = \frac{1}{i^2} f''_c(0) = n^2 p^2 + npq$

c) $D = M_2(f) - (M(f))^2 = npq$

Fie p probabilitatea cu care apare A într-o experiență independentă. Atunci în n experiențe independente apariția de k ori a evenimentului A se poate face în C_n^k feluri, egal cu numărul de

variante în care sunt plasate cele k experiențe unde apare A printre toate cele n experiențe. Probabilitatea fiecărui șir de n experiențe, cu k apariții ale evenimentului A , este $p^k q^{n-k}$, $q = 1 - p$. Sumând probabilitățile tuturor acestor șiruri favorabile găsim că probabilitatea ca evenimentul A să apară de k ori este $C_n^k p^k (1-p)^{n-k}$. Această distribuție a fost studiată intensiv de Bernoulli și se mai numește și repartiție Bernoulli. Pentru $n=20$ și $p=3/4$ probabilitățile $C_n^k p^k (1-p)^{n-k}$ arată astfel:



Se vede că probabilitățile care diferă substanțial de 0 apar pentru valori ale lui k în jurul lui np , care în acest caz este 15. În lecția 5 vom demonstra că acest lucru are loc pentru orice $p \in (0, 1)$ și n suficient de mare.

4.2 Repartiția Poisson

Spunem că f este o variabilă aleatoare de tip Poisson dacă ia valori întregi pozitive și $p(f = k) = \frac{\lambda^k e^{-\lambda}}{k!}$ unde $\lambda > 0$. Diagrama unei asemenea variabile aleatoare este

$$f = \begin{pmatrix} 0 & 1 & 2 & \dots & n & \dots \\ e^{-\lambda} & e^{-\lambda} \frac{\lambda}{1!} & e^{-\lambda} \frac{\lambda^2}{2!} & \dots & e^{-\lambda} \frac{\lambda^n}{n!} & \dots \end{pmatrix}$$

λ se numește parametrul variabilei aleatoare. Funcția caracteristică este

$$f_c(t) = \sum_{k=0}^{\infty} e^{itk} e^{-\lambda} \frac{\lambda^k}{k!} = e^{-\lambda} \sum_{k=0}^{\infty} \frac{(e^{it}\lambda)^k}{k!} = e^{-\lambda} e^{\lambda e^{it}} = e^{\lambda(e^{it}-1)}$$

Prin urmare:

a) $M(f) = \frac{1}{i} f'_c(0) = \lambda$

$$b) M_2(f) = \frac{1}{i^2} f_c''(0) = \lambda^2 + \lambda$$

$$c) D = M_2 - M^2 = \lambda$$

Probabilitățile $e^{-\lambda} \frac{\lambda^k}{k!}$ apar în felul următor:

Propoziția 4.1 Fie un șir de variabile aleatoare binomiale

$$f_n = \begin{pmatrix} 0 & 1 & \dots & k & \dots & n \\ q^n & C_n^1 p q^{n-1} & \dots & C_n^k p^k q^{n-k} & \dots & p^n \end{pmatrix}$$

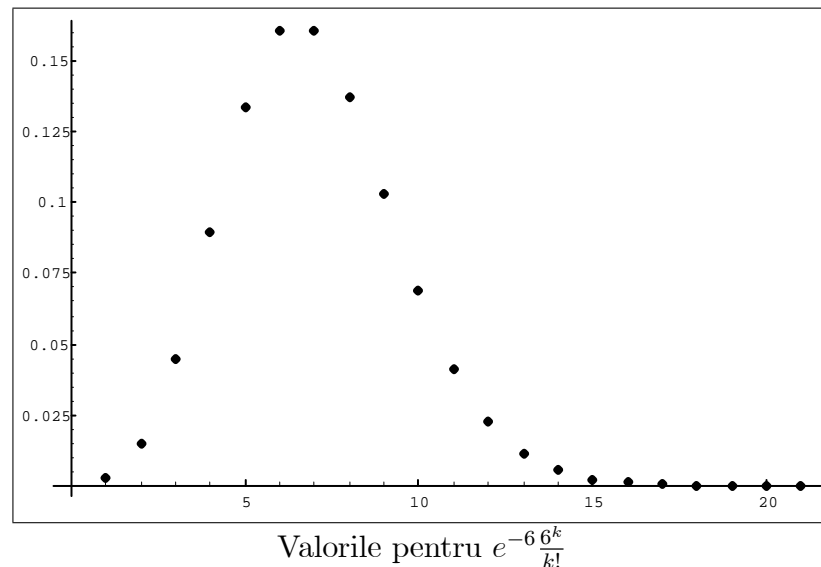
în care $np \rightarrow \lambda$ când $n \rightarrow \infty$. Atunci $p(f_n = k) \rightarrow e^{-\lambda} \frac{\lambda^k}{k!}$ când $n \rightarrow \infty$.

Demonstrație. Vom demonstra această afirmație pentru $np = \lambda$. Avem

$$\begin{aligned} p(f_n = k) &= C_n^k p^k q^{n-k} = \frac{n!}{k!(n-k)!} p^k q^{n-k} = \\ &= \frac{n(n-1)(n-2)\dots(n-k+1)}{k!} p^k q^{n-k} = \frac{n^k \left(1 - \frac{1}{n}\right) \left(1 - \frac{2}{n}\right) \dots \left(1 - \frac{k-1}{n}\right)}{k!} \frac{\lambda^k}{n^k} \left(1 - \frac{\lambda}{n}\right)^{n-k} \rightarrow \\ &\rightarrow \frac{\lambda^k}{k!} e^{-\lambda} \end{aligned}$$

QED.

Prin urmare $\lambda^k e^{-\lambda}/k!$ reprezintă probabilitatea ca un anumit eveniment A să apară de k ori într-un șir foarte mare de n experiențe independente, probabilitatea de apariție a evenimentului A într-o singură experiență fiind foarte mică, de tipul λ/n . Valoarea medie a numărului de apariții este $np = \lambda$. Valorile $p(f = k)$ pentru $k = 0, 1, \dots, 10$, la o valoare a parametrului $\lambda = 6$ se reprezintă astfel:



Se vede că valorile $p(f = k)$ care diferă substanțial de zero se concentrează în jurul valorii $k = \lambda = 6$.

Propoziția 4.2 Fie f și g două variabile Poisson independente, de parametri λ și μ . Atunci $f + g$ este o variabilă Poisson de parametru $\lambda + \mu$.

Demonstrație. Deoarece f și g sunt independente avem

$$\begin{aligned}(f + g)_c(t) &= f_c(t) \cdot g_c(t) = e^{\lambda(e^{it}-1)} \cdot e^{i\mu(e^{it}-1)} = \\ &= e^{i(\lambda+\mu)(e^{it}-1)}\end{aligned}$$

Deoarece funcția caracteristică determină complet funcția de repartiție rezultă că $f + g$ este o variabilă Poisson de parametru $\lambda + \mu$.

QED

4.3 Repartiția uniformă

Spunem că o v.a. are o repartiție uniformă dacă densitatea ei este:

$$\rho(x) = \begin{cases} \frac{1}{b-a} & x \in [a, b] \\ 0 & x \notin [a, b] \end{cases}$$

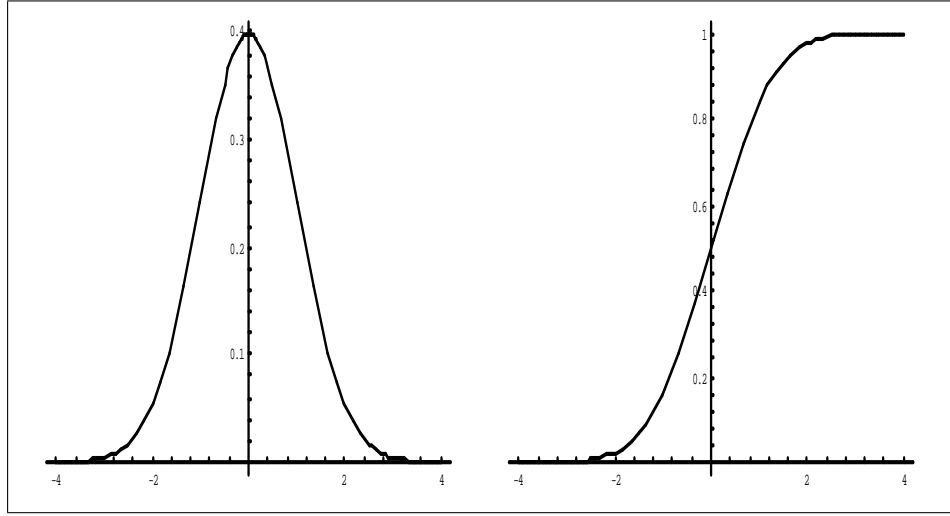
unde $a < b$ sunt reale. Funcția caracteristică este $\int_a^b e^{itx} \frac{1}{b-a} dx = \frac{e^{ibt} - e^{iat}}{it}$, media este $\frac{b+a}{2}$, etc.

4.4 Repartiția Normală

Spunem că o variabilă aleatoare f admite o repartiție normală de parametri m și $\sigma > 0$ dacă densitatea ei de probabilitate este:

$$\rho(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-m)^2}{2\sigma^2}}$$

O variabilă aleatoare cu repartiția normală, de medie m și dispersie σ^2 o vom numi de tipul $N(m, \sigma)$. Graficul ρ pentru $m = 0$ și $\sigma = 1$ al densității de probabilitate și al funcției de repartiție, arată astfel:



Densitatea de probabilitate și funcția de repartiție a unei variabile normale

Deoarece $\int_{-\infty}^{\infty} \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-m)^2}{2\sigma^2}} dx = (\text{prin schimbarea } t = \frac{x-m}{\sigma}) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-\frac{t^2}{2}} dt = 1$, rezultă că ρ este o densitate de probabilitate.

Vom determina la început funcția caracteristică.

$$\begin{aligned} f_c(t) &= \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{itx} e^{-\frac{(x-m)^2}{2\sigma^2}} dx = \\ (\text{schimbarea } y &= \frac{x-m}{\sigma}) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{it(m+\sigma y)} e^{-\frac{y^2}{2}} dy = \\ &= e^{itm - \frac{\sigma^2 t^2}{2}} \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-\frac{(y-it\sigma)^2}{2}} dy = e^{itm - \frac{\sigma^2 t^2}{2}} \end{aligned}$$

Am folosit

$$\int_{-\infty}^{\infty} e^{-\frac{(y-it\sigma)^2}{2}} dy = \int_{-\infty}^{\infty} e^{-\frac{y^2}{2}} dy = \sqrt{2\pi} \quad (4.1)$$

. Acest lucru se arată astfel:

$$\begin{aligned} 0 &= 2\pi i \sum \text{Res} \left(e^{-\frac{z^2}{2}}, z_k \right) = \\ &= \int_{-R}^R e^{-\frac{x^2}{2}} dx + \int_R^{R-it\sigma} e^{-\frac{z^2}{2}} dz - \int_{-R}^{-R-it\sigma} e^{-\frac{z^2}{2}} dz + \int_{-R-it\sigma}^{-R} e^{-\frac{z^2}{2}} dz \end{aligned} \quad (4.2)$$

Suma reziduurilor este în dreptunghiul $[-R, R, R-it\sigma, -R-it\sigma]$. Pe latura $[R, R-it\sigma]$ putem scrie $z = R-ist\sigma$ cu $0 \leq s \leq t\sigma$ deci $z^2 = R^2 - s^2 t^2 \sigma^2 - i2Rst\sigma$ deci $\int_R^{R-it\sigma} e^{-\frac{z^2}{2}} dz = e^{-R^2/2} \int_0^{t\sigma} e^{s^2 t^2 \sigma^2/2} e^{i2Rst\sigma} (-it\sigma) ds \rightarrow 0$ când $R \rightarrow \infty$. Pe latura $[-R-it\sigma, -R]$ se procedează analog și se găsește limita integrale tot zero. Trecând la limită $R \rightarrow \infty$ în (4.2) găsim prima egalitate din (4.1). A doua egalitate se știe din anul I.

Cunoscând funcția caracteristică se pot determina caracteristicile numerice:

- a) $M(f) = \frac{1}{i} f'_c(0) = m$
- b) $M_2(f) = \frac{1}{i^2} f''_c(0) = \sigma^2 + m^2$

c) $D(f) = M_2 - M^2 = \sigma^2$

d) $\mu_{2k+1} = 0 \quad \mu_{2k} = \frac{(2k)!}{2^k k!} \sigma^{2k}$. (exercițiu)

În anumite condiții avem următoarea relație între două variabile normale:

Propoziția 4.3 a) Fie f_1 și f_2 două variabile aleatoare normale independente, de parametri m_1, σ_1 respectiv m_2, σ_2 . Atunci $f_1 + f_2$ este normală de parametri $m = m_1 + m_2$, $\sigma = \sqrt{\sigma_1^2 + \sigma_2^2}$.

b) Dacă f_1, f_2, \dots, f_n sunt v.a. independente și normale de tip $N(m, \sigma)$, atunci $\frac{f_1 + f_2 + \dots + f_n}{n}$ este v.a. normală de tip $N\left(m, \frac{\sigma}{\sqrt{n}}\right)$.

Demonstrație. a) $f_{1c}(f) = e^{im_1 t - \frac{\sigma_1^2 t^2}{2}} \quad f_{2c}(t) = e^{im_2 t - \frac{\sigma_2^2 t^2}{2}}$

$$(f_1 + f_2)_c(t) = f_{1c}(t) f_{2c}(t) = e^{i(m_1 + m_2)t - \frac{(\sigma_1^2 + \sigma_2^2)t^2}{2}} = e^{imt - \frac{\sigma^2 t^2}{2}}$$

cu $m = m_1 + m_2$ $\sigma = \sqrt{\sigma_1^2 + \sigma_2^2}$.

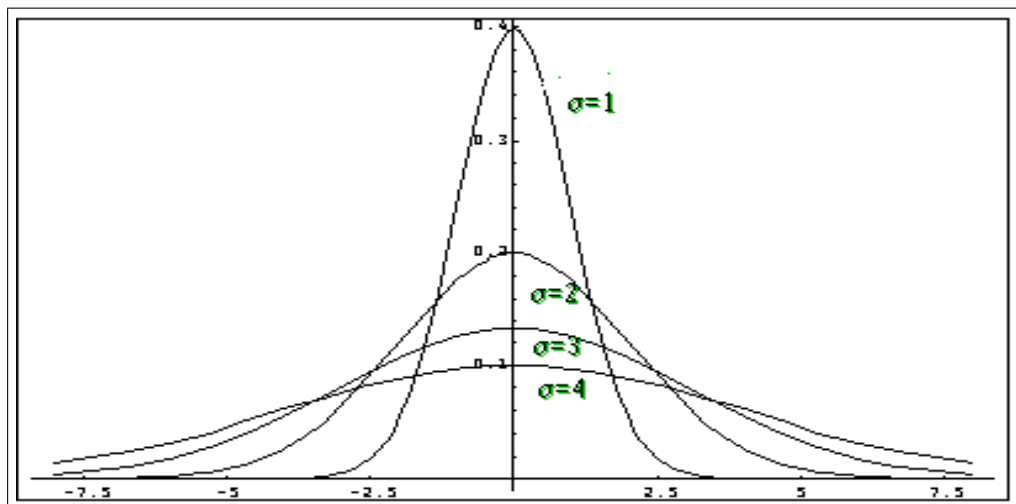
b) Ca la punctul a) găsim $(f_1 + f_2 + \dots + f_n)_c(t) = e^{i \cdot n \cdot m \cdot t - \frac{n \sigma^2}{2} t^2}$. Din teorema 3, lecția 2, rezultă

$$\left(\frac{f_1 + f_2 + \dots + f_n}{n} \right)_c(t) = (f_1 + f_2 + \dots + f_n)_c\left(\frac{t}{n}\right) = e^{imt - \frac{1}{2} \left(\frac{\sigma}{\sqrt{n}}\right)^2 t^2}$$

Cu aceasta b) este demonstrat.

QED.

Pentru $\sigma = 1, 2, 3, 4$ graficele densităților apar în continuare. Se vede că pentru σ mare (adică dispersie mare) graficul este mai împrăștiat (aplatizat) în jurul mediei $m = 0$.



Graficele densităților unor variabile normale
pentru diverse valori ale lui σ

Legea normală a apărut în legătură cu teoria erorilor de măsurare. Să presupunem că o mărime z este determinată prin măsurători. În mod normal se vor face erori. Fie $p(z, x)dx$ probabilitatea ca să obținem rezultatul între x și $x+dx$, dacă valoarea exactă este z . $p(z, x)$ este deci densitatea de probabilitate a rezultatului măsurătorii. Dacă facem n măsurători independente atunci probabilitatea ca să obținem rezultatele în intervalele $[x_1 + dx_1]$, $[x_2 + dx_2]$, ..., $[x_n + dx_n]$ este

$$p(z, x_1)dx_1 p(z, x_2)dx_2 \cdots p(z, x_n)dx_n = F(z, x_1, x_2, \dots, x_n)dx_1 \cdots dx_n \quad (4.3)$$

Care este expresia cea mai potrivită pentru funcția $p(z, x)$? În mare Gauss a plecat de la următoarele ipoteze pentru a determina pe $p(z, x)$:¹

1) $p(z, x) = \phi(z - x)$ (adică erori se fac cu aceeași probabilitate în stânga ca și în dreapta valorii exacte și probabilitatea de a obține prin măsurare o anumită valoare x depinde doar de depărtarea față de valoarea exactă z).

2) Coeficientul lui $dx_1 dx_2 \cdots dx_n$ în n măsurători independente (4.3) este maxim pentru $z = (x_1 + x_2 + \dots + x_n)/n$.

3) ϕ este cel puțin de clasă C^2 .

Cu aceste ipoteze găsim astfel legea erorilor $p(z, x)$:

$F(z, x_1, \dots) = \phi(z - x_1)\phi(z - x_2) \cdots \phi(z - x_n) > 0$ este maxim în același timp cu $\ln(F(z, x_1, x_2, \dots))$.

Prin urmare:

$$(\ln(F(z, x_1, \dots)))'_x = \frac{\phi'(z - x_1)}{\phi(z - x_1)} + \dots + \frac{\phi'(z - x_n)}{\phi(z - x_n)} = 0$$

pentru $z = (x_1 + x_2 + \dots + x_n)/n$. Notând $g(x) = \frac{\phi'(z-x)}{\phi(z-x)}$ avem

$$g(x_1) + g(x_2) + \dots + g(x_n) = 0$$

de fiecare dată când

$$x_1 + x_2 + \dots + x_n = nz$$

adică

$$g(x_1) + g(x_2) + \dots + g(x_{n-1}) + g(nz - x_1 - x_2 - \dots - x_{n-1}) \equiv 0$$

Prin derivare după x_1 obținem:

$$g'(x_1) - g'(nz - x_1 - \dots - x_{n-1}) \equiv 0$$

și datorită arbitrarității valorilor x_1, \dots, x_{n-1} găsim că

$$g'(x) \equiv a = \text{const}$$

¹Demonstrația este adaptată după H. Poincaré-Calcul des probabilités, Gauthiers-Villars, Paris, 1912

de unde

$$g(x) = ax + b$$

Deoarece $0 = g(x_1) + g(x_2) + \dots + g(nz - x_1 - \dots - x_{n-1}) = anz + nb$ rezultă $b = -az$, deci:

$$\frac{\phi'(z-x)}{\phi(z-x)} = a(x-z)$$

sau

$$\frac{\phi'(t)}{\phi(t)} = -at$$

de unde

$$\begin{aligned} \ln(\phi(t)) &= -a\frac{t^2}{2} + c \\ \phi(t) &= e^c e^{-a\frac{t^2}{2}} \end{aligned}$$

Deci:

$$p(z, x) = \phi(z-x) = k e^{-a\frac{(z-x)^2}{2}}$$

Deoarece $p(z, x)$ este o densitate de probabilitate trebuie ca

$$\int_{-\infty}^{\infty} p(z, x) dx = 1$$

Aceasta implică $a > 0$ și $k = \sqrt{\frac{a}{2\pi}}$. Dacă luăm $a = \frac{1}{\sigma^2}$ atunci densitatea p devine:

$$p(z, x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-z)^2}{2\sigma^2}}$$

adică ceea ce am numit legea normală.

Alte motive pentru care legea normală este importantă vom vedea la considerarea teoremelor de tip limită centrală. Fiind frecvent folosită valorile funcției de repartiție pentru $m = 0$, $\sigma = 1$, adică $\int_{-\infty}^t \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx$ au fost tabelate. Mai precis s-a introdus funcția

$$\Phi(t) = \int_0^t \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx$$

. Cu ajutorul acestei funcții avem

$$p(a \leq f < b) = \int_a^b \rho(x) dx = \Phi(b) - \Phi(a)$$

pentru o v.a. de tipul $N(0, 1)$. Φ este o funcție antisimetrică, $\Phi(a) = -\Phi(-a)$. Valorile funcției Φ s-au tabelat pentru diverse valori ale lui t . Dacă f este o v.a. normală de parametri m și σ atunci variabila $Z = \frac{f-m}{\sigma}$ este normală cu media 0 și dispersia 1 (exercițiu), deci cu funcția de repartiție tabelată.

4.5 Repartiția exponențială negativă

Se numește astfel o repartiție cu densitatea

$$\rho(x) = \begin{cases} 0 & x < 0 \\ \lambda e^{-\lambda x} & x \geq 0 \end{cases}$$

unde $\lambda > 0$.

Plecând de la definiție se găsește $f_c(t) = \frac{\lambda}{\lambda - it}$ $M = \frac{1}{\lambda}$ $M_2 = \frac{2}{\lambda^2}$ $D = \frac{1}{\lambda^2}$. Toate aceste calcule sunt lăsate ca exercițiu.

4.6 Repartiția Gamma

Se numește astfel o repartiție cu densitatea

$$\rho_{\alpha,\beta}(x) = \begin{cases} 0 & x \leq 0 \\ \frac{1}{\Gamma(\alpha+1)\beta^{\alpha+1}} x^\alpha e^{-\frac{x}{\beta}} & x > 0 \end{cases}$$

unde $\alpha > -1$ și $\beta > 0$. Cel mai frecvent se folosește cazul $\beta = 1$ $\alpha + 1 = m$, când avem

$$\rho(x) = \rho_{m+1,0}(x) = \begin{cases} 0 & x \leq 0 \\ \frac{1}{\Gamma(m)} x^{m-1} e^{-x} & x > 0 \end{cases}$$

Reamintim că funcția gamma se definește astfel

$$\Gamma(x) = \int_0^\infty e^{-t} t^{x-1} dt$$

pentru $x > 0$. Următoarele proprietăți se cunosc de la cursul de analiză:

$$\begin{aligned} \Gamma(x+1) &= x\Gamma(x) \\ \Gamma\left(\frac{1}{2}\right) &= \sqrt{\pi} \end{aligned}$$

Funcția caracteristică se calculează astfel:

$$\begin{aligned}
f_c(t) &= \frac{1}{\Gamma(\alpha+1)\beta^{\alpha+1}} \int_0^\infty e^{itx} x^\alpha e^{-\frac{x}{\beta}} dx \\
&= \frac{1}{\Gamma(\alpha+1)\beta^{\alpha+1}} \int_0^\infty \sum_{n=0}^\infty \frac{(itx)^n}{n!} x^\alpha e^{-\frac{x}{\beta}} dx \\
&= \frac{1}{\Gamma(\alpha+1)\beta^{\alpha+1}} \sum_{n=0}^\infty \frac{(it)^n}{n!} \int_0^\infty x^{n+\alpha} e^{-\frac{x}{\beta}} dx \\
&= \sum_{n=0}^\infty (it)^n \frac{1}{\Gamma(\alpha+1)\beta^{\alpha+1}} \frac{\beta^{n+\alpha+1}}{n!} \int_0^\infty y^{n+\alpha} e^{-y} dy \\
&= \sum_{n=0}^\infty (it)^n \frac{1}{\Gamma(\alpha+1)\beta^{\alpha+1}} \frac{\beta^{n+\alpha+1}}{n!} \Gamma(n+\alpha+1) \\
&= \sum_{n=0}^\infty (it)^n \frac{\beta^{n+\alpha+1} (n+\alpha)(n+\alpha-1)\cdots(\alpha+1)\Gamma(\alpha+1)}{\Gamma(\alpha+1)\beta^{\alpha+1}n!} \\
&= \sum_{n=0}^\infty (-it\beta)^n \cdot \frac{(-\alpha-1)(-\alpha-2)\cdots(\alpha-n)}{n!} = (1-it\beta)^{-\alpha-1}
\end{aligned}$$

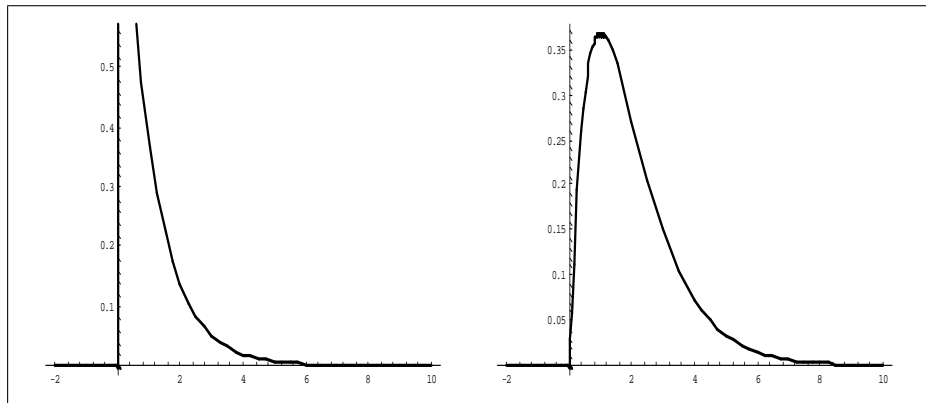
Am folosit aici dezvoltarea

$$(1+x)^\gamma = 1 + \frac{\gamma}{1!}x + \frac{\gamma(\gamma-1)}{2!}x^2 + \dots + \frac{\gamma(\gamma-1)\dots(\gamma-n+1)}{n!}x^n + \dots$$

Caracteristicile numerice sunt:

- a) $M(f) = \frac{1}{i} f'_c(0) = \beta(\alpha+1)$
- b) $M_2(f) = \frac{1}{i^2} f''_c(0) = \beta^2(\alpha+1)(\alpha+2)$
- c) $D(f) = M_2(f) - M(f)^2 = \beta^2(\alpha+1)$

În cazul $\beta = 1$ $\alpha = m - 1$ se obține $f_c(t) = (1 - mt)^{-m}$, media= m , dispersia= m . Tipul de mai sus de variabilă aleatoare îl vom numi $\gamma(\alpha, \beta)$ iar în cazul particular considerat îl vom numi $\gamma(m)$. Pentru $m = 0$ și $m = 2$ graficele densităților apar în continuare:



Densitățile $\gamma(m)$ pentru $m = 0$ și $m = 2$.

Propoziția următoare este simplu de demonstrat:

Propoziția 4.4 Dacă f și g sunt variabile aleatoare independente de tipurile $\gamma(m_1)$ respectiv $\gamma(m_2)$ atunci $f + g$ este de tipul $\gamma(m_1 + m_2)$.

Demonstrație.Exercițiu

QED.

4.7 Repartiția χ^2 (hi pătrat)

Se numete astfel o v.a. cu densitatea :

$$\rho(x) = \begin{cases} 0 & x \leq 0 \\ \frac{1}{2^{\frac{s}{2}} \sigma^s \Gamma(\frac{s}{2})} x^{\frac{s}{2}-1} e^{-\frac{x}{2\sigma^2}} & x > 0 \end{cases}$$

s este un număr natural numit numărul gradelor de libertate, iar $\sigma > 0$. Acest tip de variabilă aleatoare îl vom numi $H(s, \sigma)$; se vede imediat că este același cu $\gamma(\frac{s}{2} - 1, 2\sigma^2)$. Prin urmare avem:

- a) $f_c(t) = (1 - 2\sigma^2 t)^{-\frac{s}{2}}$
- b) $M = s\sigma^2$
- c) $M_2 = s(s+2)\sigma^4$
- c) $D = 2s\sigma^4$

Iată cum se poate ajunge la o distribuție χ^2 plecând de la distribuții normale:

Propoziția 4.5 Fie f_1, f_2, \dots, f_s variabile aleatoare independente de tipu $N(0, \sigma)$. Atunci variabila aleatoare $g = f_1^2 + f_2^2 + \dots + f_s^2$ este de tipul $H(s, \sigma)$.

Demonstrație.Fie $F_i(t)$ funcția de repartiție a variabilei f_i^2 . Avem

$$\begin{aligned} F_i(t) &= p(f_i^2 < t) = p(-\sqrt{t} < f_i < \sqrt{t}) = \frac{1}{\sqrt{2\pi}\sigma} \int_{-\sqrt{t}}^{\sqrt{t}} e^{-\frac{x^2}{2\sigma^2}} dx = \\ &= \frac{2}{\sigma\sqrt{2\pi}} \int_0^{\sqrt{t}} e^{-\frac{x^2}{2\sigma^2}} dx \end{aligned}$$

De aici prin derivare determinăm densitățile pentru f_i^2 , notate ρ_i . Avem

$$\rho_i(t) = F_i'(t) = \frac{1}{\Gamma(\frac{1}{2})(2\sigma^2)^{1/2}} t^{-\frac{1}{2}} e^{-\frac{t}{2\sigma^2}}$$

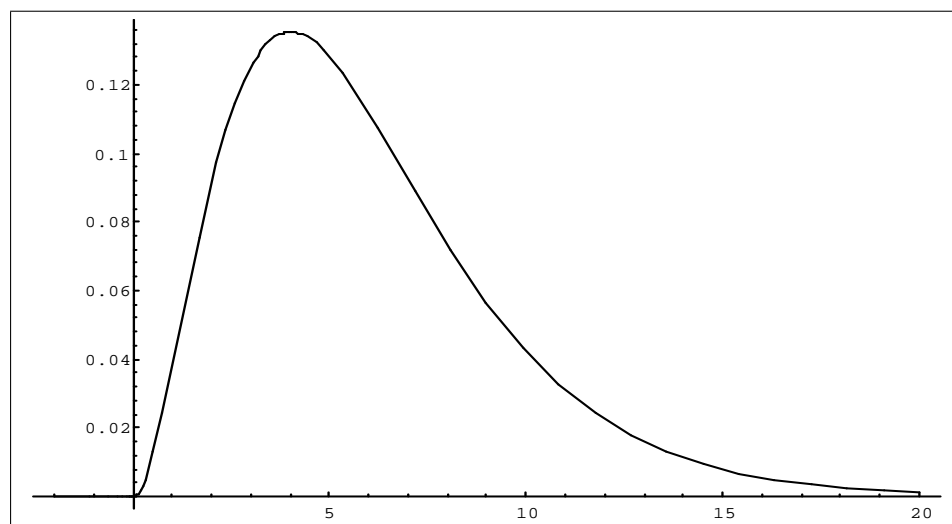
pentru $t > 0$ ceea ce pune în evidență faptul că f_i^2 sunt distribuții de tip $\gamma(-\frac{1}{2}, 2\sigma^2)$, deci $(f_i^2)_c(t) = (1 - 2\sigma^2 t)^{-\frac{1}{2}}$. Deoarece $f_1^2, f_2^2, \dots, f_s^2$ sunt independente rezultă că

$$(f_1^2 + f_2^2 + \dots + f_s^2)_c(t) = (f_1^2)_c(t) \cdot (f_2^2)_c(t) \cdots (f_s^2)_c(t) = (1 - 2\sigma^2 t)^{-\frac{s}{2}}$$

care este chiar funcția caracteristică a unei distribuții $H(m, \sigma)$.

QED.

Pentru $\sigma = 1$ și $s = 6$ grade de libertate, graficul densității arată astfel:



Densitatea χ^2 pentru $s = 6$ și $\sigma = 1$.

4.8 Repartiția Student

O v.a. f , cu densitatea de probabilitate

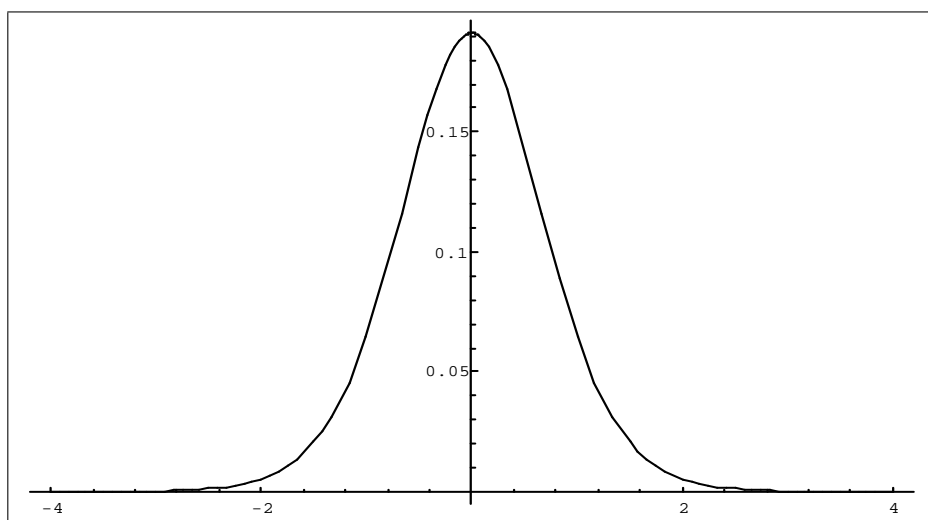
$$\rho(x) = \frac{\Gamma\left(\frac{s+1}{2}\right)}{\sqrt{\pi s} \Gamma\left(\frac{s}{2}\right)} \left(1 + \frac{x^2}{s}\right)^{-\frac{s+1}{2}}$$

se numește variabilă Student. Tipul acesta de v.a. îl vom nota $S(s)$. Demonstrațiile afirmațiilor de mai jos în legătură cu distribuțiile Student se vor da în lecția 6.

a) $M(f) = 0$ și în general $\mu_{2k+1}(f) = 0$ pentru $0 \leq k < \frac{s}{2}$. Pentru $k \geq \frac{s}{2}$, $\mu_{2k+1}(f)$ nu există.

$$\text{b) } \mu_{2k}(f) = M_{2k}(f) = \frac{s^r}{\sqrt{\pi}} \frac{\Gamma(k+\frac{1}{2})\Gamma(\frac{s}{2}-k)}{\Gamma(\frac{s}{2})} = \frac{s^k 1 \cdot 3 \cdot \dots \cdot (2k-1)}{(s-2) \cdot (s-4) \cdot \dots \cdot (s-2k)}.$$

c) Dacă ξ este o v.a. de tip normal $N(0, \sigma)$ și η este de tip χ^2 cu s grade de libertate, adică de tip $H(s, \sigma)$ atunci v.a. $\frac{\xi}{\sqrt{\frac{\eta}{s}}}$ este de tip Student $S(s)$. Pentru $s = 6$, graficul densității este următorul:

Densitatea Student pentru $s = 6$

Această repartiție este frecvent întâlnită în Statistică și este tabelată pentru diverse valori ale lui s .

4.9 Rezumat

Variabilele aleatoare anterioare apar în multe modele de probleme reale. De aceea au fost studiate mai în detaliu.

i) Modelul binomial descrie probabilitatea de apariție de k ori a unui eveniment într-un șir de n experiențe independente, dacă probabilitatea de apariție a evenimentului într-o experiență este p . Această probabilitate este $C_n^k p^k (1-p)^{n-k}$.

ii) Modelul Poisson descrie probabilitatea de apariție de k ori a unui eveniment într-un șir mare de n experiențe independente, dacă probabilitatea de apariție a evenimentului într-o experiență este foarte mică, de tipul λ/n . Numărul mediu de apariții este același, λ , iar probabilitatea este $e^{-\lambda} \lambda^k / k!$.

iii) Legea normală, cu densitatea $\rho(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-m)^2}{2\sigma^2}}$ este legată de distribuția erorilor de măsurare. m este valoarea medie a valorii măsurate iar σ^2 este măsura dispersiei valorilor găsite în jurul mediei.

iv) Distribuția χ^2 cu s grade de libertate are densitatea

$$\rho(x) = \begin{cases} 0 & x \leq 0 \\ \frac{1}{2^{\frac{s}{2}} \sigma^s \Gamma(\frac{s}{2})} x^{\frac{s}{2}-1} e^{-\frac{x}{2\sigma^2}} & x > 0 \end{cases}$$

și este distribuția sumei pătratelor a s variabile aleatoare normale independente, de medie 0 și dispersie σ^2 . Această distribuție apare frecvent în statistică.

Alte legi clasice sunt date în exerciții sau vor fi studiate mai târziu, pe măsură ce vor fi folosite. E de remarcat utilitatea funcției caracteristice în calculul momentelor.

4.10 Exerciții

1. Să se arate că

$$\mu_{2k} = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{\infty} (x-m)^{2k} e^{-\frac{(x-m)^2}{2\sigma^2}} dx = (2k-1)\sigma^2 \mu_{2k-2}$$

de unde

$$\mu_{2k} = \frac{(2k)!}{2^k k!} \sigma^{2k}$$

Indicație. Cu substituția $t = \frac{x-m}{\sigma}$ se găsește

$$\mu_{2k} = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} x^{2k} e^{-\frac{x^2}{2}} dx = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} x^{2k-1} \left(-e^{-\frac{x^2}{2}}\right)' dx = \dots$$

2. Fie o v.a. cu densitatea $\rho(x) = \frac{1}{2\lambda} e^{-\frac{|x-m|}{\lambda}}$ $\lambda > 0$ (distribuția lui Laplace).

a) Să se arate că $f_c(t) = \frac{e^{imt}}{1+\lambda^2 t^2}$.

b) Să se determine media și dispersia.

Indicație. Se aplică definițiile.

3. Să se arate că dacă o v.a. f este de tipul $N(m, \sigma^2)$, atunci $g = \frac{(f-m)^2}{2\sigma^2}$ este de tipul $\gamma(\frac{1}{2}) = \gamma(-\frac{1}{2}, 0)$.

Indicație. Se arată mai întâi că $h = \frac{f-m}{\sigma}$ este de tipul $N(0, 1)$. Pe urmă după modelul de la distribuția χ^2 se arată că $g = h^2$ e de tipul cerut.

4. Fie f_1, f_2, \dots, f_k variabile aleatoare independente, fiecare având o repartiție exponențială negativă.

a) Să se arate că $f = f_1 + f_2 + \dots + f_k$ are ca funcție caracteristică $f_c(t) = \frac{\lambda^k}{(\lambda - it)^k}$.

b) Fie variabila aleatoare g cu densitatea $\rho(x) = \frac{\lambda^k}{\Gamma(k)} x^{k-1} e^{-\lambda x}$ dacă $x \geq 0$ și $\rho(x) = 0$ dacă $x < 0$. Să se arate că g este de tipul $\gamma(\alpha, \beta)$ pentru un α și β . Sa se determine funcția caracteristică a lui g .

c) Sa se arate că $f_1 + f_2 + \dots + f_k$ și g au aceleași repartiții.

Indicație. Utilizăm funcțiile caracteristice. Funcția caracteristică a lui $f_1 + f_2 + \dots + f_n$ este produsul funcțiilor caracteristice ale variabilelor f_1, \dots, f_n . Aceste funcții caracteristice sunt

calculate în lecția prezentă.

5. Să se arate că într-o distribuție Poisson de parametru λ valoarea k cea mai probabilă satisface inegalitatea $\lambda - 1 \leq k \leq \lambda$.

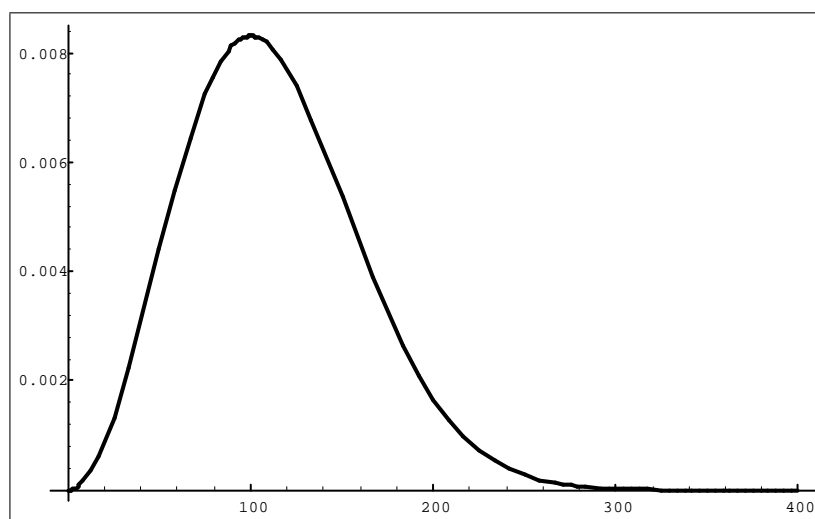
Indicație. Se face raportul a două probabilități vecine, $p(f = k)$ și $p(f = k + 1)$.

6. Viteza absolută a unei molecule într-un gaz perfect este o mărime aleatoare. Conform teoriei lui Maxwell densitatea de probabilitate a acestei viteze este:

$$\begin{aligned}\rho(x) &= \frac{4x^2}{c^3\sqrt{\pi}}e^{-\frac{x^2}{c^2}}, \text{ pentru } x > 0 \\ \rho(x) &= 0, \text{ pentru } x \leq 0\end{aligned}$$

unde c este o constantă. Să se determine viteza medie, energia cinetică medie și dispersiile lor.

Indicație. Pentru $c=100$, graficul distribuției vitezelor este:



Distribuția vitezelor într-un gaz perfect la $c=100$

Viteza medie $= \int_{-\infty}^{\infty} x\rho(x)dx$, energia cinetică medie $= \int_{-\infty}^{\infty} \frac{mx^2}{2}\rho(x)dx$, unde m este masa unei molecule. Integralele se fac prin părți sau se folosesc rezultatele de la legea normală (pb. 1).

7. Distribuția cu densitatea

$$\rho(x) = \begin{cases} 0 & x \leq 0 \\ \frac{1}{x\beta\sqrt{2\pi}}e^{-\frac{(\ln x - \alpha)^2}{2\beta^2}} & x > 0 \end{cases}$$

se numește normal logaritmică. Se cere media și dispersia unei astfel de distribuții.

Indicație. Dacă ξ este o v.a. cu astfel de lege atunci:

$$\begin{aligned} M(\xi) &= \int_0^\infty x \frac{1}{x\beta\sqrt{2\pi}} e^{-\frac{(\ln x - \alpha)^2}{2\beta^2}} dx \\ &= \frac{1}{\sqrt{2\pi}\beta} \int_{-\infty}^\infty e^{t - \frac{(t-\alpha)^2}{2\beta^2}} dt = e^{\alpha + \frac{\beta^2}{2}}. \end{aligned}$$

Analog se procedează pentru dispersie, găsiindu-se $D(f) = e^{2\alpha + \beta^2} (e^{\beta^2} - 1)$ A.N. Kolmogorov a arătat că aceasta este legea de distribuție a diametrelor particulelor într-un proces de măcinare.

8. O variabilă aleatoare f cu densitatea

$$\rho(x) = \begin{cases} 0 & x \leq 0 \\ c\alpha x^{\alpha-1} e^{-cx^\alpha} & x > 0 \end{cases}$$

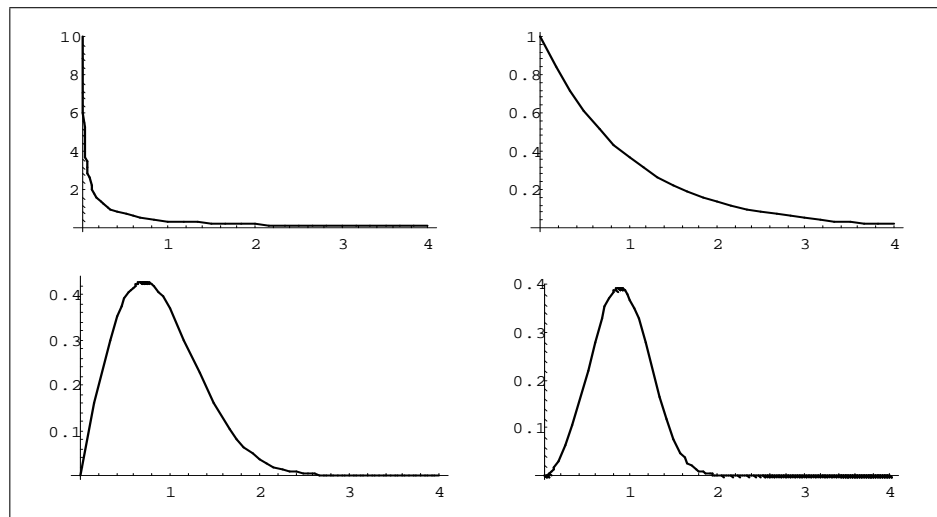
se numește variabilă Weibull.

a) Care media lui f ?

b) Care e dispersia lui f ?

c) Care este momentul de ordin k al lui f ?

Indicație. Folosind formulele din lecția 3 pentru momentul de ordin k și utilizând schimbarea $x^\alpha = t$ ajungem la o funcție Γ . Acum media și dispersia se calculează imediat. În figura următoare sunt reprezentate graficele densităților pentru $\alpha = 0,5; 1; 2; 3$ la $c=1$.



Densitatea Weibull pentru $\alpha = 0,5; 1; 2; 3$.

9. Să se arate că funcția $\rho(x)$ care maximizează expresia

$$I = - \int_{-\infty}^{\infty} \rho(x) \ln(\rho(x)) dx$$

cu condițiile

$$\int_{-\infty}^{\infty} \rho(x) dx = 1 \text{ și } \int_{-\infty}^{\infty} x^2 \rho(x) dx = \sigma^2$$

este $\rho(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{x^2}{2\sigma^2}}$. (I se numește entropia densității ρ . Prin urmare distribuția normală are entropie maximă, la o dispersie dată).

Indicație. Se știe din Calculul Variațional că funcția $\rho(x)$ care realizează extremul funcționalei I , cu cele două constrângeri, trebuie să satisfacă ecuația $\frac{\partial F}{\partial \rho} - \frac{\partial}{\partial \rho} \left(\frac{\partial F}{\partial \rho'} \right) = 0$, unde $F(x, \rho, \rho') = -\rho \ln \rho + \lambda_1 \rho + \lambda_2 x^2 \rho$.

10. Să se arate că distribuția pe $[0, \infty)$ care maximizează entropia și are media $m > 0$, este exponențială negativă $\rho(x) = \frac{1}{m} e^{-\frac{x}{m}}$.

Indicație. Trebuie maximizată integrala $I = -\int_0^\infty \rho(x) \ln \rho(x) dx$ în condițiile $\int_0^\infty \rho(x) dx = 1$, $\int_0^\infty x \rho(x) dx = m$. (vezi problema 9).

11. Să se arate că distribuția pe intervalul $[a, b]$ care maximizează entropia, este distribuția uniformă.

12. Într-o urnă se găsesc a bile albe și b bile negre. Se extrag n bile din care k sunt albe și $n - k$ sunt negre ($k \leq a$, $n - k \leq b$). Extragerea se face fără punerea bilei înapoi (schema bilei neîntoarce). Să se arate că probabilitatea acestui eveniment este $\frac{C_a^k C_b^{n-k}}{C_{a+b}^n}$. Să se arate că pentru variabila aleatoare

$$f = \begin{pmatrix} \dots & k & \dots \\ \dots & \frac{C_a^k C_b^{n-k}}{C_{a+b}^n} & \dots \end{pmatrix}$$

cu $k \in [\max(0, n - b), \min(a, n)]$ media este $M(f) = n \frac{a}{a+b}$ și dispersia este $D(f) = \frac{a+b-n}{a+b-1} n \frac{a}{a+b} \frac{b}{a+b}$.

Indicație. Se utilizează identitatea $\sum_k C_a^k C_b^{n-k} = C_{a+b}^n$ unde însumarea se face între limitele specificate în problemă pentru k (identitatea se poate obține egalând coeficienții lui x^n în $(1+x)^a (1+x)^b = (1+x)^{a+b}$). Pentru medie se utilizează $k C_a^k = a C_{a-1}^{k-1}$ iar pentru dispersie se utilizează $k(k-1) C_a^k = a(a-1) C_{a-2}^{k-2}$.

13. Față de un adversar la fel de tare ce este mai probabil să se câștige: trei partide din patru sau cinci partide din opt?

14. Pentru o variabilă aleatoare se definește asimetria prin $\gamma_1 = \frac{\mu_3}{\sqrt{\mu_2^3}}$ și excesul prin $\gamma_2 = \frac{\mu_4}{\mu_2^2} - 3$. Să se calculeze valorile acestor mărimi pentru variabilele aleatoare din această lecție.

Lecția 5

Legi limită

În cele ce urmează vom studia semnificația dispersiei, precum și importanța repartiției normale pentru probabilități. Începem cu inegalitatea lui Cebîșev.

Teorema 5.1 *Fie f o variabilă aleatoare cu medie și dispersie finite. Atunci pentru $a > 0$ avem:*

$$p(|f - M(f)| > a) \leq \frac{D(f)}{a^2}$$

sau

$$p(|f - M(f)| \leq a) \geq 1 - \frac{D(f)}{a^2} \quad (5.1)$$

Demonstrație.

$$\begin{aligned} p(|f - M(f)| > a) &= \int_{|x - M(f)| > a} dF(x) \leq \int_{|x - M(f)| > a} \frac{(x - M(f))^2}{a^2} dF(x) \leq \\ &\leq \frac{1}{a^2} \int_{|x - M(f)| > a} (x - M(f))^2 dF(x) \leq \frac{1}{a^2} \int_{-\infty}^{\infty} (x - M(f))^2 dF(x) = \frac{D(f)}{a^2} \end{aligned}$$

QED.

Putem pune acest lucru și în alt mod dacă scriem $a = \nu\sigma = \nu\sqrt{D}$. În acest caz avem $\frac{D}{a^2} = \frac{1}{\nu^2}$. Inegalitatea devine

$$\begin{aligned} p(|f - M(f)| > \nu\sigma) &\leq \frac{1}{\nu^2} \text{ sau} \\ p(|f - M(f)| \leq \nu\sigma) &\geq 1 - \frac{1}{\nu^2} \end{aligned}$$

Aceasta înseamnă că variabila aleatoare ia valori în intervalul $[M(f) - \nu\sigma, M(f) + \nu\sigma]$ cu o probabilitate mai mare ca $1 - \frac{1}{\nu^2}$. În cazul particular al unei variabile normale $f \in N(m, \sigma)$ găsim că f ia valori în intervalul $[m - 3\sigma, m + 3\sigma]$ cu probabilitatea $\geq 1 - \frac{1}{9} = 0,88..$ adică aria de sub graficul densității $\rho(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-m)^2}{2\sigma^2}}$ cuprins între $x = m - 3\sigma$ și $x = m + 3\sigma$ este cel puțin 0,88... Cu cât σ este mai mic, cu atât mai mic este intervalul în care funcția ia valori cu probabilitate mare (valori în jurul mediei m).

Observația 5.2 Estimarea dată de inegalitatea lui Cebîșev este destul de grosieră în cazuri practice și se înlocuiește cu alte estimări mai eficiente. De exemplu

$$p(-3 < f < 3) = \Phi(3) - \Phi(-3) = 2\Phi(3) = 0,9973..$$

pentru $f \in N(0, 1)$ este o estimare mai bună ca cea dată de inegalitatea lui Cebîșev

$$p(-3 < f < 3) \geq 0,88..$$

5.1 Legea numerelor mari

Următoare teoremă este o ilustrare a aplicării inegalității lui Cebîșev.

Teorema 5.3 (Legea numerelor mari) Fie $f_1, f_2, \dots, f_n, \dots$ un șir de v.a. independente două câte două și cu dispersiile mărginite de aceeași constantă: $D(f_n) \leq C$, pentru orice $n \in N$. Atunci pentru orice $\epsilon > 0$ avem:

$$\lim_{n \rightarrow \infty} p \left(\left| \frac{f_1 + f_2 + \dots + f_n}{n} - \frac{M(f_1) + \dots + M(f_n)}{n} \right| \leq \epsilon \right) = 1 \quad (5.2)$$

Mai precis

$$p \left(\left| \frac{f_1 + f_2 + \dots + f_n}{n} - \frac{M(f_1) + \dots + M(f_n)}{n} \right| \leq \epsilon \right) \geq 1 - \frac{C}{n\epsilon^2} \quad (5.3)$$

Demonstrație. Fie $g_n = \frac{f_1 + f_2 + \dots + f_n}{n}$. Avem $M(g_n) = \frac{M(f_1) + \dots + M(f_n)}{n}$ și $D(g_n) = \frac{1}{n^2}(D(f_1) + \dots + D(f_n)) \leq \frac{C}{n^2}$. Inegalitatea lui Cebîșev pentru g_n dă (5.3) iar prin trecere la limită se obține (5.2).

QED.

Teorema se interpretează astfel: oricare ar fi $\epsilon > 0$, **dacă n este suficient de mare**, atunci funcția $(f_1 + f_2 + \dots + f_n)/n$ diferă cu mai puțin de ϵ de constanta $(M(f_1) + \dots + M(f_n))/n$ pe o mulțime cu probabilitate foarte aproape de 1. Dacă f_1, \dots, f_n, \dots au aceeași medie m atunci avem:

Corolarul 5.4 Fie v.a $f_1, f_2, \dots, f_n, \dots$ cu aceeași medie m și dispersiile mărginite de C . Atunci pentru orice $\epsilon > 0$ avem:

$$\lim_{n \rightarrow \infty} p \left(\left| \frac{f_1 + f_2 + \dots + f_n}{n} - m \right| \leq \epsilon \right) = 1 \quad (5.4)$$

Un alt caz particular este:

Corolarul 5.5 (teorema lui Bernoulli) Fie μ_n numărul de apariții ale unui eveniment A în n experiențe independente. Probabilitatea de realizare a lui A într-o singură experiență este $\alpha \in [0, 1]$. Atunci pentru orice $\epsilon > 0$ avem

$$\lim_{n \rightarrow \infty} p\left(\left|\frac{\mu_n}{n} - \alpha\right| \leq \epsilon\right) = 1 \quad (5.5)$$

Observația 5.6 Cu alte cuvinte pentru un $\epsilon > 0$ oricât de mic, dacă numărul de experiențe n este suficient de mare atunci frecvența relativă $\frac{\mu_n}{n}$ este aproape sigur în jurul probabilității p de apariție a evenimentului într-o singură experiență.

Demonstrție. Fie șirul de v.a.

$$f_n(\omega) = \begin{cases} 1 & \text{dacă la experiența } n \text{ s-a realizat } A \\ 0 & \text{dacă la experiența } n \text{ s-a realizat } \bar{A} \end{cases} \quad (5.6)$$

Aici ω este un șir de experiențe independente. Fiecare v.a. f_n are diagrama

$$f_n \begin{pmatrix} 1 & 0 \\ \alpha & \beta \end{pmatrix}$$

cu $\beta = 1 - \alpha$. Toate variabilele f_n au mediile α și dispersiile $\alpha\beta$. Ținând seama că

$$(f_1 + f_2 + \dots + f_n) = \mu_n$$

prin aplicarea formulei (5.4) rezultă teorema lui Bernoulli.

O altă demunstratie se poate da astfel:

- i) μ_n are o distribuție Bernoulli cu $p(\mu_n = k) = C_n^k \alpha^k \beta^{n-k}$.
- ii) $M(\mu_n) = n\alpha$ și $D(\mu_n) = n\alpha\beta$ deci $M(\frac{\mu_n}{n}) = \alpha$ și $D(\frac{\mu_n}{n}) = \frac{\alpha\beta}{n}$.
- iii) Aplicând teorema lui Cebășev lui μ_n/n găsim

$$\begin{aligned} 1 &\geq p\left(\left|\frac{\mu_n}{n} - \alpha\right| \leq \epsilon\right) \geq 1 - \frac{D(\mu_n)}{\epsilon^2} \\ &= 1 - \frac{pq}{n\epsilon^2} \geq 1 - \frac{1}{4n\epsilon^2} \rightarrow 1 \end{aligned} \quad (5.7)$$

când $n \rightarrow \infty$, pentru că $\alpha\beta = \alpha(1 - \alpha) \leq \frac{1}{4}$ atunci când $0 \leq \alpha \leq 1$.

QED.

Observația 5.7 Avem

$$\sum_{\left|\frac{k}{n} - \alpha\right| \leq \epsilon} C_n^k \alpha^k \beta^{n-k} = p\left(\left|\frac{\mu_n}{n} - \alpha\right| \leq \epsilon\right) \geq 1 - \frac{1}{4n\epsilon^2} \quad (5.8)$$

sau

$$\sum_{\left|\frac{k}{n} - \alpha\right| > \epsilon} C_n^k \alpha^k \beta^{n-k} \leq \frac{1}{4n\epsilon^2} \quad (5.9)$$

Datorită frecvenței schemei binomiale și dificultății de a calcula $C_n^k \alpha^k \beta^{n-k}$ s-au dezvoltat tehnici de calcul aproximativ, mai rapide, pentru aceste mărimi.

5.2 Teoreme limită centrală

Avem nevoie în continuare de următoarea formulă pentru factorial:

$$n! = \sqrt{2\pi n} \left(\frac{n}{e}\right)^n \cdot e^{\frac{\theta_n}{12n}}, \quad 0 < \theta_n < 1$$

numită formula lui **Stirling**. Demonstrația acestei formule se poate găsi într-un curs de Analiză matematică.

Teorema 5.8 (formula locală de Moivre-Laplace) Fie $0 < p < 1$ și $x_k = \frac{k-np}{\sqrt{npq}}$, $q = 1 - p$. Atunci

$$\frac{C_n^k p^k q^{n-k}}{\frac{1}{\sqrt{npq}} \frac{1}{\sqrt{2\pi}} e^{-\frac{x_k^2}{2}}} \rightarrow 1 \quad (5.10)$$

când $n \rightarrow \infty$, uniform pentru $a \leq x_k \leq b$, cu a și b finite.

Demonstrație. Avem:

$$\sqrt{npq} C_n^k p^k q^{n-k} = \sqrt{npq} \sqrt{\frac{n}{2\pi k(n-k)}} \frac{n^n p^k q^{n-k}}{k^k (n-k)^{n-k}} \cdot e^{\theta_{n,k}} \quad (5.11)$$

unde $\theta_{n,k} = \frac{\theta_n}{n} - \frac{\theta_k}{k} - \frac{\theta_{n-k}}{n-k}$.

i) Din $a \leq \frac{k-np}{\sqrt{npq}} \leq b$ deducem:

$$\begin{aligned} k &\geq np + a\sqrt{npq} = np(1 + a\sqrt{\frac{q}{np}}) \\ n-k &\geq nq - b\sqrt{npq} = nq(1 - b\sqrt{\frac{q}{nq}}) \end{aligned}$$

deci:

$$\begin{aligned} 0 &< |\theta_{n,k}| < \frac{\theta_n}{n} + \frac{\theta_k}{k} + \frac{\theta_{n-k}}{n-k} < \frac{1}{12} \left(\frac{1}{n} + \frac{1}{k} + \frac{1}{n-k} \right) < \\ &< \frac{1}{12n} \left(1 + \frac{1}{p + a\sqrt{\frac{pq}{n}}} + \frac{1}{q - \sqrt{\frac{pq}{n}}} \right) \end{aligned}$$

Prin urmare $\theta_{n,k}$ tinde uniform la zero când $n \rightarrow \infty$ și $a \leq x_k \leq b$, deci $e^{\theta_{n,k}}$ tinde uniform la 1.

ii) Scriind că $k = np + x_k \sqrt{npq}$, găsim:

$$\sqrt{npq} \sqrt{\frac{n}{2\pi k(n-k)}} = \sqrt{\frac{1}{2\pi}} \sqrt{\frac{1}{\left(1 + x_k \sqrt{\frac{q}{np}}\right) \left(1 - x_k \sqrt{\frac{q}{np}}\right)}} \rightarrow \frac{1}{\sqrt{2\pi}}$$

uniform pentru $a \leq x_k \leq b$.

iii) Mai avem în (5.11) de aflat limita lui

$$A_{n,k} = \frac{n^n p^k q^{n-k}}{k^k (n-k)^{n-k}}$$

Avem:

$$\begin{aligned} \ln(A_{n,k}) &= \ln \left(\left(\frac{np}{k} \right)^k \cdot \left(\frac{nq}{n-k} \right)^{n-k} \right) \\ &= -k \ln \left(\frac{k}{np} \right) - (n-k) \ln \left(\frac{n-k}{nq} \right) \\ &= -(np + x_k \sqrt{npq}) \ln \left(\frac{np + x_k \sqrt{npq}}{np} \right) - \\ &\quad - (nq - x_k \sqrt{npq}) \ln \left(\frac{nq - x_k \sqrt{npq}}{nq} \right) \\ &= -(np + x_k \sqrt{npq}) \ln \left(1 + x_k \sqrt{\frac{q}{np}} \right) - \\ &\quad - (nq - x_k \sqrt{npq}) \ln \left(1 - x_k \sqrt{\frac{p}{nq}} \right) \end{aligned}$$

Dezvoltând Taylor cei doi logaritmi, găsim:

$$\begin{aligned} \ln(A_{n,k}) &= -(np + x_k \sqrt{npq}) \left(x_k \sqrt{\frac{q}{np}} - \frac{1}{2} \frac{x_k^2 q}{np} + O \left(\frac{1}{n^{3/2}} \right) \right) - \\ &\quad - (nq - x_k \sqrt{npq}) \left(-x_k \sqrt{\frac{p}{nq}} - \frac{1}{2} \frac{x_k^2 p}{nq} + O \left(\frac{1}{n^{3/2}} \right) \right) \end{aligned}$$

deoarece x_k sunt mărginiți de a și b . După ce facem înmulțirile și reducem termenii rezultă:

$$\ln(A_{n,k}) = -\frac{x_k^2}{2} + O \left(\frac{1}{\sqrt{n}} \right)$$

Deci

$$A_{n,k} = e^{-\frac{x_k^2}{2}} e^{O(1/\sqrt{n})} \rightarrow e^{-\frac{x_k^2}{2}}$$

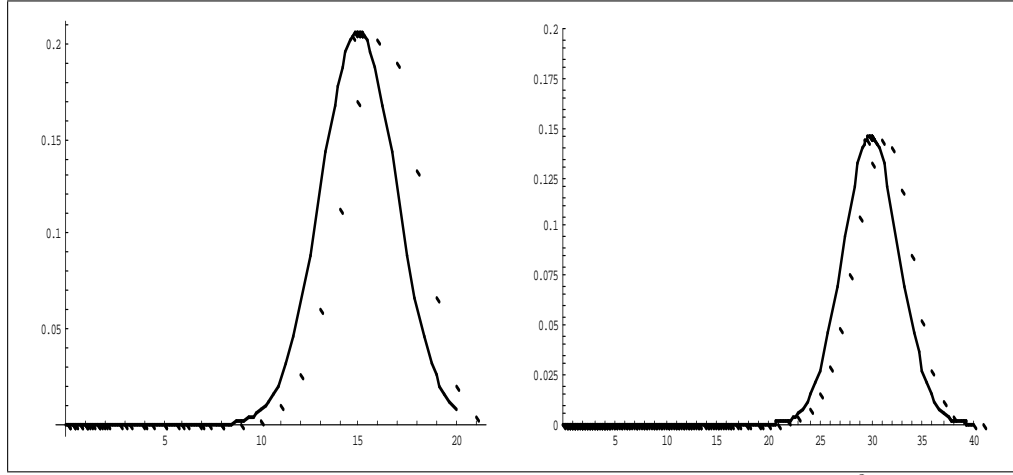
Folosind limitele de la i), ii), iii) găsim:

$$\sqrt{npq} C_n^k p^k q^{n-k} \rightarrow \frac{1}{\sqrt{2\pi}} e^{-\frac{x_k^2}{2}}$$

uniform pentru $a \leq x_k \leq b$.

QED.

În figura următoare apar pentru $n = 20$ și $n = 40$, $p = 3/4$, $q = 1/4$ atât valorile $C_n^k p^k q^{n-k}$ cât și graficele funcțiilor $\frac{1}{\sqrt{npq}} \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-np)^2}{2npq}}$ care pentru $x = x_k$ au ca valoare $\frac{1}{\sqrt{npq}} \frac{1}{\sqrt{2\pi}} e^{-\frac{x_k^2}{2}}$



Probabilitățile $C_n^k p^k q^{n-k}$ și graficele funcțiilor $\frac{1}{\sqrt{npq}} \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-np)^2}{2npq}}$ pentru $p=3/4$ și $n=20$ respectiv $n=40$

Teorema de Moivre-Laplace locală se mai poate scrie

$$\sqrt{npq} C_n^k p^k q^{n-k} = \frac{1}{\sqrt{2\pi}} e^{-\frac{x_k^2}{2}} (1 + \epsilon_{n,k})$$

cu $\epsilon_{n,k} \rightarrow 0$ când $n \rightarrow \infty$, uniform pentru k astfel încât $-\infty < a \leq x_k \leq b < \infty$.

Forma sub care se utilizează cel mai frecvent acest rezultat este:

Teorema 5.9 (formula integrală de Moivre-Laplace) Fie $0 < p < 1, q = 1 - p$. Atunci:

$$\lim_{n \rightarrow \infty} \sum_{a \leq \frac{k-np}{\sqrt{npq}} \leq b} C_n^k p^k q^{n-k} = \frac{1}{\sqrt{2\pi}} \int_a^b e^{-\frac{x^2}{2}} dx \quad (5.12)$$

limita fiind uniformă în $-\infty \leq a \leq b \leq \infty$.

Demonstrație(schită) Considerăm cazul a și b finite. Fie $x_k = \frac{k-np}{\sqrt{npq}}$. Suma din teoremă se scrie:

$$\sum_{a \leq x_k \leq b} \sqrt{npq} C_n^k p^k q^{n-k} \frac{1}{\sqrt{npq}} \quad (5.13)$$

Dar conform formulei locale $\sqrt{npq} C_n^k p^k q^{n-k} = \frac{1}{\sqrt{2\pi}} e^{-\frac{x_k^2}{2}} (1 + \epsilon_{n,k})$ cu $|\epsilon_{n,k}| \leq \epsilon_n \rightarrow 0$ independent de k , dacă $a \leq x_k \leq b$ și $x_k - x_{k-1} = \frac{1}{\sqrt{npq}}$. Deci

$$\begin{aligned} \sqrt{npq} C_n^k p^k q^{n-k} &= \frac{1}{\sqrt{2\pi}} e^{-\frac{x_k^2}{2}} + \epsilon_{n,k} e^{-\frac{x_k^2}{2}} \\ &= \frac{1}{\sqrt{2\pi}} e^{-\frac{x_k^2}{2}} + \epsilon'_{n,k} \end{aligned}$$

Deci suma (5.13) se scrie:

$$\sum_{a \leq x_k \leq b} \frac{1}{\sqrt{2\pi}} e^{-\frac{x_k^2}{2}} (x_k - x_{k-1}) + \sum_{a \leq x_k \leq b} \epsilon_{n,k} e^{-\frac{x_k^2}{2}} (x_k - x_{k-1})$$

Prima sumă de mai sus tinde când $n \rightarrow \infty$ către $\int_a^b \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx = \Phi(b) - \Phi(a)$ iar a doua e majorată de $\epsilon_n (b - a)$ deci tinde la zero când $n \rightarrow \infty$. Prin urmare teorema e demonstrată pentru a și b finite. O examinare mai atentă a situației arată că teorema este adevărată și pentru a sau b infinite, iar limita este uniformă în a și b.

QED.

Care sunt valorile n pentru care aproximarea probabilităților $C_n^k p^k q^{n-k}$ prin formulele de Moivre-Laplace este suficient de bună? În general se consideră că:

i) $n \geq 30$ $p \approx \frac{1}{2}$, atunci formulele de Moivre-Laplace dau o aproximație satisfăcătoare.

ii) $n \geq 30$, $p \leq \frac{1}{10}$ $np = \lambda \leq 10$ atunci formulele distribuției Poisson dau o aproximație suficient de bună: $C_n^k p^k q^{n-k} \approx e^{-\lambda} \frac{\lambda^k}{k!}$.

iii) $n \geq 30$ $p \leq \frac{1}{10}$ $np \geq 10$ atunci sunt satisfăcătoare formulele de Moivre-Laplace, adică $C_n^k p^k q^{n-k} \simeq \frac{1}{\sqrt{npq}} \frac{1}{\sqrt{2\pi}} e^{-\frac{x_k^2}{2}}$.

Exemplul 5.10 O firmă de asigurări a distribuit polițe de asigurare la 10.000 de persoane de aceeași vârstă și grup social contra unei sume de 20.000 lei și în caz de accident firma îi plătește persoanei 2.000.000 lei. Probabilitatea de accident pentru acest grup de persoane este $p=6/1000$.

a) Care este probabilitatea ca firma să dea faliment?

b) Care e probabilitatea ca firma să câștige cel puțin 50 milioane lei?

Soluție. a) Firma încasează 200 milioane lei. Probabilitatea de a da faliment este probabilitatea ca să plătească mai mult de 200 milioane lei, adică numărul celor accidentați să depășească 100. Avem o schemă binomială cu $n=10.000$, $p=0,006$, $np=60$, $q=0,994$. Se cere suma probabilităților $C_n^k p^k q^{n-k}$ pentru $k \geq 100$. Conform formulei de Moivre-Laplace avem:

$$\begin{aligned} \sum_{k \geq 100} C_n^k p^k q^{n-k} &= \sum_{\frac{100-np}{\sqrt{npq}} \leq \frac{k-np}{\sqrt{npq}} \leq \frac{n-np}{\sqrt{npq}}} C_n^k p^k q^{n-k} \\ &\approx \int_{\frac{100-60}{\sqrt{10000 \cdot 0,006 \cdot 0,994}}}^{\frac{10000-60}{\sqrt{10000 \cdot 0,006 \cdot 0,994}}} \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx \\ &= \Phi(1287) - \Phi(5,17) \approx 0,5 - 0,5 = 0 \end{aligned}$$

Probabilitatea de a da faliment este aproape zero.

b) Firma câștigă cel puțin 50 milioane dacă plătește sub 150 milioane, deci dacă numărul celor accidentați este sub 75. Probabilitatea cerută este deci

$$\begin{aligned}
 p &= \sum_{k \leq 75} C_n^k p^k q^{n-k} = \sum_{\substack{\frac{k-np}{\sqrt{npq}} = \frac{75-np}{\sqrt{npq}} \\ \frac{k-np}{\sqrt{npq}} = \frac{0-np}{\sqrt{npq}}}} C_n^k p^k q^{n-k} \approx \\
 &\approx \int_{x=\frac{0-10000 \cdot 0,006}{\sqrt{10000 \cdot 0,006 \cdot 0,994}}}^{x=\frac{75-10000 \cdot 0,006}{\sqrt{10000 \cdot 0,006 \cdot 0,994}}} \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx = \int_{-5,485}^{1,371} \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx \\
 &= \Phi(1,942) - \Phi(-7,769) = \Phi(1,942) + \Phi(7,769) \\
 &= 0,474 + 0,5 = 0,974
 \end{aligned}$$

Prin urmare câștigul va fi de peste 50 milioane cu o probabilitate foarte mare.

Alte exemple de aplicare a formulelor de Moivre-Laplace sunt date în exerciții.

Teoremele de Moivre-Laplace fac parte dintr-o clasă mai largă de teoreme cunoscute sub numele de teoreme limită centrală. Știm că în anumite condiții media unui șir de variabile aleatoare independente tinde spre o constantă (legea numerelor mari). Cum tinde această medie de variabile aleatoare spre o constantă? Fie șirul de v.a independente f_n dat de (5.6). Atunci $\mu_n = f_1 + f_2 + \dots + f_n$ are o distribuție Bernoulli, deci $p(\mu_n = k) = C_n^k p^k q^{n-k}$. Pe de o parte știm din legea numerelor mari că $\frac{f_1+f_2+\dots+f_n}{n} \rightarrow p$. Pe de altă parte

$$\begin{aligned}
 &p \left(a \leq \frac{(f_1 - M(f_1)) + (f_2 - M(f_2)) + \dots + (f_n - M(f_n))}{\sqrt{\sum_{k=1}^n D(f_k)}} < b \right) \\
 &= p \left(a \leq \frac{n}{\sqrt{\sum_{k=1}^n D(f_k)}} \left(\frac{f_1+\dots+f_n}{n} - \frac{M(f_1)+\dots+M(f_n)}{n} \right) < b \right) \\
 &= p \left(a \leq \frac{f_1 + f_2 + \dots + f_n - np}{\sqrt{\sum_{k=1}^n pq}} < b \right) \\
 &= p \left(a \leq \sqrt{\frac{n}{pq}} \left(\frac{f_1 + f_2 + \dots + f_n}{n} - p \right) < b \right) \\
 &= p \left(a \leq \frac{\mu_n - np}{\sqrt{npq}} < b \right) = \sum_{a \leq \frac{k-np}{\sqrt{npq}} \leq b} p(\mu_n = k) \\
 &= \sum_{a \leq \frac{k-np}{\sqrt{npq}} \leq b} C_n^k p^k q^{n-k} \rightarrow \frac{1}{\sqrt{2\pi}} \int_a^b e^{-\frac{x^2}{2}} dx. \quad (\text{cf. 5.12})
 \end{aligned}$$

Prin urmare variabila aleatoare $\frac{f_1+f_2+\dots+f_n}{n} - p$ se comportă, pentru valori mari ale lui n ca o variabilă aleatoare normală, de tip $N(0, 1)$ multiplicată cu $\sqrt{\frac{pq}{n}}$. Acest fenomen este mai general.

Anume, relația:

$$\begin{aligned}
 & p \left(a \leq \frac{(f_1 - M(f_1)) + (f_2 - M(f_2)) + \dots + (f_n - M(f_n))}{\sqrt{\sum_{k=1}^n D(f_k)}} < b \right) \\
 &= p \left(a \leq \frac{n}{\sqrt{\sum_{k=1}^n D(f_k)}} \left(\frac{f_1 + \dots + f_n}{n} - \frac{M(f_1) + \dots + M(f_n)}{n} \right) < b \right) \\
 &\xrightarrow{n \rightarrow \infty} \frac{1}{\sqrt{2\pi}} \int_a^b e^{-\frac{x^2}{2}} dx
 \end{aligned} \tag{5.14}$$

are loc dacă:

Varianta I: $(f_k)_{k \in N}$ este un șir de v.a. pe același câmp X , independente două câte două, cu aceeași funcție de repartiție, cu dispersie finită și nenulă.

Varianta II: $(f_k)_{k \in N}$ este un șir de v.a. pe același câmp X , independente două câte două și

$$\lim_{n \rightarrow \infty} \frac{\sqrt[3]{M(|f_1 - M(f_1)|^3) + \dots + M(|f_n - M(f_n)|^3)}}{\sqrt{\sum_{k=1}^n D(f_k)}} = 0$$

Mai există și alte variante de condiții necesare pentru ca relația (5.14) să aibă loc. Asemenea teoreme poartă numele de teoreme limită centrală și pun în evidență că în condiții foarte generale medierea unui număr mare de v.a. independente conduce la o constantă plus o variabilă aleatoare normală. Forma exactă este relația (5.14).

5.3 Rezumat

Prima carte de probabilități a fost publicată în 1704 de J. Bernoulli unde a fost demonstrată legea numerelor mari sub forma (5.5) și unde afirmă că s-a gândit 20 de ani la această teoremă.. Ulterior au apărut generalizări ale ei sub diverse forme, ca de exemplu (5.2). O demonstrație elegantă a acestor teoreme se bazează pe inegalitatea lui Cebâșev (5.1).

În esență legea numerelor mari spune că în condiții generale media unui număr mare de v.a. independente diferă puțin de o constantă. Acest lucru este exprimat precis prin formula (5.3).

Teoremele de tip limită centrală aduc o precizare legii numerelor mari. Anume, diferența dintre media unui șir de v.a. independente și constanta la care converge această medie este aproximativ o v.a. normală, de tip $N(0, 1)$ înmulțită cu o constantă care la rândul ei tinde la zero. Forma matematică a acestui enunț este (5.14).

În aplicații apare frecvent schema Bernoulli și deci calculul probabilităților $C_n^k p^k q^{n-k}$. Formulele de Moivre-Laplace (locală 5.10 și globală 5.12) ne permit să calculăm aproximativ aceste probabilități cu ajutorul densității normale $\frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$ (cazul local) sau primitivei ei,

funcția Φ (cazul global). În unele cazuri $C_n^k p^k q^{n-k}$ se pot aproxima prin probabilitățile legii Poisson. Limitele în care aproximațiile probabilităților $C_n^k p^k q^{n-k}$ prin legea normală sau Poisson sunt considerate satisfăcătoare sunt date în pag. ??.

Formula lui Stirling este frecvent folosită pentru evaluarea factorialelor. Există și o generalizare a ei pentru evaluarea funcției Gamma.

FORMULE UTILIZATE FRECVENT:

a) Inegalitatea lui Cebâșev: $p(|f - M(f)| \leq \epsilon) \geq 1 - \frac{D(f)}{\epsilon^2}$.

b) Legea numerelor mari: $\lim_{n \rightarrow \infty} p\left(\left|\frac{f_1 + f_2 + \dots + f_n}{n} - \frac{M(f_1) + \dots + M(f_n)}{n}\right| \leq \epsilon\right) \geq 1 - \frac{C}{n\epsilon^2}$ dacă toate dispersiile sunt mărginite de C; varianta Bernoulli: $\lim_{n \rightarrow \infty} (|\frac{\mu_n}{n} - p| \leq \epsilon) = 1$ unde μ_n este numărul de apariții ale unui eveniment A în n experiențe independente iar p este probabilitatea de apariție a lui A într-un singur experiment.

c) Formula locală de Moivre-Laplace: $\lim_{n \rightarrow \infty} \frac{\sqrt{npq} C_n^k p^k q^{n-k}}{\frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}} = 1$.

d) Formula globală de Moivre-Laplace: $\lim_{n \rightarrow \infty} \sum_{a \leq \frac{k-np}{\sqrt{npq}} \leq b} C_n^k p^k q^{n-k} = \frac{1}{\sqrt{2\pi}} \int_a^b e^{-\frac{x^2}{2}} dx$, $0 < p < 1$, $q = 1 - p$.

5.4 Exerciții

1. Se aruncă o monedă de 10000 ori. Care e probabilitatea ca banul să apară de mai puțin 4550 ori?

Indicație. Se procedează ca în exemplul din lecție.

2. Probabilitatea de a lovi o țintă este $p = 1/1000$. Care e probabilitatea ca din 5000 de lovituri, ținta să fie atinsă de cel puțin 2 ori?

Indicație. Se utilizează de Moivre-Laplace.

3. Probabilitatea ca un anumit produs să fie defect este $p = 0,005$. Care e probabilitatea ca dintr-un lot de 10000 produse luate la întâmplare să fie mai puțin de 70 defecte?

Indicație. Se utilizează formula de Moivre-Laplace.

4. O firmă de asigurări are m salariați cu un salariu mediu de s lei pe an. O persoană asigurată cu a lei primește o despăgubire de d lei dacă pățește ceva. Probabilitatea de a păți ceva este p . a) Care este probabilitatea ca veniturile companiei să fie în acel an mai mari sau egale cu c_0 . b) Care e numărul minim de asigurați de la care firma înregistrează un venit brut mai mare ca c_0 cu probabilitatea de cel puțin 0,99?

Indicație. Se utilizează formula de Moivre-Laplace.

5. Să se demonstreze că pentru $x > 0$ funcția $\phi(x) = \int_x^\infty e^{-\frac{t^2}{2}} dt$, satisface inegalitățile:

$$\frac{x}{1+x^2} e^{-\frac{x^2}{2}} \leq \phi(x) \leq \frac{1}{x} e^{-\frac{x^2}{2}}$$

Indicație. Pentru prima inegalitate se consideră $f(x) = \frac{x}{1+x^2} e^{-\frac{x^2}{2}} - \Phi(x)$. Avem $f(0) = 0$. Se derivează și se studiază variația lui f . Analog pentru a doua inegalitate.

6. Două vase identice conțin fiecare 10^{22} molecule. Vasele sunt puse în contact și moleculele se mișcă liber între vase. După omogenizare, care e probabilitatea ca în vasul A să fie cu a zece miliarde parte mai multe molecule ca în B ? (admitem că pentru fiecare moleculă probabilitatea de a fi în A este egală cu probabilitatea de a fi în B și este $1/2$).

Indicație. Se utilizează formula de Moivre-Laplace și ex. precedent.

7. Pe scara unui bloc sunt 40 apartamente și în fiecare apartament este câte un frigider de putere 0,6 kW. În medie un frigider consumă energie electrică 3 ore din 24. Care este probabilitatea ca la un moment dat să fie absorbită de la rețea, de ctre frigidere, o putere mai mare de 4 kW?

Indicație. Se utilizează de Moivre-Laplace.

8. a) Să se arate că dacă $\epsilon > 0$ și $\delta > 0$ sunt date, iar $n \geq \frac{1}{4\delta\epsilon^2}$ atunci:

$$\sum_{\left|\frac{k}{n}-x\right|>\epsilon} C_n^k x^k (1-x)^{n-k} \leq \delta$$

unde $x \in [0, 1]$.

b) Să se arate că dacă $f : [0, 1] \rightarrow \mathbb{R}$ este continuă, deci mărginită $|f| \leq M$ atunci ϵ, δ și n fiind ca la pct. a) avem:

$$-2M\delta \leq \sum_{\left|\frac{k}{n}-x\right|>\epsilon} C_n^k x^k (1-x)^{n-k} \left(f\left(\frac{k}{n}\right) - f(x) \right) \leq 2M\delta$$

c) Să se arate că în condițiile de la pct. b)

$$-\omega_{x,\epsilon} \leq \left(\sum_{\left|\frac{k}{n}-x\right|\leq\epsilon} C_n^k x^k (1-x)^{n-k} \left(f\left(\frac{k}{n}\right) - f(x) \right) \right) \leq \omega_{x,\epsilon}$$

unde $\omega_{x,\epsilon}$ este oscilația funcției f pe intervalul $[x - \epsilon, x + \epsilon]$, adică diferența între minimum și maximumul funcției pe acest interval.

d) Folosind a), b) și c) să se arate că șirul de polinoame

$$B_n(x) = \sum_{k=0..n} C_n^k x^k (1-x)^{n-k} f\left(\frac{k}{n}\right)$$

converge uniform la f pe $[0, 1]$ oricare ar fi funcția continuă f . (polinoamele lui Bernstein).

Indicație. a) Se utilizează inegalitatea lui Cebășev pentru schema Bernoulli (vezi obs. 5.9).

b) Se ține seama de și de $|f| \leq M$. c) Se ține seama de $\left|f\left(\frac{k}{n}\right) - f(x)\right| \leq \omega_{x,\epsilon}$ și $\sum_{k=0..n} C_n^k x^k (1-x)^{n-k} = 1$.

d) Se pune diferența $f(x) - B_n(x)$ ca în b) și c) apoi se ține seama că f e uniform continuă.

9. Să presupunem că f este o variabilă aleatoare nenegativă, cu media m . Să se arate că este adevărată inegalitatea

$$p(f \geq c) \leq \frac{m}{c}$$

(inegalitatea lui Markov).

10. Fie ν_r momentul centrat absolut de ordin r al unei variabile aleatoare f care are media m , definit prin $\nu_r = M(|f - m|^r)$. Să se arate că

$$p(|f - m| \geq \varepsilon) \leq \frac{\nu_r}{\varepsilon^r}$$

11. Vrem să verificăm dacă un zar este falsificat sau nu, căutând care este probabilitatea r de apariție a feței șase. Pentru aceasta alegem $\varepsilon = 0,01$, aruncăm zarul de un număr n de ori și observăm care este numărul k de apariții ale feței șase.

Să se determine o valoare cât mai mică pentru n astfel ca probabilitatea $p\left(\left|\frac{k}{n} - r\right| \leq \varepsilon\right)$ să fie cel puțin 0,99 utilizând:

a) inegalitatea lui Cebășev

b) formula integrală de Moivre-Laplace.

12. În problema acului lui Buffon (vezi lecția 3, exerciții), probabilitatea de intersecare a paralelor este $p = \frac{2l}{\pi a}$ unde l este lungimea acului iar a este distanța dintre paralele. De câte aruncări este nevoie ca probabilitatea ca $\left|\pi - \frac{2 \cdot l \cdot n}{a \cdot k}\right| < 0,01$ să fie cel puțin 0,99. (n este numărul de aruncări iar k este numărul de intersecări ale rețelei de drepte paralele).

Lecția 6

Dependența între variabilele aleatoare

6.1 Coeficientul de corelație

Definiția 6.1 Fie $f_1, f_2 : X \rightarrow R$ două v.a. cu mediile m_1, m_2 și dispersiile $D_1 = \sigma_1^2 = M((f_1 - m_1)^2)$, $D_2 = \sigma_2^2 = M((f_2 - m_2)^2)$. Se numește coeficient de corelație al variabilelor aleatoare f_1, f_2 numărul:

$$c(f_1, f_2) = \frac{M((f_1 - m_1) \cdot (f_2 - m_2))}{\sqrt{D_1 \cdot D_2}} \quad (6.1)$$

Mărimea $cov(f_1, f_2) = M((f_1 - m_1) \cdot (f_2 - m_2))$ se numește covarianța lui f_1 și f_2 .

Spunem că două funcții f_1 , și f_2 coincid aproape peste tot (a.p.t.) dacă $p(f_1 \neq f_2) = 0$. Importanța acestui coeficient se va vedea în continuare.

Teorema 6.2 (Inegalitatea Cauchy-Buniakovski). Fie $f_1, f_2 : X \rightarrow R$ două v.a. cu momente de ordinul doi finite. Atunci:

$$M^2(f_1 f_2) \leq M(f_1^2) \cdot M(f_2^2) \quad (6.2)$$

Egalitatea poate avea loc dacă și numai dacă $f_2 = a f_1$ sau $f_1 = a f_2$ a.p.t. cu $a = \text{const.}$

Demonstrație. $M((f_1 + x f_2)^2) \geq 0$ pentru orice $x \in R$, deci $M(f_1^2) + 2M(f_1 f_2)x + M(f_2^2)x^2 \geq 0$ pentru orice x . De aici rezultă că discriminantul funcției de gradul doi este negativ: $\Delta = M^2(f_1 f_2) - M(f_1^2) M(f_2^2) \leq 0$ adică tocmai inegalitatea Cauchy-Buniakovski. Dacă $\Delta = 0$ atunci există un x_0 real, soluție a ecuației de gradul doi

$$M((f_1 + x_0 f_2)^2) = 0$$

de unde rezultă că $(f_1 + x_0 f_2)^2$ ia valoarea zero pe o mulțime cu probabilitatea 1, adică $f_1 = -x_0 f_2$ a.p.t.

În cazul în care inecuația nu ar fi de gradul 2, deci $M(f_2^2) = 0$, rezultă $f_2 = 0$ a.p.t., deci $f_1 f_2 = 0$ a.p.t., deci (6.2) are din nou loc deoarece ambii membri sunt 0. În acest caz $f_2 = 0 \cdot f_1$.

Teorema 6.3 Fie $f_1, f_2 : X \rightarrow R$ două v.a. cu mediile m_1, m_2 și dispersiile $D_1 = \sigma_1^2 = M((f_1 - m_1)^2)$, $D_2 = \sigma_2^2 = M((f_2 - m_2)^2)$ finite. Atunci:

- a) $c(f_1, f_2) \in [-1, 1]$.
- b) $c(f_1, f_2) = \pm 1$ este echivalent cu $f_1 = af_2 + b$ sau $f_2 = af_1 + b$ a.p.t. pentru niște constante a și b potrivit alese.
- c) f_1 și f_2 independente implică $c(f_1, f_2) = 0$ (dar nu și invers).

Demonstrație. a) Se aplică inegalitatea Cauchy-Buniakovski pentru $f_1 - m_1$ și $f_2 - m_2$.
b) Dacă $c(f_1, f_2) = \pm 1$ atunci în inegalitatea Cauchy-Buniakovski semnul \leq devine $=$ și teorema precedentă implică $(f_2 - m_2) = a(f_1 - m_1)$ deci $f_2 = af_1 + m_2 - a \cdot m_1$ sau $(f_1 - m_1) = a \cdot (f_2 - m_2)$ de unde rezultă $f_1 = af_2 + b$ pentru a și b potrivit alese.

c) f_1 și f_2 independente implică $(f_1 - m_1)$ și $(f_2 - m_2)$ independente deci:

$$M((f_1 - m_1)(f_2 - m_2)) = M(f_1 - m_1)M(f_2 - m_2) = 0 \cdot 0 = 0$$

Prin urmare $c(f_1, f_2) = 0$.

QED.

Exemplul 6.4 Fie două șiruri finite de numere $(x_i)_{1 \leq i \leq n}$ și $(y_i)_{1 \leq i \leq n}$. Să considerăm o partiție a unei mulțimi X astfel: $X = A_1 \cup A_2 \cup A_3 \cup \dots \cup A_n$ cu toate mulțimile $A_i \neq \emptyset$. Definim pe X o probabilitate prin $p(A_i) = \frac{1}{n}$ și fie variabilele aleatoare $f_1, f_2 : X \rightarrow R$, definite prin $f_1(\omega) = x_i$ dacă $\omega \in A_i$ și $f_2(\omega) = y_i$ dacă $\omega \in A_i$. Avem :

$$\begin{aligned} m_1 &= M(f_1) = \frac{\sum_{i=1}^n x_i}{n}; & m_2 &= M(f_2) = \frac{\sum_{i=1}^n y_i}{n} \\ D_1 &= D(f_1) = \frac{\sum_{i=1}^n (x_i - m_1)^2}{n}; & D_2 &= D(f_2) = \frac{\sum_{i=1}^n (y_i - m_2)^2}{n} \\ m_{1,2} &= M((f_1 - m_1)(f_2 - m_2)) = \frac{\sum_{i=1}^n (x_i - m_1)(y_i - m_2)}{n} \\ c(f_1, f_2) &= \frac{m_{1,2}}{\sqrt{D_1 D_2}} = \frac{\sum_{i=1}^n (x_i - m_1)(y_i - m_2)}{\sqrt{\sum_{i=1}^n (x_i - m_1)^2} \sqrt{\sum_{i=1}^n (y_i - m_2)^2}} \end{aligned} \quad (6.3)$$

Conform teoremei precedente, egalitatea cu ± 1 a coeficientului 6.3 este echivalentă cu $f_2 = af_1 + b$ sau cu alte cuvinte $y_i = ax_i + b$ pentru toate valorile $1 \leq i \leq n$. Dacă $|c(f_1, f_2)|$ este suficient de aproape de unu atunci punctele (x_i, y_i) sunt aproximativ pe o dreaptă. Dreapta ce aproximează "cel mai bine" acest nor de puncte se numește dreapta de regresie a lui y în x și se determină prin metoda celor mai mici pătrate. Această dreaptă de

forma $y = ax + b$ se determină din condiția ca $\sum (y_i - ax_i - b)^2$ să fie minimă (a se vedea lecția 8 sau lecția 14). Dacă se caută dependențe de forma $y_i = be^{ax_i}$, prin logaritmare se găsește $\ln(y_i) = ax_i + \ln(b)$. O asemenea dependență există, conform celor de mai sus dacă $c(f_1, \ln(f_2)) = \pm 1$. Dacă $|c(f_1, \ln(f_2))|$ este suficient de aproape de unu atunci punctele $(x_i, \ln(y_i))$ sunt aproximativ pe o dreaptă care se determină prin metoda celor mai mici pătrate. Găsim astfel a și $\ln(b)$ care ne dau dependența $\ln(y_i) = ax_i + \ln(b)$, adică $y_i = be^{ax_i}$.

6.2 Variabile aleatoare bidimensionale

Fie (X, Ω, p) un σ câmp de probabilitate și fie $f_1, f_2 : X \rightarrow R$ două variabile aleatoare. Legătura între variabilele aleatoare este pusă în evidență prin funcția lor comună de repartiție. Să considerăm $f : X \rightarrow R^2$, definită prin $f(\omega) = (f_1(\omega), f_2(\omega)) \in R^2$. Funcția f are calitatea că pentru orice produs de intervale $I_1 \times I_2 \subset R^2$ mulțimea $f^{-1}(I_1 \times I_2) = f_1^{-1}(I_1) \cap f_2^{-1}(I_2) \in \Omega$.

Definiția 6.5 Se numește v.a. bidimensională pe câmpul de probabilitate (X, Ω, p) o funcție $f : X \rightarrow R^2$ cu proprietatea că pentru orice intervale I_1 și I_2 din R , rezultă $f^{-1}(I_1 \times I_2) \in \Omega$.

Vom mai spune uneori variabilă aleatoare dublă în loc de variabilă aleatoare bidimensională. Notând $f_1(\omega)$ și $f_2(\omega)$ cele două componente ale lui f , vedem că o v.a. bidimensională este asociată unei perechi de v.a. unidimensionale.

Definiția 6.6 Se numește funcție de repartiție bidimensională a lui f , funcția $F : R^2 \rightarrow R$ definită prin:

$$\begin{aligned} F(x_1, x_2) &= p(f^{-1}((-\infty, x_1) \times (-\infty, x_2))) \\ &= p((f_1 < x_1) \cap (f_2 < x_2)) \end{aligned}$$

Teorema 6.7 1) Funcția de repartiție bidimensională are proprietățile:

- a) F este crescătoare în fiecare argument.
- b) $F(x_1, -\infty) = 0, F(-\infty, x_2) = 0, F(\infty, \infty) = 1$.
- c) F este continuă la stânga în fiecare argument.
- d) Pentru orice $a_1 < b_1$ și $a_2 < b_2$ avem

$$\begin{aligned} & p((a_1 \leq f_1 < b_1) \cap (a_2 \leq f_2 < b_2)) \\ &= F(b_1, b_2) - F(a_1, b_2) - F(b_1, a_2) + F(a_1, a_2) \geq 0 \end{aligned} \quad (6.4)$$

2) Reciproc, orice funcție $F : R^2 \rightarrow R$ cu proprietățile a)-d) de mai sus este funcție de repartiție pentru o v.a. bidimensională potrivit aleasă.

Demonstrație.a) $x_1 < x'_1$ implică $\{f_1 < x_1 \cap f_2 < x_2\} \subset \{f_1 < x'_1 \cap f_2 < x_2\}$. Luând probabilitățile ambilor membri găsim $F(x_1, x_2) \leq F(x'_1, x_2)$. Analog se arată monotonia în x_2 .

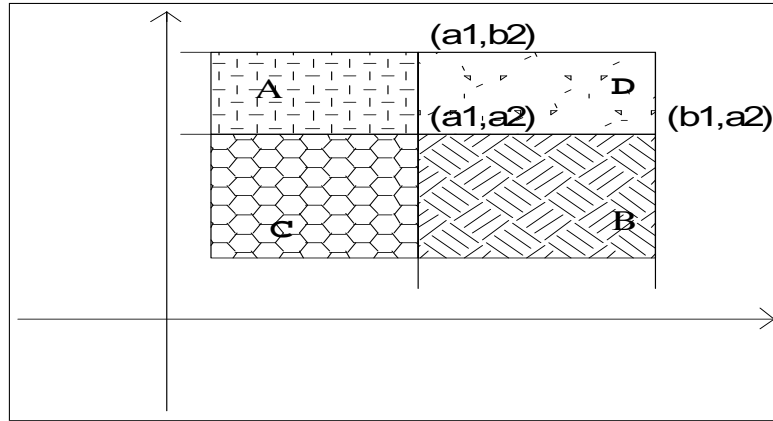
b) Prin definiție $F(x_1, -\infty) = \lim_{x_2 \rightarrow -\infty} F(x_1, x_2)$. Însă pentru orice șir monoton $x_{2,n} \rightarrow -\infty$ avem $\{f_1 < x_1 \cap f_2 < x_{2,n+1}\} \subset \{f_1 < x_1 \cap f_2 < x_{2,n}\}$ și $\bigcap_{n=1, \infty} \{f_1 < x_1 \cap f_2 < x_{2,n}\} = \emptyset$, deci $F(x_1, x_{2,n}) \rightarrow 0$, deci $F(x_1, -\infty) = 0$. Analog se arată și celelalte egalități.

c) Continuitatea la stânga se arată ca pentru v.a. unidimensionale (lecția 3).

d) Fie regiunile din $x_1 O x_2$ ca în figura următoare:

$A = \{(x_1, x_2) \mid -\infty < x_1 < a_1, a_2 \leq x_2 < b_2\}$, etc. (vezi figura). Fie evenimentele:

$A' = \{\omega \in X \mid (f_1(\omega), f_2(\omega)) \in A\}$ și analog B', C', D' . Observăm că $D' = (a_1 \leq f_1 < b_1) \cap (a_2 \leq f_2 < b_2)$. Avem $p(A') = F(a_1, b_2) - F(a_1, a_2)$, $p(C') = F(a_1, a_2)$, $p(B') = F(b_1, a_2) - F(a_1, a_2)$.



Scriem acum

$$\begin{aligned} F(b_1, b_2) &= p(A' \cup B' \cup C' \cup D') = p(A') + p(B') + p(C') + p(D') \\ &= F(a_1, b_2) - F(a_1, a_2) + F(b_1, a_2) - F(a_1, a_2) + F(a_1, a_2) + p(D') \end{aligned}$$

de unde rezultă pentru $p(D')$ formula 6.4.

QED.

Definiția 6.8 Se numește densitate de probabilitate a v.a. $f : X \rightarrow R^2$ o funcție $\rho : R^2 \rightarrow R$ pozitivă, integrabilă, astfel ca $\int_{-\infty}^{x_1} \int_{-\infty}^{x_2} \rho(t_1, t_2) dt_1 dt_2 = F(x_1, x_2)$ pentru $(x_1, x_2) \in R^2$ F fiind funcția de repartiție a lui f .

Vom mai numi ρ densitatea mixtă a variabilelor f_1 și f_2 unde $f = (f_1, f_2)$.

Teorema 6.9 1) Densitatea de probabilitate ρ are proprietățile:

a) $\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \rho(t_1, t_2) dt_1 dt_2 = 1$.

$$b) p((a_1 \leq f_1 < b_1) \cap (a_2 \leq f_2 < b_2)) = \int_{a_1}^{b_1} \int_{a_2}^{b_2} \rho(x_1, x_2) dx_1 dx_2.$$

2) Invers, orice funcție $\rho : R^2 \rightarrow R$ pozitivă și integrabilă cu proprietatea a) de mai sus este densitate de probabilitate pentru o v.a. bidimensională potrivit aleasă.

Demonstrație. a), b) sunt evidente din definiție. Punctul 2) se poate demonstra luând $X = R^2$, probabilitatea unui paralelipiped $[a, b] \times [c, d]$ se definește prin $\int_a^b \int_c^d \rho(x_1, x_2) dx_1 dx_2$, iar cele două variabile aleatoare sunt $f_1, f_2 : R^2 \rightarrow R$, $f_1(x_1, x_2) = x_1$, $f_2(x_1, x_2) = x_2$.

QED.

Cunoscând F =funcția de repartiție și ρ =densitatea de probabilitate bidimensională a lui $f = (f_1, f_2)$ putem calcula relativ ușor diverse mărimi legate de f_1 și f_2 .

1. Funcția F_1 de repartiție a lui f_1 este

$$F_1(x_1) = p(f_1 < x_1) = F(x_1, \infty) = \int_{-\infty}^{x_1} dt_1 \int_{-\infty}^{\infty} \rho(t_1, x_2) dx_2 \quad (6.5)$$

iar densitatea este

$$\rho_1(x_1) = \int_{-\infty}^{\infty} \rho(x_1, x_2) dx_2. \quad (6.6)$$

2. Media lui f_1 este

$$m_1 = M(f_1) = \int_{-\infty}^{\infty} x_1 \rho_1(x_1) dx_1 = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x_1 \rho(x_1, x_2) dx_1 dx_2 \quad (6.7)$$

3. Dispersia este

$$D(f_1) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x_1 - m_1)^2 \rho(x_1, x_2) dx_1 dx_2. \quad (6.8)$$

4. Formule analoage sunt valabile și pentru f_2 . De exemplu

$$m_2 = M(f_2) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x_2 \rho(x_1, x_2) dx_1 dx_2$$

5. Dacă $g : R^2 \rightarrow R$ este o funcție continuă, atunci $g(f_1, f_2) : X \rightarrow R$ dată de $g(f_1, f_2)(\omega) = g(f_1(\omega), f_2(\omega))$ este o v.a. Media acestei v.a. se calculează cu densitatea mixtă prin

$$M(g) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x_1, x_2) \rho(x_1, x_2) dx_1 dx_2 \quad (6.9)$$

6. În particular luând $g = (f_1 - m_1)(f_2 - m_2)$ apoi $g = (f_1 - m_1)^2$ și $g = (f_2 - m_2)^2$ găsim pentru coeficientul de corelație:

$$\begin{aligned} c(f_1, f_2) &= \frac{M((f_1 - m_1)(f_2 - m_2))}{\sqrt{M((f_1 - m_1)^2)} \sqrt{M((f_2 - m_2)^2)}} \\ &= \frac{\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x_1 - m_1)(x_2 - m_2) \rho(x_1, x_2) dx_1 dx_2}{\sqrt{D(f_1)} \sqrt{D(f_2)}} \end{aligned} \quad (6.10)$$

7. Dacă f_1 și f_2 sunt independente atunci

$$F(x_1, x_2) = p(f < x_1 \cap f_2 < x_2) = p(f_1 < x_1) \cdot p(f_2 < x_2) = F_1(x_1) \cdot F_2(x_2) \quad (6.11)$$

De aici rezultă

$$\rho(x_1, x_2) = \rho_1(x_1) \cdot \rho_2(x_2). \quad (6.12)$$

Se vede ușor că și reciproc este adevărat: dacă funcția de repartiție ori densitatea de probabilitate se descompune în produs de două funcții, una depinzând de x_1 și cealaltă depinzând de x_2 atunci f_1 și f_2 sunt independente (exercițiu).

Dintre formulele de mai sus doar (6.9) mai necesită demonstrație, celelalte fiind evidente.

a) Considerăm cazul când f_1 și f_2 sunt mărginite, deci $f = (f_1, f_2) : X \rightarrow [a, b] \times [c, d]$. Divizăm intervalele $[a, b]$ și $[c, d]$ prin $\Delta_x = (x_i)_{0 \leq i \leq n}$ respectiv $\Delta_y = (y_j)_{0 \leq j \leq m}$ și notăm $\Delta = \Delta_x \times \Delta_y$, $|\Delta| = \max\{|\Delta_x|, |\Delta_y|\}$, $D_{ij} = [x_{i-1}, x_i] \times [y_{j-1}, y_j]$, $A_{ij} = f^{-1}(D_{ij}) \subset X$. Variabila aleatoare

$$g_\Delta(\omega) = \begin{cases} g(x_i, y_j) & \text{dacă } \omega \in A_{ij} \\ 0 & \text{în caz contrar} \end{cases}$$

este simplă, și tinde la $g(f_1, f_2)$ atunci când $|\Delta_x| \rightarrow 0$ și $|\Delta_y| \rightarrow 0$. Deci

$$\begin{aligned} M(g_\Delta) &= \sum g(x_i, y_j) p(A_{ij}) \\ &= \sum g(x_i, y_j) \underbrace{\int \int_{D_{ij}} \rho(x, y) dx dy}_{p(A_{ij})} \\ &= \sum g(x_i, y_j) \underbrace{\rho(\xi_i, \eta_j) (x_i - x_{i-1}) (y_j - y_{j-1})}_{\text{din teorema de medie}} \\ &\xrightarrow{|\Delta| \rightarrow 0} \int \int_{[a, b] \times [c, d]} g(x, y) \rho(x, y) dx dy \end{aligned}$$

b) Demonstrația se poate extinde acum și la cazul când f_1 și f_2 nu sunt mărginite dar integrala (6.9) este convergentă. Argumentele fiind lungi, însă standard, nu le reproducem aici.

Repartițiile pentru f_1 și f_2 se mai numesc repartiții marginale ale lui f .

Exemplul 6.10 Fie $D = [a, b] \times [c, d]$. O variabilă aleatoare $f = (f_1, f_2)$ cu valori în D care are densitatea:

$$\rho(x_1, x_2) = \begin{cases} 0 & (x_1, x_2) \notin D \\ \frac{1}{\text{aria}(D)} & (x_1, x_2) \in D \end{cases}$$

se numește uniformă. Conform celor de la pct. 6) rezultă că f_1 și f_2 sunt independente.

Exemplul 6.11 Fie o v.a. (f_1, f_2) bidimensională cu densitatea:

$$\rho(x_1, x_2) = \frac{\sqrt{\det(A)}}{2\pi} e^{-\frac{1}{2}Q(x_1, x_2)}$$

unde $A = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}$ este o matrice pozitiv definită iar

$$Q(x_1, x_2) = \sum_{i,j=1}^2 a_{i,j} (x_i - m_i) (x_j - m_j)$$

O astfel de variabilă aleatoare se numește normală. Se cere:

a) $\sigma_1^2 = D(f_1)$ și $\sigma_2^2 = D(f_2)$; b) $c(f_1, f_2)$.

Soluție. a) Avem nevoie de $M(f_1) = \int \int x_1 \rho(x_1, x_2) dx_1 dx_2$. Se știe că matricea A este simetrică. Pentru aceste integrale este utilă următoarea schimbare de coordonate:

i) Se determină λ_1 și λ_2 valorile proprii ale matricei A , din ecuația:

$$\begin{vmatrix} a_{11} - \lambda & a_{12} \\ a_{21} & a_{22} - \lambda \end{vmatrix} = 0$$

ii) Se determină vectorii proprii ai matricei A și se normalizează. Acești vectori puși pe două coloane dau o matrice ortogonală $R = \begin{pmatrix} r_{11} & r_{12} \\ r_{21} & r_{22} \end{pmatrix}$, $R^t = R^{-1}$.

iii) Se știe că $R^t A R = \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix}$. Facem schimbarea $\begin{pmatrix} x_1 - m_1 \\ x_2 - m_2 \end{pmatrix} = R \cdot \begin{pmatrix} x'_1 \\ x'_2 \end{pmatrix}$. Prin această schimbare Q devine $Q(x_1, x_2) = \lambda_1 x_1'^2 + \lambda_2 x_2'^2$ iar iacobianul $\frac{D(x_1, x_2)}{D(x'_1, x'_2)} = \det(R) = 1$. Avem de asemenea $\det(A) = \lambda_1 \lambda_2$. Toate aceste lucruri se studiază la cursul de algebră liniară. Utilizând $\frac{1}{\sqrt{2\pi\sigma}} \int_{-\infty}^{\infty} e^{-\frac{x^2}{2\sigma^2}} dx = 1$, $\frac{1}{\sqrt{2\pi\sigma}} \int_{-\infty}^{\infty} x^2 e^{-\frac{x^2}{2\sigma^2}} dx = \sigma^2$, $\frac{1}{\sqrt{2\pi\sigma}} \int_{-\infty}^{\infty} x e^{-\frac{x^2}{2\sigma^2}} dx = 0$, (a se vedea lecția 4, legea normală) găsim:

iv)

$$\begin{aligned} M(f_1) &= \frac{\sqrt{\lambda_1 \lambda_2}}{2\pi} \int_{-\infty}^{\infty} dx'_1 \int_{-\infty}^{\infty} (m_1 + r_{11}x'_1 + r_{12}x'_2) e^{-\frac{\lambda_1}{2}x_1'^2 - \frac{\lambda_2}{2}x_2'^2} dx'_2 \\ &= \frac{\sqrt{\lambda_1}}{\sqrt{2\pi}} \cdot \frac{\sqrt{\lambda_2}}{\sqrt{2\pi}} \cdot m_1 \cdot \int_{-\infty}^{\infty} e^{-\frac{\lambda_1}{2}x_1'^2} dx'_1 \cdot \int_{-\infty}^{\infty} e^{-\frac{\lambda_2}{2}x_2'^2} dx'_2 = m_1 \end{aligned}$$

Analog găsim $M(f_2) = m_2$.

Putem acum calcula dispersiile.

a)

$$\begin{aligned}
\sigma_1^2 &= \int \int (x_1 - m_1)^2 \rho(x_1, x_2) dx_1 dx_2 \\
&= \frac{\sqrt{\lambda_1 \lambda_2}}{2\pi} \int \int (r_{11}x'_1 + r_{12}x'_2)^2 e^{\frac{-\lambda_1 x'^2_1 - \lambda_2 x'^2_2}{2}} dx'_1 dx'_2 \\
&= r_{11}^2 \cdot \frac{1}{\lambda_1} + r_{12}^2 \cdot \frac{1}{\lambda_2}
\end{aligned}$$

In mod analog găsim $\sigma_2^2 = r_{21}^2 \cdot \frac{1}{\lambda_1} + r_{22}^2 \cdot \frac{1}{\lambda_2}$

b) Covarianța lui f_1 și f_2 este:

$$\begin{aligned}
cov(f_1, f_2) &= \frac{\sqrt{\lambda_1 \lambda_2}}{2\pi} \int \int (r_{11}x'_1 + r_{12}x'_2)(r_{21}x'_1 + r_{22}x'_2) e^{\frac{-\lambda_1 x'^2_1 - \lambda_2 x'^2_2}{2}} dx'_1 dx'_2 \\
&= \frac{\sqrt{\lambda_1 \lambda_2}}{2\pi} \int \int (r_{11}r_{21}x'^2_1 + r_{12}r_{22}x'^2_2) e^{\frac{-\lambda_1 x'^2_1 - \lambda_2 x'^2_2}{2}} dx'_1 dx'_2 \\
&= r_{11}r_{21} \cdot \frac{1}{\lambda_1} + r_{12}r_{22} \cdot \frac{1}{\lambda_2}
\end{aligned}$$

. De aici rezultă $c(f_1, f_2)$ și punctul b) e terminat. Mai remarcăm că deoarece $R^t A R = \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix}$, rezultă

$$A^{-1} = R \begin{pmatrix} \frac{1}{\lambda_1} & 0 \\ 0 & \frac{1}{\lambda_2} \end{pmatrix} R^t = \begin{pmatrix} \sigma_1^2 & cov(f_1, f_2) \\ cov(f_1, f_2) & \sigma_2^2 \end{pmatrix}$$

6.3 Funcția de repartiție condiționată

In această secțiune (X, Ω, p) este un câmp de probabilitate, $f : X \rightarrow R$ o v.a. cu funcția de repartiție F , iar α, β sunt numere reale pozitive. Dacă există următoarea limită

$$\begin{aligned}
&\lim_{\alpha, \beta \rightarrow 0} \frac{p(B \cap (x - \alpha \leq f < x + \beta))}{p(x - \alpha \leq f < x + \beta)} \\
&= \lim_{\alpha, \beta \rightarrow 0} \frac{p(B \cap (x - \alpha \leq f < x + \beta))}{F(x + \beta) - F(x - \alpha)}
\end{aligned}$$

atunci o notăm $p(B|f = x)$. Putem da acum definiția:

Definiția 6.12 Dacă $x \in R$ și $B \in \Omega$, numim probabilitatea lui B condiționată de $f = x$ mărimea $p(B|f = x)$.

E posibil ca $p(f = x) = 0$ și deci să nu putem defini $p(B|f = x) = \frac{p(B \cap \{f=x\})}{p(f=x)}$ așa ca în lecția 1. În asemenea situație aplicăm trecerea la limită de mai sus. Ca funcție de B , $p(\cdot|f = x)$ este o probabilitate finit aditivă pe o subalgebră a lui Ω . Dacă nu există pericol de confuzie vom scrie $p(B|x)$ în loc de $p(B|f = x)$.

Fie $g : X \rightarrow R$ este o v.a. iar $h : X \rightarrow R^2$ $h = (f, g)$ o v.a. dublă care are funcția de repartiție H . Vom nota cu ρ_f densitățile care se referă la f , cu ρ_g densitățile care se referă la g , cu ρ densitatea variabilei bidimensionale (f, g) . Notăm

$$\begin{aligned} G(t|f = x) &= p(g < t|f = x) \\ &= \lim_{\alpha, \beta \rightarrow 0} \frac{p((g < t) \cap (x - \alpha \leq f < x + \beta))}{p((x - \alpha \leq f < x + \beta))} \\ &= \lim_{\alpha, \beta \rightarrow 0} \frac{H(x + \beta, t) - H(x - \alpha, t)}{H(x + \beta, \infty) - H(x - \alpha, \infty)} \end{aligned} \quad (6.13)$$

dacă limita există. Dacă limita există în orice $t \in R$ și are ca funcție de t proprietățile unei funcții de repartiție, atunci definim:

Definiția 6.13 Se numește funcția de repartiție a lui g condiționată de $f = x$ funcția $G(\cdot|f = x) : R \rightarrow R$ definită prin (6.13).

Se vede că dacă H este derivabilă, cu densitatea ρ avem

$$G(t|f = x) = \frac{\frac{\partial H(x, t)}{\partial x}}{\frac{\partial H(x, \infty)}{\partial x}} = \frac{\int_{-\infty}^t \rho(x, y) dy}{\int_{-\infty}^{\infty} \rho(x, y) dy} \quad (6.14)$$

Funcțiile $G(\cdot)$ și $G(\cdot|f = x)$ sunt diferite. Se utilizează frecvent deosebirea funcțiilor prin "semnătura" lor, adică prin tipul de argumente, chiar dacă pentru nume se utilizează același simbol. Dacă nu există pericol de confuzie vom scrie și $G(t|x)$ în loc de $G(t|f = x)$.

Definiția 6.14 Se numește densitate de probabilitate a lui g condiționată de $f = x$ o funcție reală $\rho_g(\cdot|f = x)$, integrabilă, pozitivă, astfel ca $\int_{-\infty}^t \rho_g(y|f = x) dy = G(t|f = x)$.

O asemenea densitate nu există întotdeauna. Dacă există atunci în punctele ei de continuitate avem:

$$\frac{\partial G(t|f = x)}{\partial t} = \rho_g(t|f = x) = \frac{\rho(x, t)}{\int_{-\infty}^{\infty} \rho(x, y) dy} \quad (\text{cf. 6.14})$$

Și aici vom scrie $\rho_g(t|x)$ sau $\rho(t|x)$ în loc de $\rho_g(t|f = x)$ dacă nu există pericol de confuzie.

În studiul probabilităților condiționate din lecția 1, cele mai importante formule studiate au fost formula probabilității totale și formula lui Bayes. Cum arată aceste formule în contextul prezent?

1. Formula probabilității totale:

Dacă $p(B|x)$ este o funcție continuă atunci

$$p(B) = \int_{-\infty}^{\infty} p(B|f=x) dF(x) = \int_{-\infty}^{\infty} p(B|f=x) \rho_f(x) dx$$

în particular

$$G(t) = \int_{-\infty}^{\infty} G(t|f=x) dF(x) = \int_{-\infty}^{\infty} G(t|f=x) \rho_f(x) dx$$

1. Formula lui Bayes

$$\rho_f(x|g=y) = \frac{\rho_f(x) \rho_g(y|f=x)}{\int_{-\infty}^{\infty} \rho_g(y|f=x) \rho_f(x) dx}$$

Dacă nu există pericol de confuzie putem scrie:

$$\rho(x|y) = \frac{\rho(x) \rho(y|x)}{\int_{-\infty}^{\infty} \rho(y|x) \rho(x) dx}$$

Nu insistăm asupra demonstrațiilor acestor formule. Ele se pot obține prin divizarea intervalului unde se găsesc valorile lui x , aplicarea formulei clasice a probabilității totale sau a lui Bayes așa cum au fost studiate în lecția 1 și trecerea la limită pentru diviziuni din ce în ce mai fine.

6.4 Distribuția sumei și cântului

Fie $f, g : X \rightarrow R$ două v.a. și $h : X \rightarrow R^2$, $h = (f, g)$ cu repartiția H și densitatea ρ . Care este repartiția sumei $f+g$ și a cântului $\frac{f}{g}$ în eventualitatea că $g \neq 0$? În general dacă r este o funcție continuă $r : R^2 \rightarrow R$ atunci repartiția lui $r(f, g)$ este:

$$F_{r(f,g)}(t) = p(r(f, g) < t) = \int \int_{r(x,y) < t} \rho(x, y) dx dy \quad (6.15)$$

Pentru sumă și produs avem formule mai precise. Alte exemple sunt date la exerciții.

6.4.1 Distribuția sumei

În condițiile de mai sus suma $f + g$ are distribuția

$$F_{f+g}(t) = p(f + g < t) = \int \int_{x+y < t} \rho(x, y) dx dy = \int_{-\infty}^t du \int_{-\infty}^{\infty} \rho(v, u-v) dv$$

Am utilizat schimbarea $u = x + y$, $v = x$. Dacă f, g sunt independente atunci $\rho(x, y) = \rho_f(x) \rho_g(y)$ și densitatea sumei devine

$$\rho_{f+g}(t) = F'_{f+g}(t) = \int_{-\infty}^{\infty} \rho_f(v) \rho_g(t-v) dv \quad (6.16)$$

ceea ce se numește produsul de convoluție al densităților ρ_f și ρ_g .

6.4.2 Distribuția câtlui

Fie f și g ca mai înainte și în plus $g \neq 0$. În aceste condiții câtlul $\frac{f}{g}$ are distribuția:

$$\begin{aligned} F_{\frac{f}{g}}(t) &= p\left(\frac{f}{g} < t\right) = \int_{y>0, x<ty} \rho(x, y) dx dy + \int_{y<0, x>ty} \rho(x, y) dx dy \\ &= \int_0^{\infty} dy \int_{-\infty}^{ty} \rho(x, y) dx + \int_{-\infty}^0 dy \int_{ty}^{\infty} \rho(x, y) dx \end{aligned}$$

Prin derivare în raport cu t găsim

$$\rho_{\frac{f}{g}}(t) = \int_{-\infty}^{\infty} |y| \rho(ty, y) dy \quad (6.17)$$

În particular pentru f și g independente găsim:

$$\rho_{\frac{f}{g}}(t) = \int_{-\infty}^{\infty} |y| \cdot \rho_f(ty) \cdot \rho_g(y) \cdot dy \quad (6.18)$$

6.5 Distribuția Student

Exemplul 6.15 Fie f normală de tipul $N(0, \sigma)$ iar g o variabilă χ^2 de tipul $H(s, \sigma)$, independentă de f . Se cer:

- Densitatea lui $\sqrt{\frac{g}{s}}$.
- Densitatea lui $h = \frac{f}{\sqrt{\frac{g}{s}}}$.
- Momentele lui h .

Soluție. Densitatea lui f este

$$\rho_f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{x^2}{2\sigma^2}}$$

iar a lui g este

$$\rho_g(y) = \frac{1}{2^{\frac{s}{2}} \sigma^s \Gamma(\frac{s}{2})} y^{\frac{s}{2}-1} e^{-\frac{y}{2\sigma^2}} \text{ pentru } y > 0 \text{ și } 0 \text{ în caz contrar.}$$

a) $F_{\sqrt{\frac{g}{s}}}(y) = p\left(\sqrt{\frac{g}{s}} < y\right) = p(g < y^2 s) = F_g(y^2 s)$, pentru $y > 0$ și 0 pentru $y \leq 0$. Prin urmare densitatea lui $\sqrt{\frac{g}{s}}$ este

$$\begin{aligned}\rho_{\sqrt{\frac{g}{s}}}(y) &= \frac{d}{dy} F_{\sqrt{\frac{g}{s}}}(y) = 2sy F'_g(y^2 s) \\ &= 2sy \rho_g(y^2 s) = \frac{2sy}{2^{\frac{s}{2}} \sigma^s \Gamma\left(\frac{s}{2}\right)} (y^2 s)^{\frac{s}{2}-1} e^{-\frac{y^2 s}{2\sigma^2}}\end{aligned}$$

pentru $y > 0$ și 0 pentru $y \leq 0$.

b) Conform cu (6.18) avem pentru $h = \frac{f}{\sqrt{g/s}}$ densitatea

$$\begin{aligned}\rho_h(t) &= \int_{-\infty}^{\infty} |y| \cdot \rho_f(ty) \cdot \rho_g(y) \cdot dy \\ &= \int_0^{\infty} y \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{t^2 y^2}{2\sigma^2}} \cdot \frac{2sy}{2^{\frac{s}{2}} \sigma^s \Gamma\left(\frac{s}{2}\right)} y^{s-2} s^{\frac{s}{2}-1} e^{-\frac{sy^2}{2\sigma^2}} dy \\ &= \frac{2s^{\frac{s}{2}}}{2^{\frac{s+1}{2}} \sqrt{\pi} \sigma^{s+1} \Gamma\left(\frac{s}{2}\right)} \int_0^{\infty} e^{-\frac{(t^2+s)y^2}{2\sigma^2}} y^s dy \quad (\text{schimbăm } sy^2 = u) \\ &= \frac{1}{2^{\frac{s+1}{2}} \sqrt{\pi} \sigma^{s+1} \Gamma\left(\frac{s}{2}\right)} \int_0^{\infty} e^{-\frac{\left(1+\frac{t^2}{s}\right)u}{2\sigma^2}} u^{\frac{s+1}{2}-1} du \\ &= \frac{1}{2^{\frac{s+1}{2}} \sqrt{\pi} \sigma^{s+1} \Gamma\left(\frac{s}{2}\right)} \frac{\Gamma\left(\frac{s+1}{2}\right)}{\left(\frac{1+\frac{t^2}{s}}{2\sigma^2}\right)^{\frac{s+1}{2}}} = \frac{\Gamma\left(\frac{s+1}{2}\right)}{\sqrt{\pi} s \Gamma\left(\frac{s}{2}\right)} \left(1 + \frac{t^2}{s}\right)^{-\frac{s+1}{2}}\end{aligned}$$

c) Media lui h , $M(h)$, este zero, h având densitatea pară. Pentru momentele lui h folosim definiția:

$$M_k(h) = \int_{-\infty}^{\infty} x^k \rho_h(x) dx = \int_{-\infty}^{\infty} (x - M(h))^k \rho_h(x) dx = \mu_k(h)$$

Rezultă momente de ordin impar zero. Înlocuind pe $\rho_h(x)$ cu expresia de mai înainte și făcând schimbarea $\frac{x^2}{s} = u$ și apoi $\frac{u}{1+u} = y$ ajungem la

$$\begin{aligned}\mu_{2k} &= \frac{\Gamma\left(\frac{s+1}{2}\right) s^{k+\frac{1}{2}}}{\sqrt{\pi} s \Gamma\left(\frac{s}{2}\right) 2} \int_0^1 y^{(k+\frac{1}{2})-1} (1-y)^{\left(\frac{s}{2}-k\right)-1} dy \\ &= \frac{\Gamma\left(\frac{s+1}{2}\right) s^{k+\frac{1}{2}} \Gamma\left(k+\frac{1}{2}\right) \Gamma\left(\frac{s}{2}-k\right)}{\sqrt{\pi} s \Gamma\left(\frac{s}{2}\right) 2 \Gamma\left(\frac{s+1}{2}\right)} \\ &= \frac{s^k \Gamma\left(k+\frac{1}{2}\right) \Gamma\left(\frac{s}{2}-k\right)}{\sqrt{\pi} \Gamma\left(\frac{s}{2}\right)} \\ &= \frac{s^k 1 \cdot 3 \cdot \dots \cdot (2k-1)}{(s-2) \cdot (s-4) \cdot \dots \cdot (s-2k)}\end{aligned}$$

dacă $2k < s$. Pentru celelalte valori ale lui k momentele nu există.

Definiția 6.16 O variabilă aleatoare cu densitatea

$$\rho(t) = \frac{\Gamma\left(\frac{s+1}{2}\right)}{\sqrt{\pi s} \Gamma\left(\frac{s}{2}\right)} \left(1 + \frac{t^2}{s}\right)^{-\frac{s+1}{2}}$$

se numește variabilă **Student**. Tipul acesta de v.a. îl notăm $S(s)$. Graficul acestei densități se găsește la lecția 4.

O asemenea variabilă aleatoare este raportul dintre o variabilă normală cu media zero și dispersia σ^2 pe de o parte și $\sqrt{\frac{g}{s}}$ unde g este o variabilă de tipul $H(s, \sigma)$ adică χ^2 cu s grade de libertate obținută ca sumă a s pătrate de variabile normale de tip $N(0, \sigma)$ (vezi lecția 4), pe de altă parte.

6.6 Distribuția Snedecor-Fisher

Exemplul 6.17 Fie f_1 și f_2 două variabile aleatoare independente de tip $H(n_1, \sigma)$ și $H(n_2, \sigma)$. Se cere densitatea variabilei $g = \frac{f_1/n_1}{f_2/n_2}$.

Soluție. f_1 și f_2 sunt variabile χ^2 . Deoarece ele sunt independente densitatea variabilei $h = (f_1, f_2)$ este

$$\rho(x, y) = \rho_{f_1}(x) \rho_{f_2}(y) = \begin{cases} \frac{1}{2^{\frac{n_1+n_2}{2}} \sigma^{n_1+n_2} \Gamma\left(\frac{n_1}{2}\right) \Gamma\left(\frac{n_2}{2}\right)} x^{\frac{n_1}{2}-1} y^{\frac{n_2}{2}-1} e^{-\frac{x+y}{2}}, & x, y > 0 \\ 0, & \text{în rest} \end{cases}$$

dacă $x \geq 0$, $y \geq 0$ și 0 în rest. Aplicând aceeași tehnică precum în exemplul precedent rezultă :

$$\rho_g(x) = \begin{cases} \left(\frac{n_1}{n_2}\right)^{\frac{n_1}{2}} \frac{\Gamma\left(\frac{n_1+n_2}{2}\right)}{\Gamma\left(\frac{n_1}{2}\right) \Gamma\left(\frac{n_2}{2}\right)} x^{\frac{n_1}{2}-1} \left(1 + \frac{n_1}{n_2}x\right)^{-\frac{n_1+n_2}{2}}, & x > 0 \\ 0, & \text{în rest} \end{cases} \quad (6.19)$$

dacă $x \geq 0$ și 0 în caz contrar.

Definiția 6.18 O variabilă aleatoare cu densitatea (6.19) se numește variabilă **Snedecor-Fisher**. Acest tip îl vom nota $S(n_1, n_2)$.

6.7 Exerciții

1. Fie șirul de date

$$\begin{array}{cccccc} x = & 1 & 3 & -1 & 4 & 7 & 9 \\ y = & 0 & 4 & -3 & 8 & 12 & 18 \end{array}$$

a) Care este coeficientul de corelație între aceste șiruri de date?

b) Dacă este mai mare de 0.95, să se determine prin metoda celor mai mici pătrate dreapta de regresie a lui y în x .

Indicație. Coeficientul de corelație (6.3) este 0,99 deci punctele sunt aproximativ pe o dreaptă.

2. Fie șirul de date

$$\begin{array}{cccccc} x = & 1 & 1,5 & 2 & 2,5 & 3 \\ y = & 3,3 & 4,23 & 5,44 & 6,98 & 8,96 \end{array}$$

Să se studieze existența unei dependențe de forma $y = be^{ax}$ între date.

Indicație. Valorile lui $\ln(y)$ sunt

$$1,19 \quad 1,44 \quad 1,69 \quad 1,94 \quad 2,19$$

Coeficientul de corelație între x și $\ln(y)$ este aproximativ 1, prin urmare $\ln(y) = ax + d$ deci $y = e^d e^{ax} = be^{ax}$. Prin metoda celor mai mici pătrate găsim $a = 0,5$ $d = 0,69$ deci $b = 2$.

3. Variabilele aleatoare independente f_1 și f_2 sunt uniform distribuite în $[0, 1]$. Se cer:

a) Distribuțiile pentru $\min(f_1, f_2)$ și $\max(f_1, f_2)$.

b) Distribuția sumei $f_1 + f_2$.

Indicație. Variabila dublă (f_1, f_2) are densitatea $\rho(x, y) = 1$ în $(0, 1) \times (0, 1)$ și 0 în rest. Fie F funcția de repartiție a lui $\min(f_1, f_2)$. Avem

$$\begin{aligned} F(x) &= p(f_1 < x \cup f_2 < x) = 1 - p(\overline{f_1 < x \cup f_2 < x}) \\ &= 1 - p(f_1 \geq x \cap f_2 \geq x) = 1 - p(f_1 \geq x) \cdot p(f_2 \geq x) \\ &= 1 - (1 - x)(1 - x) \end{aligned}$$

pentru $x \in [0, 1]$. În mod analog se calculează repartiția lui $\max(f_1, f_2)$. Pentru sumă se folosește formula (6.16)

4. Fie X și Y două v.a. cu coeficientul de corelație r . Care este coeficientul de corelație al variabilelor $aX + b$ și $cY + d$, unde $a, b, c, d \in \mathbb{R}$?

5. Numărul de mașini pe șoseaua București-Ploiești, dinspre București spre Ploiești, într-un interval de timp I este o variabilă Poisson cu parametrul λ_1 , iar numărul de mașini în

sens contrar, în același interval de timp, este o variabilă Poisson cu parametrul λ_2 . Dacă cele două v.a. sunt independente, care este distribuția numărului total de mașini pe șosea în același interval de timp?

6. Un cuplu de v.a. $f = (f_1, f_2)$ are densitatea dublă

$$\rho(x, y) = \begin{cases} \frac{4(x+3y)e^{-x-3y}}{5}, & \text{dacă } x \geq 0 \quad y \geq 0 \\ 0 & \text{în caz contrar} \end{cases}$$

i) Să se determine densitățile marginale pentru f_1 și f_2 .

ii) Să se calculeze densitatea condițională a lui X dacă $Y=y$.

ii) Să se calculeze coeficientul de corelație dintre X și Y .

7. Se trage asupra unei ținte plane așezate în $(0,0)$. Măsurătorile indică o distribuție normală a absciselor X a punctelor de impact de tip $N(0, 2)$ și la fel pentru ordonatele Y . De asemenea X și Y sunt independente. Se cer:

a) Distribuția distanței $\sqrt{X^2 + Y^2}$ de la punctul de impact până la țintă.

b) Probabilitatea ca o lovitură specificată să cadă în discul $X^2 + Y^2 \leq 4$.

c) Probabilitatea ca din 4 lovituri cel puțin una să cadă în discul $X^2 + Y^2 \leq 4$.

d) Probabilitatea ca din 100 lovituri cel puțin 25 să cadă în discul $X^2 + Y^2 \leq 4$.

Indicație. a) $\rho_1(x) = \frac{1}{\sqrt{2\pi 2}} e^{-\frac{x^2}{8}}$ $\rho_2(x) = \frac{1}{\sqrt{2\pi 2}} e^{-\frac{y^2}{8}}$ deci $\rho(x, y) = \frac{1}{8\pi} e^{-\frac{x^2+y^2}{8}}$. Fie $F(t)$ funcția de repartiție a lui $\sqrt{X^2 + Y^2}$. Avem pentru $F(t) = 0$ pentru $t < 0$ iar pentru $t \geq 0$:

$$F(t) = p(X^2 + Y^2 < t^2) = \frac{1}{8\pi} \int \int_{x^2+y^2 < t^2} e^{-\frac{x^2+y^2}{8}} dx dy$$

și se trece la coordonate polare. b) Probabilitatea este $p = F(2)$. c) Probabilitatea căutată este $1 - (1 - p)^4$. d) Răspunsul este $\sum_{k=25}^{100} C_{100}^k p^k (1-p)^{n-k}$ și se folosește formula integrală de Moivre-Laplace.

8. Perechea de v.a. (f_1, f_2) admite densitatea dublă $\rho(x, y) = \frac{1}{2\pi} \frac{1}{(1+x^2+y^2)^{3/2}}$. Se cere:

a) Să se arate că f_1 și f_2 sunt v.a. de tip Cauchy, adică au densitatea $\rho(x) = \frac{1}{\pi} \frac{1}{1+x^2}$.

b) Să se calculeze $p(f_1^2 + f_2^2 \leq 4)$.

Indicație. a) $\rho_1(x) = \int_{-\infty}^{\infty} \rho(x, y) dy = \dots$ b) Se procedează ca la punctul a) de la problema precedentă.

9. Durata de viață a unui bec ce funcționează continuu este o variabilă aleatoare cu densitatea

$$\rho(t) = \begin{cases} \lambda e^{-\lambda t} & t \geq 0 \\ 0 & t < 0 \end{cases}$$

Intr-un anumit loc trebuie să funcționeze continuu un bec și se dispune de un număr de 10 becuri pentru înlocuire în caz că cel în funcțiune se defectează. Se cer:

a) Valoarea lui λ dacă durata medie de viață a unui bec este de 20 zile (ziua va fi în problemă unitatea de măsură pentru timp).

b) Să se arate că dacă f_1 și f_2 sunt independente, cu densitatea ρ atunci $f_1 + f_2$ are densitatea

$$\rho_1(t) = \begin{cases} \lambda^2 t e^{-\lambda t} & t \geq 0 \\ 0 & t < 0 \end{cases}$$

c) Să se arate că dacă f_1, f_2, \dots, f_n sunt independente și au densitatea ρ de mai sus, atunci $f_1 + f_2 + \dots + f_n$ are densitatea

$$\rho_n(t) = \begin{cases} \frac{\lambda^n t^{n-1}}{(n-1)!} e^{-\lambda t} & t \geq 0 \\ 0 & t < 0 \end{cases}$$

d) Care e probabilitatea ca prin înlocuirea becurilor arse să poată fi menținută lumina aprinsă peste 200 zile?

Indicație. a) $M(f_1) = \frac{1}{\lambda}$ deci $\lambda = \frac{1}{20}$. b,c) Se aplică (5.13) sau se procedează ca în lecția 4, exercițiul nr. 4. d) În datele de la c) avem $\lambda = \frac{1}{20}$, $n = 10$, iar probabilitatea cerută este:

$$\begin{aligned} P &= p(f_1 + f_2 + \dots + f_{10} \geq 200) = \int_{200}^{\infty} \frac{1}{20^{10} \times 9!} t^9 e^{-\frac{t}{20}} dt \\ &= \frac{10^{10}}{9!} \int_1^{\infty} u^9 e^{-10u} du \simeq 0,457 \end{aligned}$$

10) Două variabile aleatoare independente f_1 și f_2 au distribuții normale $N(0, \sigma)$. Să se arate că $f_1^2 + f_2^2$ și $\frac{f_1}{f_2}$ sunt independente.

Indicație. (f_1, f_2) are densitatea $\frac{1}{4\pi\sigma^2} e^{-\frac{x^2+y^2}{2\sigma^2}}$. $\frac{f_1}{f_2}$ este definită a.p.t. deoarece $p(f_2 = 0) = 0$. Fie $\phi : R^2 - \{(x, 0) | x \in R\} \rightarrow R^2$ definită prin $(x, y) \xrightarrow{\phi} \left(r = x^2 + y^2, s = \frac{x}{y}\right)$. Avem $\frac{D(r,s)}{D(x,y)} = -2(s^2 + 1)$. De asemenea $x^2 = \frac{rs^2}{1+s^2}$ $y^2 = \frac{r}{1+s^2}$. Vedem că pentru orice pereche (r, s) există două perechi (x, y) ce îi corespund prin transformarea dată. Pentru $D \subset R^2$ suficient de mic $\phi^{-1}(D) = D_1 \cup D_2$ cu D_1 și D_2 disjuncte în corespondență bijectivă și diferențiabilă cu D . Avem

$$\begin{aligned} & p\left(\left(f_1^2 + f_2^2, \frac{f_1}{f_2}\right) \in D\right) \\ &= p((f_1, f_2) \in \phi^{-1}(D)) \\ &= \frac{1}{4\pi\sigma^2} \int \int_{D_1} e^{-\frac{x^2+y^2}{2\sigma^2}} dx dy + \frac{1}{4\pi\sigma^2} \int \int_{D_2} e^{-\frac{x^2+y^2}{2\sigma^2}} dx dy \\ &= 2 \cdot \frac{1}{4\pi\sigma^2} \int \int_D e^{-\frac{r}{2\sigma^2}} \frac{1}{2(s^2 + 1)} dr ds \end{aligned}$$

Prin urmare perechea de v.a. $\left(f_1^2 + f_2^2, \frac{f_1}{f_2}\right)$ are densitatea $\frac{1}{4\pi\sigma^2} e^{-\frac{r}{2\sigma^2}} \cdot \frac{1}{1+s^2}$. Conform (6.12) și remarci de după ea rezultă că $f_1^2 + f_2^2, \frac{f_1}{f_2}$ sunt independente.

11. Să se arate că dacă $f : X \rightarrow R^2$ este o v.a. bidimensională cu densitatea $\rho_f(x_1, x_2)$ și $g : R^2 \rightarrow R^2$ este o bijectie de clasă C^1 cu iacobianul nenul atunci densitatea variabilei duble $g \circ f$ este $\rho_{g \circ f}(y_1, y_2) = \rho_f(g^{-1}(y_1, y_2)) \frac{D(x_1, x_2)}{D(y_1, y_2)}$.

12. Să se calculeze mediile, dispersiile, corelația distribuțiilor marginale pentru densitatea dublă

$$\rho(x_1, x_2) = \frac{\sqrt{12}}{2\pi} e^{-2(x_1^2 + x_1 x_2 + x_2^2)}$$

Indicație. A se vedea exemplul 3.

13. O variabilă aleatoare bidimensională Z are densitatea dublă

$$\rho(x_1, x_2) = \begin{cases} \alpha(x_1 + x_2 + 2) & \text{dacă } (x_1, x_2) \in [0, 2] \times [1, 3] \\ 0, & \text{în rest} \end{cases}$$

Să se determine:

- i) α .
- ii) Densitățile distribuțiilor marginale
- iii) $p(Z \in \{(x_1, x_2) \mid x_2 > x_1 + 1\})$.
- iv) Fie Z_1 și Z_2 cele două distribuții marginale. Care este media lui $Z_1 + Z_2$.

Indicație. $\int \int_{R^2} \rho(x_1, x_2) dx_1 dx_2 = 1$. Se utilizează formula 6.9

14. O anumită lungime este împărțită în două și sunt măsurate cele două porțiuni de către două persoane în mod independent. Eroarea de măsurare este o variabilă aleatoare cu densitatea $\rho(x) = \begin{cases} \frac{3}{4}(1-x^2) & \text{pentru } x \in [-1, 1] \\ 0 & \text{în rest} \end{cases}$ pentru fiecare persoană. Care este probabilitatea ca eroarea în măsurarea întregii lungimi să fie mai mare ca 1?

Indicație. Cunoaștem densitățile a două variabile aleatoare independente. Deci rezultă imediat densitatea sumei și de aici răspunsul la întrebare. Altfel, putem folosi formula 6.9

Lecția 7

Procese aleatoare

În modelarea probabilistică a unor fenomene, presupunem implicit existența unei mulțimi X de factori necontrolabili, mulțime care nu poate fi exact definită, ce afectează desfășurarea fenomenelor, precum și existența unei probabilități pe această mulțime. Rezultatul observației unui fenomen îl reprezentăm printr-un număr sau mai multe numere. Observând de mai multe ori același fenomen, chiar reproducând condițiile de desfășurare cât mai fidel posibil, rezultatele numerice pot ieși diferite. Intervin aici acei factori necontrolabili care fac să fluctueze rezultatele numerice ale experienței într-o anumită zonă de valori. De exemplu, măsurând independent cu o ruletă obișnuită de 2m o distanță destul de lungă, să zicem în jur de 100 m, persoane diferite vor ajunge la rezultate diferite. Nu putem specifica exact ce factori contribuie la apariția diferenței între rezultate și care este contribuția fiecăruia, dar putem admite existența unei probabilități pe această mulțime X de factori necontrolabili.

Fie acea mărime numerică pe care o urmărim, egală cu ξ . După ce fixăm valorile factorilor controlabili, ξ devine o funcție $\xi : X \rightarrow R$, pe mulțimea factorilor necontrolabili, adică o variabilă aleatoare. Dacă ξ reprezintă o mărime ce evoluează în timp atunci $\xi : T \times X \rightarrow R$, deci variabila aleatoare este $\xi(t, \cdot) : X \rightarrow R$, pentru fiecare $t \in T \subset R$. Mai notăm această variabilă aleatoare prin ξ_t .

Definiția 7.1 O asemenea familie de variabile aleatoare, ce depinde de un parametru $t \in T \subset R$ se numește proces aleator sau proces stochastic.

Fie $F(t, x)$ sau funcția de repartiție a acestei variabile aleatoare, adică:

$$F(t, x) = p(\{\omega \in X \mid \xi(t, \omega) < x\})$$

sau pe scurt

$$F(t, x) = p(\xi(t, \cdot) < x) = p(\xi_t < x)$$

Uneori vom mai nota această funcție $F_t(x)$ pentru a sublinia că e funcție de repartiție în x și depinde de parametrul t . În timp ce ξ nu poate fi specificată exact, neavând o descriere a

mulțimii X , F este o funcție de variabilele reale t, x . Cum funcția de repartiție conține toate informațiile probabilistice despre o v.a., un proces stochastic va fi descris prin funcția sa de repartiție $F(t, x)$. Vom caracteriza procesul stochastic prin $F_t(x)$ sau derivata ei $\rho(t, x) = \rho_t(x) = \frac{\partial F(t, x)}{\partial x}$. Dacă ξ_t ia o valori într-o mulțime discretă $\{v_1, \dots, v_n, \dots\}$ atunci funcția de repartiție este complet determinată de $p_k(t) = p(\xi_t = v_k)$ sau prin $F_t(x) = \sum_{v_k \leq x} p_k(t)$ deci, procesul stochastic va fi descris prin aceste funcții. Mai jos studiem câteva tipuri de procese stochastice.

7.1 Procese Poisson

Să considerăm numărul particulelor cosmice care pătrund într-un anumit volum V , într-un interval de timp. Acest număr depinde de factori incontrolabili și apare ca o valoare aleatoare. Presupunem că ξ_t este un proces aleator pentru valori ale timpului $t \in [0, \infty)$ care ia doar valori naturale $0, 1, 2, \dots, k, \dots$, și anume pentru $t > 0$, ξ_t ia valoarea k dacă în intervalul de timp $(0, t]$ au pătruns în volumul V , k particule cosmice. Pentru fiecare $t \geq 0$, ξ_t este o variabilă aleatoare. Pentru $t < t + \Delta t$, variabila aleatoare $\xi_{t+\Delta t} - \xi_t$ reprezintă numărul de particule cosmice care au pătruns în volumul V în intervalul de timp $(t, t + \Delta t]$. Variabilele ξ_t nu sunt independente deoarece valoarea lui $\xi_{t+\Delta t}$ depinde de valoarea lui ξ_t . Următoarele ipoteze asupra variabilelor ξ_t apar ca naturale:

a) *Absența post efectului*, adică numărul de particule care în intervalul de timp $[a, b]$ intră în volumul V este independent de numărul de particule care au pătruns anterior în acest volum și de momentele la care au pătruns. Exprimăm acest fapt astfel: dacă $0 \leq a_1 < b_1 \leq a_2 < b_2 \leq \dots a_n < b_n$ atunci avem v.a. $\xi_{b_1} - \xi_{a_1}, \dots, \xi_{b_n} - \xi_{a_n}$ care au ca valoare numărul de particule ce au pătruns în V în intervalele de timp $(a_1, b_1], \dots, (a_n, b_n]$; ipoteza a) înseamnă că aceste v.a. sunt independente, adică

$$\begin{aligned} & p(\xi_{b_1} - \xi_{a_1} = k_1, \xi_{b_2} - \xi_{a_2} = k_2, \dots, \xi_{b_n} - \xi_{a_n} = k_n) \\ &= p(\xi_{b_1} - \xi_{a_1} = k_1) \cdot p(\xi_{b_2} - \xi_{a_2} = k_2) \cdot \dots \cdot p(\xi_{b_n} - \xi_{a_n} = k_n) \end{aligned} \quad (7.1)$$

pentru orice $k_1, \dots, k_n \in N$ și orice număr n de intervale.

b) *Staționaritatea*, înseamnă că numărul de particule ce pătrunde în volum într-un interval de timp depinde doar de lungimea intervalului de timp și nu de plasarea celui interval pe axa timpului. Exprimăm acest lucru prin:

$$p(\xi_b - \xi_a = k) = p(\xi_{a+T} - \xi_{b+T} = k) \quad (7.2)$$

oricare ar fi $0 \leq a < b$, $T \geq 0$ și oricare ar fi $k \in N$.

c) *Ordinaritatea*, adică probabilitatea ca în intervalul $[t, t + \Delta t]$ să pătrundă în volum mai mult de o particulă este zero în raport cu Δt . Mai precis:

$$p(\xi_{t+\Delta t} - \xi_t > 1) = O(\Delta t)$$

Prin $O(\Delta t)$ se înțelege o funcție reală O de argument Δt , astfel ca $\lim_{\Delta t \rightarrow 0} \frac{O(\Delta t)}{\Delta t} = 0$.

Definiția 7.2 Un proces aleator $\xi_t, t \in [0, \infty)$ unde ξ_t ia valori numere naturale și unde sunt satisfăcute condițiile a), b), c) de mai sus se numește proces Poisson.

Multe alte fenomene din natură apar ca satisfacând cerințele a), b), c) : dezintegrarea spontană a nucleelor radioactive, venirea la coadă a mașinilor la o stație de benzină, numărul de apeluri la o centrală telefonică etc.

Problema pe care vrem să o rezolvăm în continuare este de a determina probabilitățile $p_k(t) = p(\xi_{a+t} - \xi_a = k)$, adică probabilitatea ca în într-un interval de lungime t să intre k particule cosmice în volumul V . Conform cu b) aceste probabilități nu depind de a .

Teorema 7.3 Fie un proces Poisson ca mai sus. Atunci există $\lambda > 0$ astfel ca $p_1(t) = \lambda t + o(t)$ și

$$p_k(t) = \frac{(\lambda t)^k}{k!} e^{-\lambda t} \quad (7.3)$$

Demonstrație.(Schiță) Fie $r = p_0(1)$. Impărțim intervalul $[0, 1]$ prin

$$0 < \frac{1}{n} < \frac{2}{n} < \dots < \frac{n}{n}$$

, pentru un $n \in \mathbb{N}$. Din staționaritate rezultă

$$p_0\left(\frac{1}{n}\right) = p(\xi_{1/n} - \xi_0 = 0) = p(\xi_{2/n} - \xi_{1/n} = 0) = \dots = p\left(\xi_1 - \xi_{\frac{n-1}{n}} = 0\right)$$

Utilizând acum faptul că nu există post efect, avem:

$$\begin{aligned} r &= p_0(1) = p\left(\xi_{1/n} - \xi_0 = 0, \xi_{2/n} - \xi_{1/n} = 0, \xi_1 - \xi_{\frac{n-1}{n}} = 0\right) = \\ &= p(\xi_{1/n} - \xi_0 = 0) \cdot p(\xi_{2/n} - \xi_{1/n} = 0) \cdot \dots \cdot p\left(\xi_1 - \xi_{\frac{n-1}{n}} = 0\right) = \\ &= (p(\xi_{1/n} - \xi_0 = 0))^n = (p_0(1/n))^n \end{aligned}$$

De aici rezultă $p_0\left(\frac{1}{n}\right) = r^{\frac{1}{n}}$. Evenimentul că în intervalul $(0, \frac{k}{n}]$ nu intră nici o particulă în V , este echivalent cu intersecția evenimentelor " în intervalul $(\frac{i-1}{n}, \frac{i}{n}]$ nu intră vreo particulă în V ", pentru $1 \leq i \leq k$. Cum aceste evenimente sunt independente, rezultă $p_0\left(\frac{k}{n}\right) = \left(p_0\left(\frac{1}{n}\right)\right)^k = r^{\frac{k}{n}}$. Acum, deoarece $p_0(t)$ este crescătoare în t , rezultă ușor că $p_0(t) = r^t$. Cum r este o probabilitate, $r \in (0, 1)$, deci $r = e^{-\lambda}$ pentru un $\lambda > 0$, deci $p_0(t) = e^{-\lambda t}$. Am demonstrat deci formula 7.3 pentru $k=0$.

Mai departe avem

$$p_0(t) + p_1(t) + \underbrace{\sum_{k \geq 1} p_k(t)}_{=0(t) \text{ conform cu c)}} = 1$$

adică $e^{-\lambda t} + p_1(t) + 0(t) = 1$. De aici rezultă

$$p_1(t) = 1 - e^{-\lambda t} + O(t) = \lambda t + O(t)$$

Pentru a determina $p_k(t)$ este ușor de justificat formula

$$\begin{aligned} p_k(t + \Delta t) &= \sum_{i=0}^k p_i(\Delta t) p_{k-i}(t) \\ &= e^{-\lambda \Delta t} p_k(t) + (\lambda \Delta t + 0(\Delta t)) p_{k-1}(t) + \underbrace{\sum_{i=2}^k p_i(\Delta t) p_{k-i}(t)}_{=0(\Delta t) \text{ conform cu c)}} \\ &= p_k(t) - \lambda \Delta t p_k(t) + \lambda \Delta t p_{k-1}(t) + 0(\Delta t) \end{aligned}$$

pentru $k \geq 1$. Prin trecerea lui $p_k(t)$ în membrul stâng și divizarea la Δt , apoi trecerea la limită $\Delta t \rightarrow 0$, rezultă

$$\frac{dp_k(t)}{dt} = -\lambda p_k(t) + \lambda p_{k-1}(t) \quad (7.4)$$

Dacă ținem seama de condițiile $p_0(t) = e^{-\lambda t}$, $p_k(0) = 0$ (din ordinaritate), atunci sistemul 7.4 se rezolvă și se găsește $p_k(t) = \frac{(\lambda t)^k}{k!} e^{-\lambda t}$.

QED.

Probabilitățile $p_k(t) = p(\xi_t = k)$ sunt la fel ca la o v.a. Poisson (vezi lecția 4), deci media lui ξ_t este $M(\xi_t) = \lambda t$ și dispersia $D(\xi_t) = \lambda t$.

7.2 Procese Markov discrete

Definiția 7.4 Un proces stochastic se numește proces Markov discret dacă ξ_t poate lua un număr finit de valori $V = \{v_1, v_2, \dots, v_n\}$, timpul t variază într-o mulțime discretă de valori $T = \{t_1, t_2, \dots, t_n, \dots\}$ și

$$p(\xi_{t_k} = v \mid \xi_{t_{k-1}} = v_{i_{k-1}}, \xi_{t_{k-2}} = v_{i_{k-2}}, \dots, \xi_{t_1} = v_{i_1}) = p(\xi_{t_k} = v \mid \xi_{t_{k-1}} = v_{i_{k-1}})$$

pentru orice $v, v_{i_{k-1}}, \dots, v_{i_1} \in V$.

Cu alte cuvinte probabilitatea ca sistemul (variabila aleatoare ξ) să ajungă în momentul t_k în starea v_{i_k} depinde doar de starea $v_{i_{k-1}}$ de la momentul imediat anterior, t_{k-1} , și nu depinde de stările la momentele $t_{k-2}, t_{k-3}, \dots, t_1$. Fără a pierde din generalitate putem considera mulțimea stărilor $V = \{1, 2, \dots, n\}$, iar mulțimea valorilor temporale $T = \{0, 1, 2, \dots, n, \dots\}$. Vom nota cu $p_{i,j}(k) = p(\xi_k = j \mid \xi_{k-1} = i)$, adică $p_{i,j}(k)$ este probabilitatea ca în cazul când la momentul $k-1$ sistemul este în starea i , el să ajungă în starea j la momentul următor, k . Aceste

mărimi se numesc *probabilități de tranziție*. Deoarece din starea i se poate ajunge la momentul următor doar în stările $1, 2, 3, \dots, n$, trebuie să avem $p_{i,1}(k) + p_{i,2}(k) + \dots + p_{i,n}(k) = 1$. Cel mai simplu proces de acest tip este acela în care probabilitățile de tranziție $p_{i,j}(k)$ nu depind de momentul k . În acest caz avem $p_{i,j}(k) = p_{i,j}$. Matricea p cu elementele $p_{i,j}$ se numește matricea de tranziție. Într-o astfel de matrice $\sum_{j=1}^n p_{i,j} = 1$, $p_{i,j} \geq 0$.

Exemplul 7.5 *Să presupunem că o particulă se poate mișca între două bariere $a < b$ trecând prin puncte intermediare*

$$a = x_1 < x_2 < \dots < x_n = b$$

Dacă particula se găsește în poziția x_i atunci cu probabilitatea r se deplasează înainte în poziția x_{i+1} și cu probabilitatea $s = 1 - p$ se deplasează înapoi în poziția x_{i-1} . În capetele a și b particula este respinsă în poziția imediat vecină. Pentru 4 poziții posibile, matricea de tranziție este

$$p = \begin{pmatrix} 0 & 1 & 0 & 0 \\ s & 0 & r & 0 \\ 0 & s & 0 & r \\ 0 & 0 & 1 & 0 \end{pmatrix}$$

Fie acum $p_{i,j}^{(k)} = p(\xi_k = j \mid \xi_0 = i)$, adică $p_{i,j}^{(k)}$ este probabilitatea ca sistemul să se găsească la momentul k în starea j , dacă la momentul inițial 0 s-a găsit în starea i . Fie $p^{(k)}$ matricea de tranziție de la momentul 0 la momentul k , având elementele $p_{i,j}^{(k)}$. Din formula probabilității totale găsim $p(\xi_k = j \mid \xi_0 = i) = \sum_{l=0}^n p(\xi_{k-1} = l \mid \xi_0 = i) \cdot p(\xi_k = j \mid \xi_{k-1} = l, \xi_0 = i)$
 $= \sum_{l=0}^n p(\xi_{k-1} = l \mid \xi_0 = i) \cdot p(\xi_k = j \mid \xi_{k-1} = l)$
 altfel spus

$$p_{i,j}^{(k)} = \sum_{l=0}^n p_{i,l}^{(k-1)} \cdot p_{l,j}$$

Din această formulă rezultă prin inducție

$$p^{(k)} = p^k$$

adică $p^{(k)}$ este puterea de ordinul k a matricei de tranziție.

Exemplul 7.6 *Pentru matricea de tranziție din exemplul 1 găsim*

$$p^2 = \begin{pmatrix} s & 0 & r & 0 \\ 0 & s + r \cdot s & 0 & r^2 \\ s^2 & 0 & r + r \cdot s & 0 \\ 0 & s & 0 & r \end{pmatrix}$$

Pentru $r=3/4$ și $s=1/4$ matricea de tranziție devine:

$$p = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 1/4 & 0 & 3/4 & 0 \\ 0 & 1/4 & 0 & 3/4 \\ 0 & 0 & 1 & 0 \end{pmatrix}$$

Pentru valori mari ale lui k găsim:

$$p^{(100)} = \begin{pmatrix} 0,0769231 & 0 & 0,923077 & 0 \\ 0 & 0,307692 & 0 & 0,692308 \\ 0,0769231 & 0 & 0,923077 & 0 \\ 0 & 0,307692 & 0 & 0,692308 \end{pmatrix}$$

Cum arată probabilitățile $p_{i,j}^{(k)}$ pentru valori mari ale lui k ? Următoarea teoremă aduce lămuriri în această privință.

Teorema 7.7 Dacă pentru $m \in N$ are loc inegalitatea $p_{i,j}^{(m)} > 0$ pentru orice i, j , atunci există $\lim_{k \rightarrow \infty} p^{(k)}$ și

$$\lim_{k \rightarrow \infty} p_{i,j}^{(k)} = p_j$$

independent de i .

Altfel spus, probabilitățile de tranziție de la starea i în momentul 0, la starea j în momentul k , se stabilizează pentru $k \rightarrow \infty$ la valori independente de starea la momentul 0.

Pentru demonstrație se poate consulta bibliografia.

Exemplul 7.8 Fie matricea de tranziție

$$p = \begin{pmatrix} 1/8 & 2/8 & 2/8 & 3/8 \\ 1/4 & 1/4 & 1/4 & 1/4 \\ 1/10 & 5/10 & 2/10 & 2/10 \\ 7/20 & 3/20 & 5/20 & 5/20 \end{pmatrix}$$

Utilizând calculatorul găsim

$$p^7 = \begin{pmatrix} 0,214013 & 0,283041 & 0,238095 & 0,264851 \\ 0,21402 & 0,283038 & 0,238095 & 0,264846 \\ 0,214014 & 0,28304 & 0,238095 & 0,26485 \\ 0,214024 & 0,283037 & 0,238095 & 0,264844 \end{pmatrix}$$

$$p^{100} = \begin{pmatrix} 0,214018 & 0,283039 & 0,238095 & 0,264848 \\ 0,214018 & 0,283039 & 0,238095 & 0,264848 \\ 0,214018 & 0,283039 & 0,238095 & 0,264848 \\ 0,214018 & 0,283039 & 0,238095 & 0,264848 \end{pmatrix}$$

Se vede cum de la p^7 se stabilizează primele trei zecimale ale probabilităților limită, la p^{100} fiind stabilizate cel puțin 6 zecimale.

7.3 Procese de naștere și moarte

Să presupunem că într-un proces stochastic ξ_t ia valori numere naturale și că următoarele condiții sunt îndeplinite:

i) $p(\xi_{t+\Delta t} = k+1 | \xi_t = k) = \lambda_k \cdot \Delta t + O(\Delta t)$, $\lambda > 0$, (aceasta este probabilitatea unui proces de naștere în intervalul $(t, t + \Delta t)$)

ii) $p(\xi_{t+\Delta t} = k-1 | \xi_t = k) = \mu_k \cdot \Delta t + O(\Delta t)$, $\mu > 0$, $k \geq 1$ (aceasta este probabilitatea unui proces de moarte în intervalul $(t, t + \Delta t)$).

iii) $p(|\xi_{t+\Delta t} - \xi_t| > 1) = O(\Delta t)$ (aceasta este probabilitatea a mai mult de o naștere sau o moarte în intervalul $(t, t + \Delta t)$)

Definiția 7.9 Un proces stochastic cu valori naturale și care îndeplinește condițiile i)-iii) de mai sus se numește proces de naștere și moarte.

Dacă $\mu_k = 0$ atunci procesul se numește proces de naștere iar dacă $\lambda_k = 0$ se numește proces de moarte. Deoarece pentru $k \geq 1$ avem

$$\begin{aligned} 1 &= \sum_{n=0}^{\infty} p(\xi_{t+\Delta t} = n | \xi_t = k) \\ &= p(\xi_{t+\Delta t} = k-1 | \xi_t = k) + p(\xi_{t+\Delta t} = k | \xi_t = k) \\ &\quad + p(\xi_{t+\Delta t} = k+1 | \xi_t = k) + p(|\xi_{t+\Delta t} - k| \geq 2 | \xi_t = k) \end{aligned}$$

rezultă imediat din i). ii), iii) că

$$\text{iv) } p(\xi_{t+\Delta t} = k | \xi_t = k) = 1 - \lambda_k \Delta t - \mu_k \Delta t + O(\Delta t).$$

Problema care se pune la aceste procese constă în determinarea probabilităților $p(\xi_t = k) = p_k(t)$ **cunoscând** $p_k(0) = p(\xi_0 = k)$.

Soluție. Deoarece evenimentele $\{\xi_t = k\}_{k \in \mathbb{N}}$ sunt disjuncte avem conform formulei probabilității totale pentru $k \geq 1$

$$p(\xi_{t+\Delta t} = k) = p(\xi_t = k) \cdot p(\xi_{t+\Delta t} = k | \xi_t = k) \quad (7.5)$$

$$+ p(\xi_t = k+1) \cdot p(\xi_{t+\Delta t} = k | \xi_t = k+1) \quad (7.6)$$

$$+ p(\xi_t = k-1) \cdot p(\xi_{t+\Delta t} = k | \xi_t = k-1) \quad (7.7)$$

$$+ p(|\xi_t - k| \geq 2) \cdot p(\xi_{t+\Delta t} = k | |\xi_t - k| \geq 2) \quad (7.8)$$

Folosind relația iv) în (7.5), relația ii) în (7.6), relația i) în (7.7), relația iii) în (7.8) găsim:

$$\begin{aligned} p_k(t + \Delta t) &= p_k(t) \cdot (1 - \lambda_k \Delta t - \mu_k \Delta t + O(\Delta t)) \\ &\quad + p_{k+1}(t) \cdot (\mu_{k+1} \Delta t + O(\Delta t)) \\ &\quad + p_{k-1}(t) \cdot (\lambda_{k-1} \Delta t + O(\Delta t)) \\ &\quad + O(\Delta t) \end{aligned}$$

Trecând $p_k(t)$ în membrul stâng, divizând prin Δt , apoi trecând la limită $\Delta t \rightarrow 0$ găsim

$$p'_k(t) = \lambda_{k-1} p_{k-1}(t) - (\lambda_k + \mu_k) p_k(t) + \mu_{k+1} p_{k+1}(t) \quad (7.9)$$

Pentru cazul $k = 0$ o analiză similară conduce la:

$$p'_0(t) = -\lambda_0 p_0(t) + \mu_1 p_1(t) \quad (7.10)$$

Rezolvarea acestui sistem infinit de ecuații diferențiale este anevoioasă, însă în practică după un proces de tranziție haotic urmează o stabilizare, în care $p_k(t)$ sunt constante. Ecuațiile în acest caz sunt

$$\begin{cases} 0 = -\lambda_0 p_0 + \mu_1 p_1 \\ 0 = \lambda_{k-1} p_{k-1} - (\lambda_k + \mu_k) p_k + \mu_{k+1} p_{k+1} \end{cases} \quad (7.11)$$

Din prima ecuație rezultă $p_1 = \frac{\lambda_0}{\lambda_1} p_0$. Folosind celelalte ecuații găsim succesiv $p_1 = \frac{\lambda_0 \lambda_1}{\mu_1 \mu_2}$, ș.a.m.d. În general obținem:

$$p_n = \frac{\lambda_0 \lambda_1 \dots \lambda_{n-1}}{\mu_1 \mu_2 \dots \mu_n} p_0 \quad (7.12)$$

Valoarea lui p_0 se determină din condiția

$$\sum_{i=0}^{\infty} p_i = 1 \quad (7.13)$$

și acest lucru este posibil numai dacă $\sum_n \frac{\lambda_0 \lambda_1 \dots \lambda_{n-1}}{\mu_1 \mu_2 \dots \mu_n} < \infty$.

În acest fel poate fi modelat fenomenul de așteptat la coadă.

7.3.1 Model de așteptare cu o singură stație de deservire și un număr mare de unități ce au nevoie de serviciile stației

Se poate imagina o astfel de situație la o benzinărie cu o singură stație la care vin aleator mașini pentru alimentare. Mașinile provin dintr-o populație mare astfel încât coada nu este influențată de diminuarea numărului de mașini care necesită alimentare.

Prima problemă care se pune este modelarea caracterului întâmplător al venirilor în stație și al plecărilor din stație. Fie λ numărul mediu de venituri în unitatea de timp. Deci într-un timp Δt vor fi în medie $\lambda \Delta t$ venituri. Dacă numărul de mașini din care se vine la coadă este n iar nevoia de benzină este întâmplătoare și independentă de la o mașină la alta atunci probabilitatea de a veni la coadă k mașini în intervalul $(t, t + \Delta t)$ este $C_n^k p^k q^{n-k}$ unde p este probabilitatea ca o mașină să aibă nevoie de benzină iar $q = 1 - p$ (vezi legea binomială). Numărul mediu de venituri este np iar pe de altă parte este $\lambda \cdot \Delta t$. Deci $np = \lambda \Delta t \Rightarrow p = \frac{\lambda \Delta t}{n}$. În aceste condiții știm lecția 4 că $C_n^k p^k q^{n-k} \rightarrow e^{-\lambda \Delta t} \frac{(\lambda \Delta t)^k}{k!}$ atunci când $n \rightarrow \infty$. Cum n este mare, putem considera că probabilitatea ca la coadă să vină k mașini în intervalul $(t, t + \Delta t)$ este $e^{-\lambda \Delta t} \frac{(\lambda \Delta t)^k}{k!}$. **Numărul de venituri la coadă în intervalul de timp Δt poate fi considerat ca o variabilă aleatoare de tip Poisson:**

$$\begin{pmatrix} 0 & 1 & \dots & k & \dots \\ e^{-\lambda \Delta t} & e^{-\lambda \Delta t} \frac{\lambda \Delta t}{1} & \dots & e^{-\lambda \Delta t} \frac{(\lambda \Delta t)^k}{k!} & \dots \end{pmatrix}$$

Prin urmare :

a) Probabilitatea ca în intervalul $(t, t + \Delta t)$ să vină în stație o mașină este $e^{-\lambda \Delta t} \frac{\lambda \Delta t}{1} = \lambda \cdot \Delta t + O(\Delta t)$.

b) Probabilitatea ca să nu vină în stație nici o mașină în intervalul $(t, t + \Delta t)$ este $e^{-\lambda \Delta t} = 1 - \lambda \cdot \Delta t + O(\Delta t)$.

c) Probabilitatea ca în intervalul $(t, t + \Delta t)$ să vină în stație mai mult de o mașină este $\sum_{k \geq 2} e^{-\lambda \Delta t} \frac{(\lambda \Delta t)^k}{k!} = O(\Delta t)$.

În mod analog, fie μ numărul mediu de mașini deservite de o stație în unitatea de timp (presupunând că are continuu de lucru). Într-un interval de lungime Δt vor fi în medie deservite $\mu \cdot \Delta t$ mașini. Numărul real al celor deservite poate fi mai mare sau mai mic decât media și depinde de factori întâmplători, ca necesarul de benzină, întârzieri produse de șofer, etc. Admitem că:

a') Probabilitatea ca în intervalul $(t, t + \Delta t)$ să plece din stație o mașină, dacă există vreuna, este $\mu \Delta t + O(\Delta t)$

b') Probabilitatea ca în intervalul $(t, t + \Delta t)$ să nu plece din stație nici o mașină, dacă există vreuna, este $1 - \mu \Delta t + O(\Delta t)$

c') Probabilitatea ca în intervalul $(t, t + \Delta t)$ să plece din stație mai mult de o mașină este $O(\Delta t)$.

Ipotezele făcute asupra plecărilor din stație sunt asemănătoare cu condițiile a), b), c) îndeplinite de probabilitățile de venire în stație.

Fie acum familia de variabile aleatoare $\xi_t, t \in T = [0, \infty)$, unde valoarea lui ξ_t este egală cu numărul de mașini în stație la momentul t . Ipotezele făcute asupra venirilor și plecărilor, independente unele de altele, se mai scriu:

i)

$$\begin{aligned}
p(\xi_{t+\Delta t} = k+1 | \xi_t = k) &= \underbrace{(\lambda \cdot \Delta t + O(\Delta t))}_{\text{o venire}} \cdot \underbrace{(1 - \mu \Delta t - O(\Delta t))}_{\text{nici o plecare}} + \\
+ \underbrace{O(\Delta t)}_{\text{mai mult de o venire sau o plecare}} &= \lambda \cdot \Delta t + O(\Delta t)
\end{aligned}$$

ii)

$$\begin{aligned}
p(\xi_{t+\Delta t} = k-1 | \xi_t = k) &= \underbrace{(\mu \cdot \Delta t + O(\Delta t))}_{\text{o plecare}} \cdot \underbrace{(1 - \lambda \Delta t - O(\Delta t))}_{\text{nici o venire}} + \\
+ \underbrace{O(\Delta t)}_{\text{mai mult de o plecare sau venire}} &= \mu \cdot \Delta t + O(\Delta t)
\end{aligned}$$

din condiția a') de la plecări și condiția a) de la venituri (o plecare și nici o venire).

iii) $p(|\xi_{t+\Delta t} - \xi_t| > 1) = O(\Delta t)$ din condițiile c) de la venituri și c') de la plecări

Prin urmare avem de a face cu un proces de naștere și moarte. Mărimile $p_k(t) = p(\xi_t = k)$ satisfac deci ecuațiile (7.9)-(7.10), iar la stabilizare ecuațiile (7.11). Soluția la stabilizare este (7.12) unde $\lambda_k = \lambda$ și $\mu_k = \mu$, adică

$$p_k = \left(\frac{\lambda}{\mu}\right)^k p_0 \quad (7.14)$$

iar din $1 = p_0 \sum_{k=0}^{\infty} \left(\frac{\lambda}{\mu}\right)^k = p_0 \frac{1}{1-\frac{\lambda}{\mu}}$ rezultă

$$p_0 = 1 - \frac{\lambda}{\mu}$$

Stabilizarea se poate face doar cu condiția $0 < \frac{\lambda}{\mu} < 1$ astfel încât seria $\sum \left(\frac{\lambda}{\mu}\right)^k$ să fie convergentă.

7.3.2 Model de așteptare cu o singură stație

iar numărul de unități care au nevoie de serviciile stației este limitat la o valoare dată.

Presupunem acum că populația de unde provin mașinile (unitățile) în așteptare este limitată la un număr de N exemplare. Coada care se formează depinde de cât de des au nevoie mașinile de serviciile stației și cât de prompte sunt aceste servicii. În ce privește capacitatea de deservire a stației, nu apar elemente care să modifice ipotezele a'), b'), c') anterioare. În ce privește venirea în stație, ele apar în mod normal proporționale cu numărul mașinilor

în activitate (deci care nu sunt așezate la coadă). Vom face ipoteza că mașinile au nevoie independent de serviciile stației, și dacă în stație sunt k mașini, atunci probabilitatea ca să mai vină una în următoarele Δt unități de timp este $(N - k) \lambda \Delta t + o(\Delta t)$ dacă $k < N$ și zero dacă $k = N$. Astfel familia de variabile aleatoare ξ_t care au ca valoare numărul de mașini în stație este un proces de naștere și moarte cu $\lambda_k = (N - k) \lambda$, iar $\mu_k = \mu$ pentru $0 \leq k \leq N$. Prin urmare, conform formulelor (7.12) rezultă

$$p_k = \frac{N(N-1) \cdot \dots \cdot (N-k+1) \lambda^k}{\mu^k} p_0, \quad 0 \leq k \leq N \quad (7.15)$$

Valoarea lui p_0 se determină din condiția ca $\sum_{k=0}^N p_k = 1$.

Un asemenea proces ar putea modela de exemplu coada la reparații într-o întreprindere cu N mașini dacă serviciul de reparații are doar un mecanic. Un model mai realist este:

7.3.3 Model de așteptare cu n stații de deservire și cu N unități ce trebuie deservite ($1 < n < N$)

În acest caz avem de a face cu un proces de naștere și moarte în care $\lambda_k = N - k$ pentru $0 \leq k \leq N$ (probabilitățile de a se veni la stație pentru alimentare, reparații etc. sunt proporționale cu numărul unităților în serviciu), și $\mu_k = k\mu$ dacă $1 \leq k < n$, $\mu_k = n\mu$ dacă $n \leq k \leq N$ (probabilitatea ca o unitate să părăsească stația este proporțională cu numărul celor în curs de deservire). Conform formulei (7.12) găsim

$$p_k = \begin{cases} C_N^k \left(\frac{\lambda}{\mu}\right)^k p_0 & \text{pentru } 1 \leq k < n-1 \\ \frac{k!}{n!n^{k-n}} C_N^k \left(\frac{\lambda}{\mu}\right)^k p_0 & \text{pentru } n \leq k \leq N \end{cases} \quad (7.16)$$

Si aici p_0 se determină din condiția ca $\sum p_k = 1$.

7.4 Procese aleatoare staționare

Un proces aleator se numește staționar dacă proprietățile sale statistice sunt invariante la o translație a timpului. Fie F_t funcția de repartiție a v.a. ξ_t și fie ρ_i densitatea de probabilitate dacă există. În continuare în această secțiune timpul t aparține mulțimii $T = (-\infty, \infty)$ sau mulțimii $T = [0, \infty)$, în afară de cazul când se specifică o altă situație.

Definiția 7.10 *Procesul aleator ξ_t se numește staționar de ordinul întâi dacă $F_t = F_{t+\tau}$ pentru orice τ .*

Prin urmare pentru procese staționare de ordinul unu, funcția de repartiție este aceeași pentru toate variabilele ξ_t .

Să considerăm acum pentru două valori t_1, t_2 vectorul aleator (ξ_{t_1}, ξ_{t_2}) . Acest vector are o funcție de repartiție mixtă (vezi lecția 6) definită prin $F_{t_1, t_2}(x, y) = p(\xi_{t_1} < x, \xi_{t_2} < y)$.

Definiția 7.11 Spunem că procesul aleator ξ_t este staționar de ordinul doi dacă funcția de repartiție mixtă este invariantă la o translație a timpului. Mai precis $F_{t_1, t_2} = F_{t_1+\tau, t_2+\tau}$ pentru orice τ .

În mod analog se poate defini un proces staționar de ordinul n , prin funcția de repartiție n dimensională a vectorului aleator $(\xi_{t_1}, \xi_{t_2}, \dots, \xi_{t_n})$. Dacă un proces aleator este staționar de ordin n , atunci el este staționar de orice ordin $0 \leq k \leq n$.

Pentru un proces staționar media, dispersia, momentele de diverse ordine ale variabilei ξ_t nu depind de t . Variabilele aleatoare ξ_t nu sunt în general independente. Dacă m este media variabilei ξ_t (independent de t) atunci coeficientul de corelație

$$c(\xi_{t_1}, \xi_{t_2}) = \frac{M((\xi_{t_1} - m) \cdot (\xi_{t_2} - m))}{\sqrt{M((\xi_{t_1} - m)^2) \cdot M((\xi_{t_2} - m)^2)}}$$

(vezi lecția 6) este și el invariant la o translație în timp, adică $c(\xi_{t_1}, \xi_{t_2}) = c(\xi_{t_1+\tau}, \xi_{t_2+\tau})$ pentru orice τ . Însă în multe aplicații ale probabilităților se utilizează doar media, dispersia pentru variabilele aleatoare individuale și coeficientul de corelație pentru perechile de variabile aleatoare. Aceste mărimi pot fi invariante la o translație a timpului fără ca funcția de repartiție să fie. De aceea s-au studiat în mod special procesele aleatoare pentru care aceste mărimi sunt invariante la o translație a timpului, fără să se ceară și invarianța funcției de repartiție.

Definiția 7.12 Un proces aleator $(\xi_t)_{t \in T}$ se numește staționar în sens larg dacă

- a) media $M(\xi_t) = m$ este independentă de $t \in T$.
- b) dispersia $M(\xi_t) = \sigma^2$ este independentă de $t \in T$.
- c) $c(\xi_{t_1}, \xi_{t_2}) = c(\xi_{t_1+\tau}, \xi_{t_2+\tau})$ pentru orice $t_1, t_2, t_1 + \tau, t_2 + \tau \in T$.

Observația 7.13 Condițiile de mai sus sunt echivalente cu

- a') media $M(\xi_t) = m$ este independentă de $t \in T$.
 - b') $M(\xi_t \xi_{t+\tau})$ este funcție doar de τ (deci independentă de t)
- Demonstrația echivalenței celor două seturi de condiții este lăsată ca exercițiu.

În continuare vom nota $R(t_1, t_2) = M(\xi_{t_1} \xi_{t_2})$ și o vom numi funcția de autocorelație a procesului aleator. Dacă procesul este staționar în sens larg, această funcție este invariantă la o translație în timp, deci depinde doar de diferența $t_2 - t_1$. Vom nota în acest caz funcția de autocorelație prin $R(\tau) = M(\xi_t \xi_{t+\tau}) = M(\xi_0 \xi_\tau)$. Câteva din proprietățile funcției de autocorelație pentru procese staționare în sens larg, apar în continuare:

1.

$$R(0) = M_2(\xi_t)$$

pentru că $R(0) = M(\xi_t \xi_{t+0}) = M_2(\xi_t)$.

2.

$$R(-\tau) = R(\tau)$$

pentru că $R(-\tau) = M(\xi_t \xi_{t-\tau}) = M(\xi_{t-\tau} \xi_t) = R(\tau)$.

3.

$$R(0) \geq |R(\tau)|$$

pentru că din $M((\xi_t \pm \xi_{t+\tau})^2) \geq 0$ rezultă $M(\xi_t^2) + M(\xi_{t+\tau}^2) \pm 2M(\xi_t \xi_{t+\tau}) \geq 0$, sau $2R(0) \pm 2R(\tau) \geq 0$.

4. Pentru orice $n \in \mathbb{N}$ și $\tau_1, \tau_2, \dots, \tau_n \in R$ forma pătratică

$$\sum_{i,j=1}^n R(\tau_i - \tau_j) x_i x_j$$

este pozitiv definită pentru că este egală cu $M((\sum_{i=1}^n x_i \xi_{t+\tau_i})^2) \geq 0$, oricare ar fi t astfel ca $t + \tau_i \in T$ pentru orice i .

Definiția 7.14 Procesul aleator staționar în sens larg, $(\xi_t)_{t \in T}$, se numește continuu dacă $\lim_{t \rightarrow 0} M((\xi_t - \xi_0)^2) = 0$. Acest lucru este echivalent cu $\lim_{t \rightarrow t_0} M((\xi_t - \xi_{t_0})^2) = 0$ pentru orice $t_0 \in T$.

Se mai spune că procesul aleator este continuu în medie pătratică.

5. Dacă procesul aleator staționar în sens larg $(\xi_t)_{t \in T}$ este continuu atunci funcția de autocorelație $R(\tau)$ este continuă (exercițiu).

6. Funcția de autocorelare a unui proces staționar în sens larg, continuu, se poate pune sub forma

$$R(\tau) = \int_{-\infty}^{\infty} e^{i\tau\omega} dF(\omega)$$

unde $F(x)$ este mărginită, monoton crescătoare, continuă la stânga. Dacă F admite o derivată, f , ceea ce se întâmplă dacă $|R|$ scade destul de rapid când $|\tau| \rightarrow \infty$, atunci $R(\tau) = \int_{-\infty}^{\infty} e^{i\tau\omega} f(\omega) d\omega$. Demonstrația acestui nu face obiectul acestui curs. $F(\omega)$ se numește funcția spectrală a procesului (ξ_t) iar $f(\omega)$ se numește densitatea spectrală. Din formulele de inversiune ale transformării Fourier găsim că

$$f(\omega) = \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{-i\tau\omega} R(\tau) d\tau$$

Mai mult, ținând seama că $R(\tau)$ este pară se poate arăta că

$$R(\tau) = \int_{-\infty}^{\infty} \cos(\tau\omega) dF(\omega)$$

7. Orice proces aleator continuu staționar în sens larg se poate aproxima prin procese aleatoare standard. Anume, pentru orice $A > 0$ și orice $\varepsilon > 0$ există un număr finit n de

variabile aleatoare $\phi_1, \eta_1, \phi_2, \eta_2, \dots, \phi_n, \eta_n$ independente două câte două și există n numere reale $\lambda_1, \lambda_2, \dots, \lambda_n$ astfel ca pentru orice $t \in [-A, A]$ are loc relația

$$M \left(\left(\xi_t - \sum_{k=1}^n (\phi_k \cos \lambda_k t + \eta_k \sin \lambda_k t) \right)^2 \right) < \varepsilon$$

7.5 Exerciții

1. Care este probabilitatea ca în cazul unui serviciu cu 5 stații de deservire a 20 de unități și cu raportul între coeficienții de venire și deservire $\rho = \frac{\lambda}{\mu} = \frac{1}{5}$, să nu poată fi la un moment dat satisfăcută o cerere? Dar dacă $\rho = \frac{4}{5}$? (o asemenea situație poate fi într-o întreprindere unde o 5 mecanici întrețin 20 de mașini).

Soluție. Avem un proces de așteptare cu $n=5$ stații și $N=20$ unități, deci putem aplica formulele (7.16), pentru $k = n$, de unde

$$p = \frac{\sum_{k=5}^{20} \frac{k!}{n! \cdot n^{k-n}} C_{20}^k \rho^k}{\sum_{k=0}^4 C_{20}^k \rho^k + \sum_{k=5}^{20} \frac{k!}{n! \cdot n^{k-n}} C_{20}^k \rho^k}$$

Cu un computer, pentru $\rho = 1/5$ găsim $p = 0,275$, iar pentru $\rho = 4/5$ găsim $p = 0,9993$.

2. Să presupunem că numărul de stații este n iar solicitările vin dintr-o populație foarte mare, deci numărul de venituri nu este influențat de numărul de unități în așteptare. Care este probabilitatea ca o solicitare să poată fi satisfăcută imediat? Caz particular $n = 5$, $\rho = \frac{\lambda}{\mu} = \frac{1}{2}$?

Soluție. În acest caz avem un proces de naștere și moarte în care $\lambda_k = \lambda$ și în care $\mu_k = k \cdot \mu$ pentru $0 \leq k \leq n$ și $\mu_k = n\mu$ pentru $k \geq n$ (probabilitatea de a ieși o unitate din sistem este proporțională cu numărul unităților în curs de deservire). Găsim, conform cu (7.12)

$$p_k = \begin{cases} \frac{\rho^k}{k!} p_0, & \text{pentru } k \leq n \\ \frac{\rho^k}{n! n^{k-n}} p_0, & \text{pentru } k > n \end{cases}$$

Condiția $\sum p_k = 1$ conduce la

$$p_0 = \frac{1}{\sum_{k=0}^n \frac{\rho^k}{k!} + \frac{n^n}{n!} \sum_{k=n+1}^{\infty} \left(\frac{\rho}{n}\right)^k} = \frac{1}{\sum_{k=0}^n \frac{\rho^k}{k!} + \frac{\rho^{n+1}}{n!(n-\rho)}}$$

În cazul concret al problemei, o solicitare poate fi satisfăcută dacă numărul de unități din sistem este mai mic de n . Probabilitatea căutată este

$$p = \sum_{k=0}^{n-1} p_k = \frac{\sum_{k=0}^{n-1} \frac{\rho^k}{k!}}{\sum_{k=0}^n \frac{\rho^k}{k!} + \frac{\rho^{n+1}}{n!(n-\rho)}}$$

Pentru $n = 5$, $\rho = 1/2$ găsim $p = 0,998$. Dacă venirile sunt mai pronunțate, de exemplu $\rho = \frac{\lambda}{\mu} = 3$ (probabilitatea de solicitare a unui serviciu într-un interval scurt Δt este de trei ori mai mare ca probabilitatea de ieșire de la un server în lucru), atunci găsim o probabilitate de satisfacere imediată a cererii mai mică, $p = \frac{263}{343} = 0,7638$.

3. La o centrală telefonică vin apeluri aleatoare, independente. Centrala poate deservi simultan $n = 50$ de cereri, iar apelurile care vin când centrala este ocupată se anulează. Presupunând că raportul $\rho = \frac{\lambda}{\mu}$ dintre probabilitatea $\lambda\Delta t$ de venire a unui apel într-un timp scurt Δt și probabilitatea $\mu\Delta t$ de eliberare într-un timp scurt Δt a unui circuit ocupat este $\rho = 50$, să se calculeze probabilitatea ca un apel telefonic să fie anulat.

Soluție. Și în acest caz avem un proces de naștere și moarte. Fie p_k probabilitatea ca în centrală să fie k circuite ocupate. Coeficientul λ_k pentru "nașterea unui apel" este fix, egal cu λ . Coeficientul μ_k pentru "moartea unui apel" este $\mu_k = k \cdot \mu$, pentru că dacă sunt k circuite ocupate atunci probabilitatea de eliberare a unui circuit crește de k ori față de cazul unui singur circuit ocupat. Conform cu formulele (7.12) avem

$$\begin{aligned} p_k &= \frac{\lambda_0 \lambda_1 \dots \cdot \lambda_{k-1}}{\mu_1 \mu_2 \dots \cdot \mu_k} p_0 = \frac{\lambda^k}{k! \mu^k} p_0 = \frac{\rho^k}{k!} p_0, \text{ pentru } 0 \leq k \leq n \\ p_k &= 0, \text{ pentru } k > n \end{aligned}$$

Din condiția $\sum_{k=0}^n p_k = 1$ rezultă

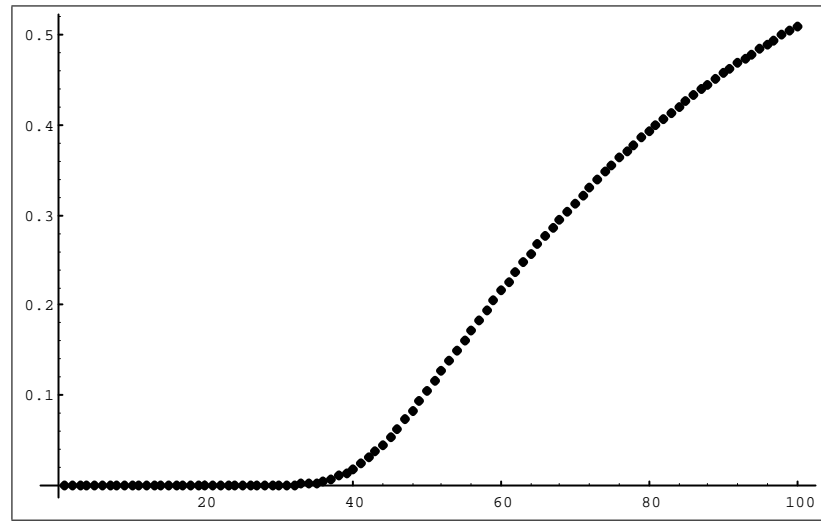
$$p_0 = \frac{1}{\sum_{k=0}^n \frac{\rho^k}{k!}}$$

Probabilitatea cerută în problemă este

$$p = p_n = \frac{\frac{\rho^n}{n!}}{\sum_{k=0}^n \frac{\rho^k}{k!}}$$

unde trebuie luat $n = 50$, $\rho = 50$. Cu calculatorul se obține $p = 0.104787$.

Pe măsură ce solicitarea centralei crește, adică ρ crește, devine din ce în ce mai mare probabilitatea ca un apel să fie anulat. Mai jos apare graficul acestei probabilități în funcție de gradul de încărcare ρ al centralei.



Creșterea probabilității p (verticală) de refuz a unui apel
cu creșterea factorului ρ (orizontală) de încărcare al centralei

4. O particulă se mișcă între doi pereți doar prin pozițiile $x_1 < x_2 < x_3 < x_4$ astfel: dacă este într-o poziție interioară atunci cu probabilitatea $\frac{3}{4}$ face un salt în poziția din față și cu probabilitatea $\frac{1}{4}$ în poziția imediat din spate; dacă se găsește în prima sau în ultima poziție, rămâne definitiv acolo.

Să se descrie comportarea particulei după un număr mare de salturi.

Indicație. Matricea probabilităților de trecere din o poziție în alta este:

$$p = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 1/4 & 0 & 3/4 & 0 \\ 0 & 1/4 & 0 & 3/4 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

Cu calculatorul găsim de exemplu

$$p^{100} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0.30769 & 4.4679 \cdot 10^{-37} & 0 & 0.692308 \\ 0.07692 & 0 & 4.4679 \cdot 10^{-37} & 0.923077 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

Aspectul matricei p^{100} pune în evidență că după un număr mare de salturi particula este captată de un perete sau altul, probabilitățile de a rămâne în poziții intermediare fiind foarte mici. Care este valoarea exactă a matricei p^∞ ?

5. Fie $(\eta_i)_{i \in \mathbb{Z}}$ o familie de variabile aleatoare independente, care au toate repartiția

$$\begin{pmatrix} 1 & -1 \\ \frac{1}{2} & \frac{1}{2} \end{pmatrix}$$

Fie $f : \mathbb{R} \rightarrow \mathbb{R}$ prin $f(t) = \begin{cases} 1, & t \in (-\frac{b}{2}, \frac{b}{2}] \\ 0, & \text{în caz contrar} \end{cases}$. Definim acum procesul aleator $(\xi_t)_{t \in \mathbb{R}}$ prin $\xi_t = \sum_{i \in \mathbb{Z}} f(t - i \cdot b) \eta_i$. Să se calculeze funcția de autocorelație $R(t_1, t_2)$. Este procesul staționar în sens larg?

Soluție. Pe intervalul $(-\frac{b}{2}, \frac{b}{2}]$, $\xi_t = \eta_0$; pe intervalul $(\frac{b}{2}, \frac{3b}{2}]$, $\xi_t = \eta_1$; ... pe intervalul $I_i = (ib - \frac{b}{2}, ib + \frac{b}{2}]$, $\xi_t = \eta_i$. Prin urmare $\xi_{t_1} \xi_{t_2} = \eta_i \eta_j$ dacă $t_1 \in I_i$ și $t_2 \in I_j$. Deoarece variabilele η_i sunt independente rezultă că $M(\xi_{t_1} \xi_{t_2}) = M(\eta_i \eta_j)$ este diferită de zero doar dacă $i = j$, adică $t_1, t_2 \in I_i$, caz în care avem $R(t_1, t_2) = M(\eta_i^2) = 1$. Funcția de autocorelație nu este invariantă la o translație în timp deci procesul nu este staționar în sens larg.

6. Fie procesul aleator $\xi_t = \sum_{k=1}^n a_k (\phi_k \cos \lambda_k t + \eta_k \sin \lambda_k t)$ unde a_k și $\lambda_k \in \mathbb{R}$, iar ϕ_k și η_k sunt variabile aleatoare independente, de medie 0 și dispersie 1, pentru $k=1, 2, \dots, n$. Să se arate că procesul este staționar în sens larg și să se determine autocorelația procesului.

Răspuns. $R(\tau) = \sum a_k^2 \cos \lambda_k \tau$.

Partea II

Statistică

Lecția 8

Statistica descriptivă

În cele ce urmează vom încerca să explicăm ce este statistica, cum diferă ea de teoria probabilităților, ce o leagă de aceasta, care sunt părțile ei componente și cum începe demersul practic într-o problemă de statistică (adică vom spune câteva cuvinte despre statistica descriptivă).

Atunci când omul nu a mai putut *intui* a început să *măsoare*. Măsurătorile și observațiile au devenit prima treaptă spre înțelegerea legilor *naturii*. Dar, în acest fel, omul nu mai poate să cunoască direct *realitatea*, el poate numai să o *aproximeze* succesiv prin *modele fizice* și apoi prin *modele matematice*. Dar aceste modele nu descriu *exact* Realitatea. Ele o aproximează și apar așa numitele erori. Unele erori sunt previzibile, altele însă sunt întâmplătoare (aleatoare). Aceste ultime erori (aleatoare) au și ele legile lor de manifestare. Apar deci fenomenele aleatoare descrise prin variabilele aleatoare. Teoria probabilităților pleacă de la ipoteza că se cunosc exact aceste variabile aleatoare (prin funcțiile de probabilitate, prin funcțiile de repartiție, prin funcțiile caracteristice, etc.). *Statistica* pleacă de la măsurătorile brute și caută să regăsească modelul probabilistic teoretic exact care se află în spatele acestor măsurători. Partea "empirică" a statisticii care se ocupă de prelucrarea datelor obținute prin măsurători sau observații se numește *statistică descriptivă*. Aparatul matematic al teoriei probabilităților, pus în funcțiune pentru a studia și interpreta aceste date, în dorința de a recupera modelul probabilistic real, care guvernează fenomenul măsurat sau observat, formează *inferența statistică*. După ce cercetătorul capătă informații suficient de clare despre fenomenul probabilistic studiat, el va trebui să *actioneze* optim potrivit acestor informații. Apare deci *teoria deciziei statistice*, care este o ramură importantă a statisticii.

8.1 Statistica unei variabile

Multimea de obiecte studiată se numește *populație*. Un obiect separat dintr-o populație dată se numește *individ* sau *membru al populației*. Trăsătura comună a tuturor membrilor populației care ne interesează în studiul nostru se numește *caracteristică*. Caracteristicile pot fi *cantitative* (înălțime, greutate, notă la examen, abscisa unui punct în plan, etc...) sau

calitative (culoarea ochilor, sex, loc de nastere, etc...). Oricum statistica lucrează cu numere, caracteristicilor calitative li se atasază *coduri numerice*.

Exemplul 8.1 *Ne interesează statistica ploilor în Bucuresti pe anul 1995, zilnic. Aici populația este multimea zilelor din anul 1995, un individ al populației este o zi anume din acest an, de exemplu 3 ianuarie, iar caracteristica calitativă este faptul că a plouat sau nu în acea zi. Dacă a plouat punem 1 și dacă nu, putem 0. Numerele 1 și 0 reprezintă coduri în statistica respectivă.*

Presupunem în continuare că avem numai caracteristici cantitative ale unor populații, mai exact avem multimi brute de numere reale, sau tabele de numere reale. Privim aceste numere atasate unei populații ca fiind valori ale unei variabile aleatoare X . Vom spune pe scurt: "fie populația X ".

Exemplul 8.2 *O mașină produce piese cilindrice fine, cu diametru standard fixat $\phi = 3$ cm. Fiecare piesă are o abatere de la acest diametru, măsurată în microni. Aceste abateri formează o "populație" în sensul de mai sus, mai bine zis valorile unei variabile aleatoare X . Noi nu putem să precizăm de la început ce abatere va avea o piesă luată la întâmplare, dar putem face o selecție de n piese și putem măsura abaterile lor: x_1, x_2, \dots, x_n . Fiecare x_i reprezintă o valoare a v.a. X care, teoretic vorbind, are o densitate de probabilitate $\rho(x)$ și o funcție de repartiție $F(x)$.*

Definiția 8.3 *O mulțime de n observații independente asupra unei caracteristici numerice X a unei populații P , care ne dă n valori x_1, x_2, \dots, x_n , se numește selecție de volum n . Șirul de valori $(x_i)_{1 \leq i \leq n}$ îl vom numi serie statistică discretă.*

În exemplul de mai sus facem o selecție de volum n din multimea pieselor și construim așa numita funcție de repartiție empirică $F_n^*(x)$.

Definiția 8.4 *Se numește funcție de repartiție empirică asociată unei variabile aleatoare X și unei selecții $\{x_1, x_2, \dots, x_n\}$, funcția $F_n^* : R \rightarrow R$*

$$F_n^*(x) = \frac{\text{nr. de valori } x_j < x}{n} = \frac{k_x}{n}$$

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

Teorema 8.5 *Fie P o populație statistică și X variabila aleatoare atasată ei cu funcția de repartiție $F(x)$. Pentru o selecție de volum n : $\{x_1, x_2, \dots, x_n\}$ construim ca mai sus funcția de repartiție empirică $F_n^*(x)$. Atunci*

$$\Pr ob \{ |F(x) - F_n^*(x)| \geq \epsilon \} \rightarrow 0$$

când $n \rightarrow \infty$, pentru orice $\epsilon > 0$, fixat. Altfel spus $F_n^*(x) \rightarrow F(x)$ în probabilitate.

Demonstratie Să notăm cu $p = \text{Prob}\{X < x\} = F(x)$, si cu $F_n^*(x) = \frac{k_x}{n}$ (vezi definitia de mai sus). Notăm cu η_1, \dots, η_n v.a. construite astfel: η_j are valoarea 1 dacă $x_j < x$ si 0 în caz contrar. Variabilele η_1, \dots, η_n sunt independente (ca valori ale unor observații independente) și au distribuția

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

Avem $M(\eta_i) = p$, $D(\eta_i) = p(1-p)$. Este clar că v.a. $Y_n = \frac{\eta_1 + \dots + \eta_n}{n}$ are media p și dispersia

$$D(Y_n) = \frac{1}{n^2}(D(\eta_1) + \dots + D(\eta_n)) = \frac{p(1-p)}{n}$$

(a se vedea proprietățile mediei și dispersiei, Lecția 2). Aplicăm acum inegalitatea lui Cebîșev lui Y_n si găsim că

$$\begin{aligned} & \text{Prob}(|F_n^*(x) - F(x)| \geq \epsilon) \\ &= \text{Prob}(|Y_n - p| \geq \epsilon) \leq \frac{D(Y_n)}{\epsilon^2} = \frac{p^2(1-p)^2}{n\epsilon^2} \end{aligned}$$

Cum partea dreaptă tinde la 0 când $n \rightarrow \infty$ rezultă că $F_n^*(x) \rightarrow F(x)$ în probabilitate, când $n \rightarrow \infty$.

QED.

În urma oricărei selecții de volum n dintr-o populație de numere se obține un sir 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_k , după criterii mai mult sau mai puțin subiective. Asociem fiecărui interval I_j mijlocul lui, M_j . Punctului M_j îi asociem *frecvența relativă* a v.a. empirice pe intervalul I_j , adică câtul dintre numărul n_j al acelor x_i care se află în I_j și n (volumul întregii selecții): n_j/n . Este clar că în felul acesta obținem o v.a. $\tilde{X}_n \left(\frac{\tilde{x}_j}{n_j/n} \right) j = 1, 2, \dots, k$,

unde \tilde{x}_j este abscisa punctului M_j . Graficul funcției de probabilitate al v.a. \tilde{X}_n se numește *histogramă* asociată selecției x_1, \dots, x_n și împărțirii în intervale I_1, \dots, I_k . Dacă unim printr-o linie poligonală punctele de coordonate $(\tilde{x}_j, n_j/n)$ obținem un *poligon al frecvențelor* ce aproximează de fapt graficul funcției densitate de probabilitate al v.a. X pe un interval finit $(I_1 \cup I_2 \cup \dots \cup I_k)$ care conține numerele x_1, \dots, x_n .

Pentru o selecție dată $\{x_1, \dots, x_n\}$ se introduc diferiți *indicatori empirici* care dau anumite informații despre întreaga populație.

Definiția 8.6 Fie $\{x_1, x_2, \dots, x_n\}$ o selecție de volum n .

i) $m^* = \frac{x_1 + x_2 + \dots + x_n}{n}$ se numește *media empirică*.

ii) $m_r^* = \frac{x_1^r + x_2^r + \dots + x_n^r}{n}$ se numește *momentul empiric de ordinul r* .

iii) $\mu_k^* = \frac{(x_1 - m^*)^k + (x_2 - m^*)^k + \dots + (x_n - m^*)^k}{n}$ se numește *momentul empiric centrat de ordin k* .

iv) $S^{*2} = D^* = \sigma^{*2} = \frac{(x_1 - m^*)^2 + (x_2 - m^*)^2 + \dots + (x_n - m^*)^2}{n}$ se numește dispersia empirică sau varianta empirică. σ^* se numește deviația standard.

v) $S'^{*2} = \frac{(x_1 - m^*)^2 + (x_2 - m^*)^2 + \dots + (x_n - m^*)^2}{n-1}$ se numește dispersia empirică modificată.

vi) Valoarea $\alpha \in R$ astfel ca numărul de valori $x_i \leq \alpha$ este egal cu numărul de valori $x_i \geq \alpha$, se numește mediană. Dacă există mai multe asemenea valori pentru α , atunci ele formează un interval și mediana este prin definiție mijlocul acestui interval.

vii) Valoarea x_i cu frecvența maximă de apariție se numește modul selecției. (este posibil să nu fie unic)

viii) Se numește prima cvartilă a selecției, cel mai mic x astfel ca numărul de valori $x_j \leq x$ să fie $\geq \frac{1}{4}n$. A treia cvartilă este cea mai mică valoare x_i astfel ca numărul de valori $x_j \leq x_i$ să fie $\geq \frac{3}{4}n$. Analog se definește a p -a cuantilă de ordin q ca cea mai mică valoare x_i astfel ca numărul de valori $x_j \leq x_i$ să fie $\geq \frac{p}{q}n$.

Observația 8.7 • În cazul când datele sunt grupate pe intervale, definițiile de mai sus se referă la mijloacele intervalelor, fiecare mijloc fiind considerat de atâtea ori câte valori se află în el.

• În general dacă o valoare x_i se repetă atunci vom nota cu n_i numărul de apariții, și cu $f_i = \frac{n_i}{n}$ frecvența relativă. Formulele de mai sus pot fi scrise $m^* = \frac{\sum n_i x_i}{n} = \sum f_i x_i$, $\sigma^{*2} = \sum f_i (x_i - m^*)^2$, etc. Insumarea se face acum numai după valorile x_i distincte. Seria statistică o vom nota în acest caz $(x_i, n_i)_{1 \leq i \leq p}$, punând în evidență de câte ori apare fiecare valoare.

• Media și mediana descriu "centrul" valorilor de selecție iar dispersia este o măsură a împrăstierii acestor valori în jurul "centrului". Modul indică în ce zonă sunt cele mai probabile valori. Cuantilele indică în ce zone se află un anumit procent de valori.

Propoziția 8.8 Următoarele formule au loc:

$$S^{*2} = \frac{n \sum x_i^2 - (\sum x_i)^2}{n^2} \quad (8.1)$$

$$S'^{*2} = \frac{n \sum x_i^2 - (\sum x_i)^2}{n(n-1)} \quad (8.2)$$

$$S'^{*2} = \frac{\sum x_i^2}{n-1} - \left(\frac{n}{n-1} \right) m^{*2} \quad (8.3)$$

Demonstrație. Sunt calcule simple lăsate ca exercițiu.

Observația 8.9 În general în calcule nu utilizăm notațiile m^*, D^* , etc. ci m, D, \dots Am introdus aici notațiile m^*, D^*, \dots pentru a le distinge de m =media teoretică, D =dispersia teoretică, etc., care se vor introduce în lecția următoare.

Exemplul 8.10 O firmă este interesată de timpul mediu al convorbirilor telefonice și de distribuția acestor timpi față de timpul mediu (dispersia) pe durata a 40 convorbiri telefonice consecutive. Timpii s-au rotunjit în minute și rezultatul sondajului a dat următorii timpi: 4, 6, 4, 4, 7, 2, 3, 1, 2, 1, 1, 4, 9, 8, 11, 12, 3, 2, 1, 1, 3, 9, 4, 5, 7, 7, 9, 10, 10, 1, 2, 2, 3, 11, 12, 10, 1, 1, 3, 4. Să se facă și o histogramă a frecvențelor relative și un grafic al funcției de repartiție pentru acest sondaj.

Soluție Facem mai întâi următorul tabel :

timpul de convorbire t_i	numărul convorbirilor n_i	frecv. relativă $f_i = n_i/n$	frecv. cumulată $F_{(t_i)}$
1 min	8	8/40	8/40
2 min	5	5/40	$\frac{8}{40} + \frac{5}{40} = \frac{13}{40}$
3 min	5	5/40	$\frac{8}{40} + \frac{5}{40} + \frac{5}{40} = \frac{18}{40}$
4 min	6	6/40	24/40
5 min	1	1/40	25/40
6 min	1	1/40	26/40
7 min	3	3/40	29/40
8 min	1	1/40	30/40
9 min	3	3/40	33/40
10 min	3	3/40	36/40
11 min	2	2/40	38/40
12 min	2	2/40	40/40 = 1

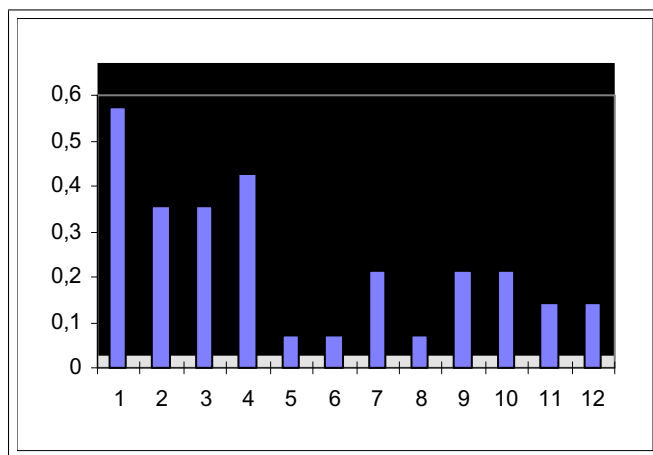
- media convorbirilor este

$$\begin{aligned}
 m^* &= \sum t_i \cdot n_i/n = \frac{1}{n} \sum t_i \cdot f_i \\
 &= \frac{1}{40} (1 \cdot 8 + 2 \cdot 5 + 3 \cdot 5 + 4 \cdot 6 + 5 \cdot 1 + 6 \cdot 1 \\
 &\quad + 7 \cdot 3 + 8 \cdot 1 + 9 \cdot 3 + 10 \cdot 3 + 11 \cdot 2 + 12 \cdot 2) \\
 &= 5
 \end{aligned}$$

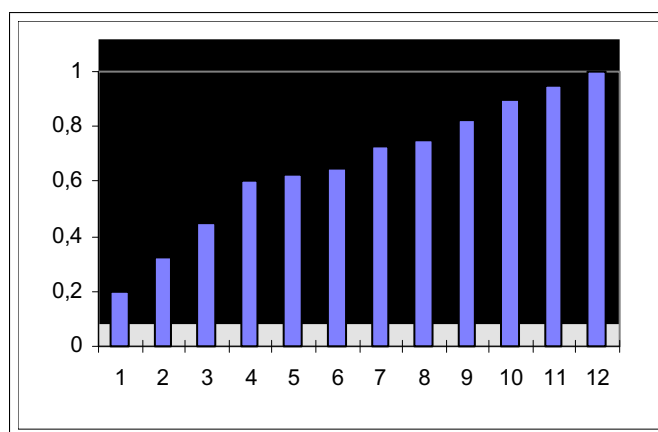
- dispersia empirică este

$$\begin{aligned}
 &\sum (t_i - m^*)^2 \cdot n_i/n \\
 &= \left(\frac{1}{40} (1 - m^*)^2 \cdot 8 + (2 - m^*)^2 \cdot 5 + (3 - m^*)^2 \cdot 5 + \dots + (12 - m^*)^2 \cdot 2 \right) \\
 &= 13,179
 \end{aligned}$$

- mediana este 4, prima cvartilă este 2, a treia cvartilă este 8,25, modul este 1.
- histograma frecvențelor este



- *histograma frecvențelor cumulate este:*



Exemplul 8.11 Se dă o selecție de 150 de numere $\{x_1, x_2, \dots, x_{150}\}$ cu media de selecție $m=102,42$. Aceste numere se grupează în 8 intervale $[81, 5; 87, 5)$, $[87, 5; 93, 5)$, ..., $[123, 5; 129, 5)$, de lungime 6 unități. Ele se repartizează în aceste intervale după cum urmează: în primul interval avem 2 numere ($n_1=2$), în al doilea 23 de numere ($n_2=23$), $f_3=22$, $n_4=65$, $n_5=20$, $n_6=10$, $n_7=0$, $n_8=8$. a) Să se calculeze media selecției. b) Să se calculeze dispersia selecției.

Soluție a) este lăsat ca exercițiu; se găsește $m^* = 102,42$.

b) Se face următorul tabel de calcule:

x_j (mijlocul int.)	n_j	$x_j - m^*$	$(x_j - m^*)^2$	$(x_j - m^*)^2 n_j$
84,5	2	-17,92	321,1264	642,2528
90,5	23	-11,92	142,0864	3267,9872
96,5	22	-5,92	35,0464	771,0208
102,5	65	0,08	0,0064	0,4160
108,5	20	6,08	36,9664	739,3280
114,5	10	12,08	145,9264	1459,2640
120,5	0	18,08	326,8864	0,0006
126,5	8	24,08	579,8464	4638,7712

				11519,0400

Găsim $S^{*2} = \frac{\sum (x_j - m)^2 n_j}{n} = \frac{11519,04}{150} = 76,79$. Pentru verificare putem folosi formula $S^{*2} = \frac{\sum x_j^2 n_j}{n} - m^{*2}$ care este mai comodă, dar cere o coloană separată cu calculul lui x_j^2 .

8.2 Statistica a două variabile

Să presupunem că avem două caracteristici numerice care se urmăresc, de exemplu înălțimea și greutatea. Prin testare se găsește următoarea situație: x_i sunt greutățile, y_j sunt înălțimile observate (grupate pe intervale), iar la întretăierea coloanei i cu linia j se află numărul de cazuri observate, $n_{i,j}$.

$\frac{x_i \rightarrow}{y_j \downarrow}$	43	48	53	58	$n_{.,j}$
152	20	8	2	0	30
157	2	18	1	4	25
162	0	1	10	4	15
167	0	1	4	15	20
$n_{i,.}$	22	28	16	23	N=80

Notăm o asemenea serie de observații prin $(x_i, y_j, n_{i,j})_{\substack{1 \leq i \leq p \\ 1 \leq j \leq q}}$. Avem de exemplu la $x_2=48$ și $y_1=152$ un număr de $n_{2,1}=8$ cazuri înregistrate.

Se definesc următoarele mărimi:

i) $n_{i,.} = \sum_j n_{i,j}$, $n_{.,j} = \sum_i n_{i,j}$, $N = \sum_{i,j} n_{i,j}$. Seria $(x_i, n_{i,.})$ se numește seria marginală în x , iar seria $(y_j, n_{.,j})$ se numește seria marginală în y . $f_{i,.} = \frac{n_{i,.}}{N}$ și $f_{.,j} = \frac{n_{.,j}}{N}$ se numesc frecvențe marginale, iar $f_{i,j} = \frac{n_{i,j}}{N}$ se numește frecvența dublă.

ii) $m_x^* = \frac{\sum_{i,j} n_{i,j} x_i}{N} = \frac{\sum_i n_{i,.} x_i}{N}$ și $m_y^* = \frac{\sum_{i,j} n_{i,j} y_j}{N} = \frac{\sum_j n_{.,j} y_j}{N}$ se numesc medii marginale.

ii)

$$\sigma_x^{*2} = \frac{\sum_{i,j} n_{i,j} (x_i - m_x^*)^2}{N} = \frac{\sum_i n_{i,.} (x_i - m_x^*)^2}{N} = \frac{\sum_i n_{i,.} x_i^2}{N} - m_x^{*2}$$

și

$$\sigma_y^{*2} = \frac{\sum_{i,j} n_{i,j} (y_j - m_y^*)^2}{N} = \frac{\sum_j n_{.,j} (y_j - m_y^*)^2}{N} = \frac{\sum_j n_{.,j} y_j^2}{N} - m_y^{*2}$$

se numesc dispersii (varianțe) marginale.

iv) Covarianța seriei este numărul

$$\text{cov}(x, y) = \frac{\sum_{i,j} n_{i,j} (x_i - m_x^*) (y_j - m_y^*)}{N} = \frac{\sum_{i,j} n_{i,j} x_i y_j}{N} - m_x^* m_y^*.$$

v) Coeficientul de corelație liniară al seriei este $\rho_{x,y} = \frac{\text{cov}(x,y)}{\sigma_x^* \sigma_y^*}$.

În cazul de mai sus găsim

$$m_x^* = \frac{22 \cdot 43 + 28 \cdot 48 + 16 \cdot 53 + 23 \cdot 58}{80} = 49,438;$$

$$m_y^* = \frac{30 \cdot 152 + 25 \cdot 157 + 15 \cdot 162 + 20 \cdot 167}{80} = 157,313;$$

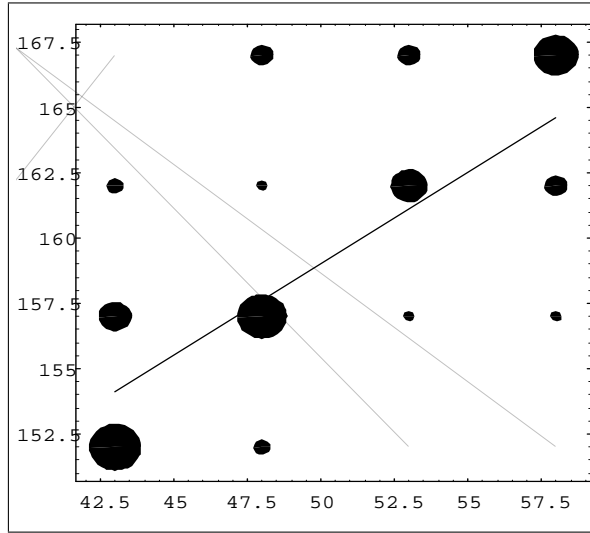
$$\sigma_x^{*2} = \frac{22 \cdot (43 - 49,438)^2 + 28 \cdot (48 - 49,438)^2 + 16 \cdot (53 - 49,438)^2 + 23 \cdot (58 - 49,438)^2}{80} = 28,246;$$

$$\sigma_y^{*2} = \frac{30 \cdot (152 - 157,313)^2 + 25 \cdot (157 - 157,313)^2 + 15 \cdot (162 - 157,313)^2 + 20 \cdot (167 - 157,313)^2}{80} = 26,465;$$

$$\text{cov}(x, y) = \frac{20 \cdot (43 - 49,438) \cdot (152 - 157,313) + 8 \cdot (48 - 49,38) \cdot (152 - 157,313) + \dots + 15 \cdot (58 - 49,438) \cdot (167 - 157,313)}{80} = 18,926;$$

$$\rho_{x,y} = \frac{18,926}{\sqrt{28,246 \cdot 26,465}} = 0,692.$$

Reprezentarea grafică a datelor se face prin discuri pline: în punctul (x_i, y_j) se pune un disc cu aria proporțională cu numărul de observații care au dat greutatea x_i și înălțimea y_j . Se obține histograma:



(dreapta nu face parte din histogramă; vezi în continuare)

Două selectii de același volum n din două populații diferite, $\{x_1, \dots, x_n\}$ și $\{y_1, \dots, y_n\}$ se zic corelate prin funcția $y = f(x)$ dacă $y_k = f(x_k)$, pentru $k = 1, 2, \dots, n$.

Dacă $f(x) = ax + b$, corelația se zice *liniară*. Am văzut în lecția 6 că $\rho_{x,y} = \pm 1$ este echivalent cu faptul că punctele (x_i, y_j) sunt de-a lungul unei drepte. Dacă $|\rho_{x,y}|$ este apropiat de 1 atunci datele (x_i, y_j) sunt aproximativ pe o dreaptă $y = ax + b$. Reluăm aici, în varianta folosită în aplicații acest lucru.

Teorema 8.12 Fie $(x_i, y_j, n_{i,j})_{\substack{1 \leq i \leq p \\ 1 \leq j \leq q}}$ o selecție dublă. Atunci:

a) $-1 \leq \rho_{x,y} \leq 1$ și semnul egal apare dacă și numai dacă punctele (x_i, y_j) , pentru $n_{i,j} \neq 0$, sunt coliniare.

b) Ecuația dreptei $y = ax + b$, unde coeficienții a, b sunt determinați de condiția ca expresia

$$\Phi(a, b) = \sum_{(x_i, y_j) = \text{observat}} (y_j - ax_i - b)^2 = \sum_{i,j} n_{i,j} (y_j - ax_i - b)^2$$

să fie minimă, este:

$$y = \frac{\text{cov}(x, y)}{\sigma_x^{*2}} (x - m_x^*) + m_y^* \quad (8.4)$$

Această dreaptă se numește dreapta de regresie a lui y în x .

Demonstrație.

a) Deoarece $n_{i,j} \geq 0$, atunci expresia $E = \sum_{i,j} n_{i,j} (t(y_j - m_y^*) + (x_i - m_x^*))^2$ este pozitivă pentru orice $t \in \mathbb{R}$. Ridicând la pătrat, găsim: $t^2 \sigma_y^{*2} + 2\text{cov}(x, y)t + \sigma_x^{*2} \geq 0$ pentru orice $t \in \mathbb{R}$. Prin urmare $\Delta = 4\text{cov}^2(x, y) - 4\sigma_y^{*2}\sigma_x^{*2} \leq 0$. Prin împărțire cu $\sigma_x^{*2}\sigma_y^{*2}$ găsim $\rho_{x,y}^2 \leq 1$, adică $-1 \leq \rho_{x,y} \leq 1$. Dacă $\rho = \pm 1$ atunci $\Delta = 0$ deci există t_0 astfel ca $E = 0$, deci fiecare paranteză

este egală cu 0, deci pentru orice i, j , pentru care $n_{i,j} \neq 0$ avem $t_0 (y_j - m_y^*) + x_i - m_x^* = 0$, deci punctele (x_i, y_j) cu $n_{i,j} \neq 0$ sunt coliniare.

b) $\frac{\partial \Phi(a,b)}{\partial a} = 0$, $\frac{\partial \Phi(a,b)}{\partial b} = 0$ formează un sistem liniar în a și b cu soluțiile $a = \frac{cov(x,y)}{\sigma_x^{*2}}$, $b = m_y^* - am_x^*$.

QED.

Analog se determină dreapta de regresie a lui x în y . Cele două drepte sunt distincte. Ele coincid doar dacă datele (x_i, y_j) sunt coliniare. In cazul de mai sus găsim $y = 0,67x + 124,188$, dreaptă care este reprezentată pe histograma datelor.

Dacă $f(x) = ax^2 + bx + c$, corelatia se zice *parabolică*. Coeficienții se determină din condiția ca $\Phi(a, b, c) = \sum_{i,j} n_{i,j} (y_j - ax_i^2 - bx_i - c)^2$ să fie minimă.

Dacă $f(x) = ae^{bx}$, corelatia se zice *exponentială* și se reduce prin logaritmare tot la o corelație liniară: $\ln(f(x)) = \ln(a) + bx$. Coeficienții $\alpha = \ln(a)$ și b se determină din condiția ca $\Phi(\alpha, b) = \sum_{i,j} n_{i,j} (\ln(y_j) - \alpha - bx_i)^2$ să fie minimă.

Dacă $f(x) = ax^b$, atunci avem prin logaritmare $\ln(f(x)) = \ln(a) + b \ln(x)$, și, la fel ca mai sus se deduc coeficienții $\alpha = \ln(a)$ și b din condiția $\Phi(\alpha, b) = \sum_{i,j} n_{i,j} (\ln(y_j) - \alpha - b \ln(x_i))^2$ să fie minimă.

In multe situații pentru fiecare x_i avem doar o valoare y pe care o notăm y_i , deci valorile (x_i, y_i) sunt pe graficul unei funcții. Determinarea unei funcții care ajustează datele respective prin metoda celor mai mici pătrate constă în propunerea unui model de funcție, $f(x, a, b, ..)$, și determinarea parametrilor $a, b, ..$ din condiția $\Phi(a, b, ..) = \sum_i (y_i - f(x_i, a, b, ..))^2$ să fie minimă.

8.3 Exerciții

1. S-a făcut un sondaj preelectoral pe un esantion de 100 persoane. Am notat cu A, B, C, D, E candidatii, cu F răspunsul "nedecis" și cu G răspunsul "nu intentionez să votez". Să se construiască o histogramă cu funcția de distribuție și altă histogramă cu funcția de repartiție (frecvența cumulată) pentru acest sondaj dacă răspunsurile sunt date în următorul tabel: C, A, A, B, E, F, F, C, C, C, A, B, A, A, A, E, F, A, B, G, D, B, B, C, F, G, G, D, D, D, B, A, B, B, B, F, G, B, C, A, E, C, C, D, G, A, A, E, E, E, C, D, D, E, G, G, A, B, B, A, F, F, G, G, G, G, A, A, A, B, B, C, C, A, A, D, D, E, F, G, A, B, C, C, D, A, E, F, A, B, F, G, A, B, C, D, A, A, B, E.

Soluție Aici trebuie mai întâi să codificăm numeric literele (opțiunile electoratului) A, B, C, D, F, G. De exemplu, propunem următoarea codificare:

$G \longleftrightarrow 0$

$F \longleftrightarrow 1$

$A \longleftrightarrow 4$

$B \longleftrightarrow 5$

$C \longleftrightarrow 6$

$D \longleftrightarrow 7$

2. Două grupe de 10 studenți A și B au obținut următoarele note la examenul de statistică:

A : 8, 5, 6, 6, 7, 9, 4, 3, 5, 6

B : 9, 6, 7, 8, 6, 10, 5, 4, 6, 7.

Să se găsească cea mai bună corelație liniară între cele două selecții. Să se găsească valoarea deviației pătratice. Să se facă același lucru pentru o corelație de tip parabolic și să se compare deviațiile pătratice.

3. Fie selecția $\{0, 1, -1, -1, -2, 1, 1, -1, 2, 3, 1, 4, 3, -1, 0, 0, 3, -1, -2, -2\}$ dintr-o populație anume. Fie X v.a. care guvernează populația. Să se aproximeze cu ajutorul selecției numărul $P(0 \leq X \leq 2)$. Se cere graficul funcției de repartiție pentru această selecție și o histogramă a frecvențelor.

4. S-a făcut un sondaj asupra pretului (în centi) galonului de benzină premium asupra a 30 stații luate la întâmplare. De aici a rezultat selecția: 65, 58, 64, 68, 52, 48, 59, 59, 56, 63, 61, 66, 52, 57, 60, 62, 55, 55, 64, 71, 61, 63, 46, 53, 60, 57, 58, 57, 54, 58. Se cere graficul poligonului de frecvență (relativă) dacă: a) grupăm datele în intervale de lungime 3, cu 60 centrul unui asemenea interval; b) grupăm datele în intervale de lungime 5, cu 60 ca centru al unui asemenea interval. Calculați pentru aceste grupări media și dispersia de selecție. Găsiți graficul funcțiilor de repartiție empirice.

5. La un concurs 12 studenți au obținut următoarele punctaje: 18, 15, 19, 27, 13, 30, 24, 11, 5, 16, 17, 20. Calculați media, mediana, deviația standard și deviația absolută medie. Construiți funcția empirică de frecvență cumulată (f. de repartiție) și interpretați rezultatele obținute.

Lecția 9

Statistici. Estimarea parametrilor

Amintim că o populație P este o mulțime de obiecte din care se fac selecții finite (de volum $n < \infty$). Populația se poate identifica cu mulțimea tuturor observațiilor potențiale pe care le putem face asupra obiectelor ei. Pentru fiecare obiect al selecției se testează valoarea unei caracteristici numerice, X . Admitem că pe P există o probabilitate și că X este o variabilă aleatoare. Distribuția (funcția de repartiție) a *v.a.* X se numește *distribuția populației după caracteristica X* .

Exemplul 9.1 *Intr-o magazie grâul este amestecat cu neghină. Populația P este aici totalitatea bobelor din magazie (câteva sute de milioane). Fie $X : P \rightarrow R$,*

$$X(\text{bob}) = \begin{cases} 1 & \text{daca e grau} \\ 0 & \text{daca e neghina} \end{cases}$$

Probabilitatea p ca un bob să fie de tip A este definită prin:

$$p(A) = \frac{\text{nr. de boabe de tipul } A}{\text{nr. de boabe din magazie}}, \quad A \subset P$$

Aici p nu se poate determina experimental exact din cauza numărului mare de boabe, dar teoretic p există. Valoarea medie a lui X înmulțită cu 100 este procentul de boabe de grâu din magazie, lucru important.

Exemplul 9.2 *Să presupunem că mai multe persoane, sau aceeași persoană în mai multe rânduri, măsoară independent o lungime, de aproximativ 1 km folosind o ruletă de 2 m. Evident că se vor obține rezultate diferite datorită unei game largi de cauze incontrolabile. Putem în acest caz considera P ca mulțimea tuturor complexelor de cauze necontrolabile care influențează rezultatul măsurătorii sau putem considera P ca mulțimea tuturor măsurătorilor posibile. Oricum P nu este o mulțime pe care o putem explicita ca în cazul precedent. Admitem însă că pe P există o probabilitate iar o măsurătoare înseamnă o manifestare a unui complex ω de cauze necontrolabile care conduc la un rezultat $X(\omega)$, în cazul nostru X fiind o lungime.*

Prin urmare caracteristica lungime apare ca o funcție $X : P \rightarrow R$. Admitem că X este o variabilă aleatoare, adică (vezi lecțiile 2, 3) $\{\omega \in P \mid X(\omega) < r\} \subset P$ este o mulțime pe care este definită probabilitatea p .

Statistica Matematică se ocupă, printre altele, cu problema determinării repartiției unei variabile aleatoare X ca în exemplele de mai sus, prin experimente. În general n experimente conduc la n valori numerice x_1, \dots, x_n . Ce operații trebuie făcute cu valorile x_1, \dots, x_n pentru a găsi caracteristici ale lui X și ce încredere putem avea în rezultatele obținute?

În continuare prezentăm felul în care putem considera rezultatele x_1, x_2, \dots, x_n ale lui X în n experiențe independente ca valori a n variabile aleatoare independente X_1, X_2, \dots, X_n . La o primă lectură se poate sări peste această parte, remarcându-se doar concluziile.

Fie P ca mai înainte spațiul probabilizat al cauzelor incontolabile, fie $\Omega \subset P(P)$ σ -algebra submulțimilor lui P pentru care e definită σ probabilitatea p . Notăm P^∞ șirurile de elemente din P . Deci $\omega \in P^\infty$ dacă și numai dacă $\omega = (\omega_k)_{k \in \mathbb{N}}$ și pentru orice k , $\omega_k \in P$. Următoarele submulțimi ale lui P^∞ :

$$\begin{aligned} A &= A_1 \times A_2 \times \dots \times A_n \times P \times P \times \dots \\ &= \{(\omega_k)_{k \in \mathbb{N}} \mid \omega_k \in A_k \text{ pentru } 1 \leq k \leq n\} \end{aligned} \quad (9.1)$$

unde $A_k \in \Omega$ pentru orice k , se numesc paralelipede. Aici n nu este fixat ci poate fi orice număr natural. Fie Ω^∞ submulțimile lui P^∞ care sunt reuniuni finite de paralelipede. Se arată că aceste mulțimi formează o algebră. Pe această algebră putem defini o unică probabilitate p' astfel ca pentru paralelipede să avem:

$$p'(A) = p(A_1) \cdot p(A_2) \cdot \dots \cdot p(A_n) \quad (9.2)$$

unde p este probabilitatea pe P . Definiția seamănă cu definiția volumului unui paraleliped în funcție de lungimile laturilor sale. Asemenea probabilități nu sunt suficiente pentru nevoile de calcul. E nevoie de o proprietate de continuitate de genul: $B_1 \subset B_2 \subset \dots \subset B_k \subset \dots$ cu $B = \bigcup_{k=1, \infty} B_k$ implică $p(B) = \lim_{n \rightarrow \infty} p(B_k)$. Construcția unei astfel de probabilități pe P^∞ se realizează astfel:

a) Se extinde Ω^∞ la cea mai mică σ algebră (deci algebră de mulțimi închisă și la reuniuni numărabile) notată $\Omega^{(\infty)}$

b) Probabilitatea p' definită pe Ω^∞ se extinde unic la o σ probabilitate pe $\Omega^{(\infty)}$ notată $p^{(\infty)}$

Probabilitatea $p^{(\infty)}$ se numește probabilitate produs. Detaliile de construcție nu fac obiectul acestui curs. Putem remarca asemănarea construcției probabilității produs cu a volumului corpurilor plecând de la lungime. Așa cum în afară de reuniuni finite de paralelipede există și alte corpuri cu volum, tot așa apar în $\Omega^{(\infty)}$ și alte mulțimi care au probabilitate, în afară de reuniunile finite de tipul (9.1).

Definiția 9.3 Mulțimea P^∞ împreună cu $\Omega^{(\infty)} \subset P(P^\infty)$ și cu probabilitatea $p^{(\infty)} : \Omega^{(\infty)} \rightarrow R$ se numește produsul infinit al câmpului de probabilitate (P, Ω, p) .

Observația 9.4 *In lecția 2 am introdus produsul finit al unor câmpuri de probabilitate. Față de cazul considerat acolo, aici avem două lucruri în plus:*

- a) *pentru a avea o σ probabilitate pe produs trebuie extinsă algebra de mulțimi formată din reuniuni finite de mulțimi paralelipipedice la o σ algebră*
- b) *Am luat în considerație o infinitate de factori în produs .*

Observația 9.5 *In principiu nu e nevoie de o cunoaștere detaliată a produsului de câmpuri de probabilitate. Este suficient să știm că el există și că probabilitatea unei mulțimi paralelipipedice este produsul probabilităților factorilor (formula 9.2).*

Fie acum X o v.a. pe P , $X : P \rightarrow R$. În aceste condiții pe P^∞ avem un șir de v.a. definite prin:

$$X_i : P^\infty \rightarrow R, \quad X_i(\omega) = X_i((\omega_k)_{k \in N}) = X(\omega_i)$$

pentru orice $i \in N$. Prin urmare X_i aplicată unui șir este valoarea lui X pe componenta a i a șirului $\omega = (\omega_k)_{k \in N} \in P^\infty$. Aceste v.a. sunt independente și la fel distribuite (adică au aceeași funcție de repartiție, deci aceleași caracteristici numerice). Pe produsul finit $P^n = P \times P \dots \times P$ avem în mod analog variabilele aleatoare X_i definite prin formula de mai sus dar cu $\omega \in P^n$. Ele sunt independente și la fel distribuite.

In concluzie, mai multe măsurători ale unei mărimi apar în statistică astfel:

- a) *Urmărim o componentă numerică a unui fenomen, să zicem notată cu X .*
- b) *Acea caracteristică depinde de o seamă de factori dintr-o mulțime P , în general neexplicită.*
- c) *Admitem că pe P există o probabilitate p , iar $X:P \rightarrow R$ este o variabilă aleatoare.*
- d) *Prin n experiențe independente găsim pentru X valorile x_1, x_2, \dots, x_n .*
- e) *x_1, x_2, \dots, x_n apar ca valorile a n variabile aleatoare X_1, X_2, \dots, X_n definite pe spațiul produs P^n sau pe P^∞ . Aceste v.a. sunt independente și la fel distribuite ca X . Vom numi X_1, X_2, \dots, X_n variabile aleatoare de selecție asociate lui X . X_i reprezintă rezultatul experienței i . În cele ce urmează vom considera toate variabilele X_i definite pe aceeași mulțime P^∞ .*

Ne ocupăm în continuare de operațiile pe care le facem cu rezultatele x_1, x_2, \dots, x_n pentru a obține caracteristici ale variabilei aleatoare X . Vom folosi doar acele operații în care mulțimea P , care nu este explicită, nu intervine efectiv. Probabilitatea pe P^∞ o vom nota uneori cu Prob alteori cu p .

Definiția 9.6 *Se numește statistică un șir $(G_n)_{n \in N}$ de variabile aleatoare $G_n : P^\infty \rightarrow R$.*

Toate statisticile utilizate de noi vor fi de forma următoare:

i) *se dă un șir de funcții $g_n : R^n \rightarrow R$*

ii) *avem o v.a. $X : P \rightarrow R$*

Fie $X_1, X_2, \dots, X_n, \dots$ variabilele aleatoare asociate. Definim statistica:

$$g_n(X_1, X_2, \dots, X_n) : P^\infty \rightarrow R$$

, astfel

$$g_n(X_1, X_2, \dots, X_n)(\omega) = g_n(X_1(\omega), X_2(\omega), \dots, X_n(\omega)) \text{ pentru } \omega \in P^\infty$$

Uneori vom folosi termenul de statistică pentru șirul de variabile aleatoare construite mai sus pe P^∞ .

Exemple de statistici frecvent folosite

Fie X o v.a. cu funcția de repartiție $F : R \rightarrow R$, și X_1, X_2, \dots, X_n variabilele de selecție asociate. Vom folosi următoarele notații:

a1) $m = \int_{-\infty}^{\infty} x dF(x)$ = media lui X (Lecția 3)

b1) $M(X_1, X_2, \dots, X_n) = \bar{X}_{(n)} = \frac{X_1 + X_2 + \dots + X_n}{n} : P^\infty \rightarrow R$, o v.a. numită *media de selecție*. Uneori o vom nota cu $\bar{X}_{(n)}$ pentru a pune în evidență dependența de n , alteori o vom nota simplu \bar{X} . Astfel,

c1) $m^* = m^*(x_1, x_2, \dots, x_n) = \frac{x_1 + x_2 + \dots + x_n}{n}$ este valoarea mediei de selecție pentru rezultatele x_1, x_2, \dots, x_n obținute în cele n experiențe, numită și *media empirică* (Lecția 8).

ak) $m_k = \int_{-\infty}^{\infty} x^k dF(x)$ = momentul de ordin k al lui X (Lecția 3)

bk) $M_k(X_1, X_2, \dots, X_n) = \frac{X_1^k + X_2^k + \dots + X_n^k}{n}$, o v.a. numită *momentul de selecție de ordin k* .

ck) $m_k^* = m_k^*(x_1, x_2, \dots, x_n) = \frac{x_1^k + x_2^k + \dots + x_n^k}{n}$ = *momentul empiric de ordin k* (Lecția 8)

ak0) $\mu_k = \int_{-\infty}^{\infty} (x - m)^k dF(x)$ = momentul centrat de ordin k al lui X (Lecția 3)

bk0) $M_{0k}(X_1, X_2, \dots, X_n) = \frac{(X_1 - \bar{X})^k + (X_2 - \bar{X})^k + \dots + (X_n - \bar{X})^k}{n} : P^\infty \rightarrow R$, o v.a. numită *momentul centrat de ordin k* .

ck0) $\mu_k^* = \mu_k^*(x_1, x_2, \dots, x_n) = \frac{(x_1 - m^*)^k + (x_2 - m^*)^k + \dots + (x_n - m^*)^k}{n}$ = *momentul centrat de ordin k , empiric* (Lecția 8).

a20) $D = \sigma^2 = \int_{-\infty}^{\infty} (x - m)^2 dF(x)$ = dispersia lui X (Lecția 3).

b20) $D(X_1, \dots, X_n) = S^2(X_1, \dots, X_n) = \frac{(X_1 - \bar{X})^2 + (X_2 - \bar{X})^2 + \dots + (X_n - \bar{X})^2}{n} : P^\infty \rightarrow R$, o v.a. numită *dispersia de selecție*.

$$S'^2(X_1, X_2, \dots, X_n) = \frac{(X_1 - \bar{X})^2 + (X_2 - \bar{X})^2 + \dots + (X_n - \bar{X})^2}{n - 1}$$

se numește *dispersia de selecție modificată*.

c20) $s^{*2} = D^* = D^*(x_1, x_2, \dots, x_n) = \frac{(x_1 - m^*)^2 + (x_2 - m^*)^2 + \dots + (x_n - m^*)^2}{n} =$ *dispersia empirică de selecție* (Lecția 8), iar $s'^{*2} = \frac{(x_1 - m^*)^2 + (x_2 - m^*)^2 + \dots + (x_n - m^*)^2}{n - 1}$ se numește *dispersia empirică modificată*.

E de așteptat ca v.a. de selecție de mai sus să aproximeze într-un fel sau altul marimile corespunzătoare ale variabilei X .

Propoziția 9.7 Avem relația

$$S^2(X_1, X_2, \dots, X_n) = M_2(X_1, X_2, \dots, X_n) - \bar{X}^2 \quad (9.3)$$

Demonstratie

$$\begin{aligned}
S^2(X_1, X_2, \dots, X_n) &= \frac{1}{n} \sum (X_i - \bar{X})^2 = \frac{1}{n} \sum (X_i^2 - 2\bar{X}X_i + \bar{X}^2) \\
&= \frac{1}{n} \sum X_i^2 - \frac{2}{n} \cdot \bar{X} \cdot \sum X_i + \frac{1}{n} \cdot n\bar{X}^2 \\
&= \frac{1}{n} \sum X_i^2 - 2 \cdot \bar{X} \cdot \bar{X} + \bar{X}^2 \\
&= M_2(X_1, X_2, \dots, X_n) - \bar{X}^2
\end{aligned}$$

QED.

Fie X o v.a. și $X_1, X_2, \dots, X_n \dots$ variabilele de selecție asociate. Fie de asemenea $A \in \mathbb{R}$.

Definiția 9.8 Se numește estimator sau funcție de estimatie pentru A , o statistică $(g_n)_{n \in \mathbb{N}}$ astfel ca pentru orice $\epsilon > 0$ să avem:

$$\lim_{n \rightarrow \infty} \text{Prob}(|g_n(X_1, X_2, \dots, X_n) - A| > \epsilon) = 0$$

Cu alte cuvinte, $\epsilon > 0$ fiind dat, pentru valori mari ale lui n este foarte puțin probabil ca variabila aleatoare $g_n(X_1, X_2, \dots, X_n)$ să ia valori în afara intervalului $[A - \epsilon, A + \epsilon]$, adică este foarte puțin probabil ca numărul

$g_n(x_1, x_2, \dots, x_n)$ să fie în afara intervalului $[A - \epsilon, A + \epsilon]$. În aceste condiții, după un număr de n experiențe, considerăm pe $g_n(x_1, x_2, \dots, x_n)$ ca o aproximație bună pentru A . Este posibil să ne înșelăm, dar probabilitatea de a ne înșela este mică, pentru n mare. Statistica nu ne oferă răspunsuri sigure ci doar aproximații în care putem avea un grad mai mic sau mai mare de încredere. Se acceptă acele aproximații în care avem un grad mai mare de încredere.

Definiția 9.9 O statistică $(g_n(X_1, \dots, X_n))_{n \in \mathbb{N}}$ se numește corectă sau deplasată relativ la valoarea A dacă avem:

$$1) \lim_{n \rightarrow \infty} M(g_n(X_1, X_2, \dots, X_n)) = A.$$

$$2) \lim_{n \rightarrow \infty} D(g_n(X_1, X_2, \dots, X_n)) = 0.$$

și se numește *absolut corectă sau nedepășată* dacă în plus $M(g_n(X_1, X_2, \dots, X_n)) = A$.

Condițiile 1) și 2) din definiția de mai sus pun în evidență situații în care o statistică oarecare $(g_n)_{n \in \mathbb{N}}$ este un estimator pentru o valoare A . Teorema de mai jos pune în evidență importanța condițiilor din definiția anterioară.

Teorema 9.10 Dacă statistica $(g_n(X_1, X_2, \dots, X_n))_{n \in \mathbb{N}}$ este corectă relativ la A atunci ea este un estimator al lui A , adică pentru orice $\epsilon > 0$ avem

$$\lim_{n \rightarrow \infty} \text{Prob}(|g_n(X_1, X_2, \dots, X_n) - A| > \epsilon) = 0.$$

Demonstrație. Conform cu inegalitatea lui Cebâșev (Lecția 5) pentru un $\epsilon > 0$ avem:

$$\begin{aligned} \text{Prob}(|g_n(X_1, X_2, \dots, X_n) - M(g_n(X_1, X_2, \dots, X_n))| > \epsilon) &\leq \\ &\leq \frac{D(g_n(X_1, X_2, \dots, X_n))}{\epsilon^2} \end{aligned}$$

Acum ținând seama de 1) și 2) din definiția corectitudinii rezultă

$$\text{Prob}(|g_n(X_1, X_2, \dots, X_n) - A| > \epsilon) \rightarrow 0$$

când $n \rightarrow \infty$, deci statistica $(g_n(X_1, X_2, \dots, X_n))_{n \in N}$ este un estimator al lui A.

QED.

Arătăm acum că funcțiile de selecție introduse cu ocazia notațiilor precedente sunt estimatori pentru valorile corespunzătoare ale variabilei X.

Teorema 9.11 a) *Statistica media de selecție:*

$$g_n(X_1, X_2, \dots, X_n) = \bar{X}_{(n)} = (X_1 + X_2 + \dots + X_n) / n$$

estimează media $m = M(X)$ a v.a. X absolut corect.

b) *Statistica*

$$h_n(X_1, X_2, \dots, X_n) = \frac{1}{n} \sum X_i^r$$

estimează absolut corect momentul de ordin r, m_r , al v.a. X.

c) *Statistica*

$$S^2(X_1, X_2, \dots, X_n) = \frac{1}{n} \sum (X_i - \bar{X}_{(n)})^2$$

estimează corect dar nu absolut corect dispersia v.a. X, $\sigma^2 = D(X)$.

d) $S'^2(X_1, X_2, \dots, X_n) = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_{(n)})^2$, adică dispersia de selecție modificată, aproximează absolut corect dispersia v.a. X.

Demonstrație Trebuie verificate condițiile 1 și 2 din definiția statisticii corecte.

a) Verificăm 1): $M(g_n) = \frac{1}{n} \sum M(X_i) = \frac{1}{n} \cdot n \cdot m = m = M(X)$.

Verificăm 2): $D(g_n) = \frac{1}{n^2} \sum D(X_i) = \frac{D(X)}{n}$, deoarece X_1, X_2, \dots, X_n sunt independente. Prin urmare $D(g_n) \rightarrow 0$, când $n \rightarrow \infty$.

b) Verificăm 1): $M(h_n) = \frac{1}{n} \sum M(X_i^r) = \frac{1}{n} \cdot n \cdot m_r = m_r$

Verificăm 2): $D(h_n) = \frac{1}{n^2} \sum D(X_i^r) = \frac{1}{n} D(X^r) \rightarrow 0$, când $n \rightarrow \infty$. Retinem formula:

$$D(\bar{X}_{(n)}) = D(X)/n \quad (9.4)$$

c) încercăm să verificăm 1):

$$\begin{aligned}
 & M(S^2(X_1, X_2, \dots, X_n)) \stackrel{\text{din 9.3}}{=} M(M_2(X_1, X_2, \dots, X_n) - \bar{X}_{(n)}^2) \\
 &= M(M_2(X_1, X_2, \dots, X_n)) - \frac{1}{n^2} M\left(\left(\sum X_i\right)^2\right) \\
 &= M(X^2) - \frac{1}{n^2} M\left(\left(\sum X_i\right)\left(\sum X_j\right)\right) \\
 &= M(X^2) - \frac{1}{n^2} M\left(\sum X_i^2\right) - \frac{1}{n^2} M\left(\sum_{i \neq j} X_i X_j\right) \\
 &= M(X^2) - \frac{n}{n^2} M(X^2) - \frac{1}{n^2} \sum_{i \neq j} M(X_i X_j) \\
 &\stackrel{X_i, X_j \text{ independente.}}{=} \frac{n-1}{n} M(X^2) - \frac{1}{n^2} \cdot n(n-1) M(X)^2 \\
 &= \frac{n-1}{n} (M(X^2) - M(X)^2) = \frac{n-1}{n} D(X) \rightarrow D(X)
 \end{aligned}$$

Prin urmare S^2 nu estimează absolut corect dispersia v.a. X . Retinem formula deja găsită:

$$M(S^2(\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n)) = \frac{n-1}{n} D(\mathbf{X}) \quad (9.5)$$

E clar că $M(S'^2(X_1, X_2, \dots, X_n)) = D(X)$.

Verificăm 2): lăsăm ca exercitiu pentru cititor verificarea formulei

$$D(S'^2) = \frac{1}{n} \left(m_4 - \frac{n-3}{n-1} [D(X)]^2 \right) \quad (9.6)$$

unde m_4 este momentul de ordin 4 al v.a. X .

Se vede clar de aici că $D(S'^2) \rightarrow 0$, când $n \rightarrow \infty$. De asemenea $D(S^2) = \left(\frac{n-1}{n}\right)^2 D(S'^2) \rightarrow 0$ când $n \rightarrow \infty$.

QED.

Observația 9.12 *In practică se folosește S'^2 în locul lui S^2 deoarece dă rezultate mai bune după cum ne arată teorema 2. Totuși formula (9.6) ne spune că pentru n suficient de mare și statistica S^2 poate fi folosită ca estimator al dispersiei v.a. X . Din definiția lui S'^2 și din formula (9.6) găsim formula utilă:*

$$D(S^2) = \frac{(n-1)^2}{n^3} \left(m_4 - \frac{n-3}{n-1} [D(X)]^2 \right) \quad (9.7)$$

Exercițiul 9.13 Fie selectia $\{0, 1, 1, 0, 1, 1, 2, 0, 0, 2\}$. Să se estimeze absolut corect dispersia populației din care provine această selecție.

Soluție. Media este estimată absolut corect de media empirică $m^* = 8/10 = 0,8$.
Dispersia este estimată absolut corect de dispersia modificată empirică

$$\begin{aligned} s'^{*2} &= \frac{1}{9}[(0 - 0,8)^2 + (1 - 0,8)^2 + (1 - 0,8)^2 + (0 - 0,8)^2 + (1 - 0,8)^2 \\ &\quad + (1 - 0,8)^2 + (2 - 0,8)^2 + (0 - 0,8)^2 + (0 - 0,8)^2 + (2 - 0,8)^2] \\ &= 0,56 \end{aligned}$$

Observația 9.14 Deoarece dispersia se mai numește și varianță vom folosi și noi uneori varianța de selecție pentru dispersia de selecție.

9.1 Principiul verosimilității maxime

Presupunem că P este o populație unde se urmărește caracteristica numerică X , care este o variabilă aleatoare cu densitatea de probabilitate $f(x; \theta)$, θ fiind un parametru necunoscut. Cunoaștem doar forma matematică a funcției $f(x; \theta)$. De exemplu dacă știm că X este o v.a. normală cu media θ , necunoscută dar cu dispersia σ^2 cunoscută, atunci $f(x; m) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-m)^2}{2\sigma^2}}$.

Pentru determinarea lui θ facem o selecție care dă rezultatele $\{x_1, \dots, x_n\}$ și încercăm pe baza lor să estimăm pe θ . Deoarece v.a. de selecție X_1, \dots, X_n sunt independente, probabilitatea ca X_1 să ia valori în intervalul $[x_1, x_1 + dx_1)$, X_2 să ia valori în $[x_2, x_2 + dx_2)$, ..., X_n să ia valori în $[x_n, x_n + dx_n)$ este dată de $f(x_1; \theta) \cdot f(x_2; \theta) \cdots f(x_n; \theta) dx_1 dx_2 \cdots dx_n = L(x_1, \dots, x_n; \theta) dx_1 dx_2 \cdots dx_n$. Această funcție L se numește *funcția de verosimilitate* și va fi folosită pentru estimarea lui θ .

Dacă X ia valori discrete, atunci $f(x, \theta)$ este probabilitatea ca X să ia valoarea x . De exemplu, în cazul distribuției Poisson, $f(x; \theta) = e^{-\theta} \cdot \frac{\theta^x}{x!}$, cu $x \in \mathbf{N}$, reprezintă probabilitatea ca $X = x$, iar θ este parametrul necunoscut (pe care urmează să-l estimăm!). Probabilitatea ca în n selecții independente să se obțină rezultatele x_1, x_2, \dots, x_n este $f(x_1, \theta) \cdot f(x_2, \theta) \cdots f(x_n, \theta) = L(x_1, x_2, \dots, x_n; \theta)$ care se numește și în acest caz funcția de verosimilitate.

Funcția L este determinată de volumul selecției n și depinde de θ . Metoda verosimilității maxime constă în următorul principiu (axiomă): *valoarea "cea mai verosimilă" (cea mai potrivită în acest sens!) a parametrului θ este aceea pentru care funcția $L(x_1, \dots, x_n; \theta)$ este maximă.* După cum știm de la Analiza matematică, această cerință are loc dacă avem:

$$\frac{\partial L(x_1, \dots, x_n; \theta)}{\partial \theta} = 0 \quad (9.8)$$

adică θ este un punct critic pentru $L(x_1, \dots, x_n; \theta)$.

Ecuatia (9.8) în practică se dovedeste dificilă. De aceea cel mai des se folosește observația: $L(x_1, \dots, x_n; \theta)$ este maximă dacă și numai dacă $\ln L(x_1, \dots, x_n; \theta)$ este maximă (funcția logaritmică este strict crescătoare). Deci (9.8) este echivalentă cu :

$$\frac{\partial \ln L(x_1, \dots, x_n; \theta)}{\partial \theta} = 0 \quad (9.9)$$

care poartă numele de *ecuație a verosimilității maxime*. Rezolvăm ecuația (9.9), sau ecuația (9.8) și găsim $\theta = \Theta_n(x_1, \dots, x_n)$. Ca estimator (funcție de estimare) pentru θ luăm variabila aleatoare $\Theta_n(X_1, X_2, \dots, X_n)$, care, pentru selecția $\{x_1, x_2, \dots, x_n\}$ dă rezultatul $\Theta_n(x_1, x_2, \dots, x_n)$.

Se poate demonstra că în condiții foarte generale, pentru selecții mari, statistica $\Theta(X_1, X_2, \dots, X_n)$ obținută prin metoda verosimilității maxime, are o distribuție aproximativ normală, cu media egală cu θ =valoarea adevărată a parametrului și dispersia

$$D(\Theta) = \frac{1}{-n \int_{-\infty}^{\infty} \left(\frac{\partial^2 \ln(f(x, \theta))}{\partial \theta^2} \right) f(x; \theta) dx} = \frac{1}{n \int_{-\infty}^{\infty} \left(\frac{\partial \ln f(x; \theta)}{\partial \theta} \right)^2 f(x; \theta) dx}$$

Dacă distribuția este discretă atunci integralele din formula precedentă devin sume.

Exemplul 9.15 *Presupunem că populația are distribuția Poisson (cazul evenimentelor rare). Funcția de probabilitate este $f(k; \lambda) = e^{-\lambda} \cdot \frac{\lambda^k}{k!}$, $k = 0, 1, 2, \dots$. Ne interesează să estimăm parametrul λ prin metoda verosimilității maxime. Pentru aceasta facem o selecție $\{x_1, x_2, \dots, x_n\} \subset \{0, 1, 2, \dots\}$.*

$$L(x_1, \dots, x_n; \lambda) = f(x_1; \lambda) \cdot f(x_2; \lambda) \cdot \dots \cdot f(x_n; \lambda) = e^{-n\lambda} \cdot \frac{\lambda^{\sum x_k}}{x_1! x_2! \dots x_n!}$$

$$\ln L(x_1, \dots, x_n; \lambda) = -n\lambda + \left(\sum x_k \right) \ln \lambda - \sum \ln(x_k!)$$

$\frac{\partial \ln L}{\partial \lambda} = 0$ ne furnizează $\lambda = \frac{\sum x_k}{n}$, deci un estimator pentru λ este $\Lambda_n(X_1, \dots, X_n) = \frac{X_1 + X_2 + \dots + X_n}{n}$ adică media de selecție. Deoarece λ este media lui X (variabilă Poisson), Λ_n este un estimator absolut corect pentru λ .

Este el oare și cel mai eficient, în sensul că are dispersia cea mai mică? Este greu de răspuns la această întrebare. Totuși avem un rezultat puternic care face oarecare lumină:

Teorema 9.16 (Rao-Cramer) *Dacă statistica $g_n(X_1, \dots, X_n)$ dă un estimator eficient (cu dispersia minimă, în multimea tuturor estimatorilor absolut corecți pentru θ), atunci*

$$D(g_n(X_1, \dots, X_n)) = \frac{1}{n \int_{-\infty}^{\infty} \left(\frac{\partial \ln f(x; \theta)}{\partial \theta} \right)^2 f(x; \theta) dx}$$

sau

$$D(g_n(X_1, \dots, X_n)) = \frac{1}{n \sum_{x=0}^{\infty} \left(\frac{\partial \ln f(x; \theta)}{\partial \theta} \right)^2 f(x; \theta)} \quad (9.10)$$

dacă distribuția este cu valori discrete.

Fără demonstrație.

Ne întoarcem la exemplul anterior. Stim că $D(\Lambda_n) = D\left(\frac{\sum X_k}{n}\right) = \frac{D(X)}{n} = \frac{\lambda}{n}$. $f(x, \lambda) = \frac{e^{-\lambda} \lambda^x}{x!}$, deci $\frac{\partial \ln f(x, \lambda)}{\partial \lambda} = -1 + \frac{x}{\lambda}$. De aici rezultă

$$\begin{aligned} & \sum_{x=0}^{\infty} \left(\frac{\partial \ln f(x, \lambda)}{\partial \lambda} \right)^2 f(x, \lambda) \\ &= \sum_{x=0}^{\infty} \left(1 - 2\frac{x}{\lambda} + \frac{x^2}{\lambda^2} \right) e^{-\lambda} \frac{\lambda^x}{x!} \\ &= e^{-\lambda} \left(\underbrace{\sum_{x=0}^{\infty} \frac{\lambda^x}{x!}}_{=e^\lambda} - 2 \underbrace{\sum_{x=1}^{\infty} \frac{\lambda^{x-1}}{(x-1)!}}_{=e^\lambda} + \sum_{x=1}^{\infty} \underbrace{\frac{\lambda^{x-2}}{(x-1)!} (x-1+1)}_{=e^\lambda + \frac{1}{\lambda} e^\lambda} \right) \\ &= \frac{1}{\lambda} \end{aligned}$$

Prin urmare $\frac{1}{n \sum_{x=0}^{\infty} \left(\frac{\partial \ln f(x; \lambda)}{\partial \lambda} \right)^2 f(x; \lambda)} = \frac{\lambda}{n} = D(\Lambda_n)$. Rezultă din teorema Rao-Cramer că statistica medie de selecție este și un estimator eficient pentru λ . Putem spune acum că $\lambda = (\sum x_k) / n$ este o estimatie "foarte bună" în toate sensurile.

Exemplul 9.17 Să se estimeze parametrul p al unei distribuții Bernoulli

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

prin metoda verosimilității maxime.

Soluție.

$$f(x, p) = \begin{cases} p & \text{dacă } x=1 \\ 1-p & \text{dacă } x=0 \end{cases}$$

Funcția de verosimilitate este $L(x_1, x_2, \dots, x_n) = p^{n_1} (1-p)^{n-n_1}$ unde n_1 este numărul de realizări ale lui 1. $\frac{\partial \ln L}{\partial p} = 0$ devine $\frac{\partial (n_1 \ln p + (n-n_1) \ln(1-p))}{\partial p} = 0$ adică $\frac{n_1}{p} - \frac{n-n_1}{1-p} = 0$ care are soluția $p = \frac{n_1}{n}$. Valoarea n_1 este valoarea variabilei $X_1 + X_2 + \dots + X_n$, unde X_i este variabila de

selecție a cărei valoare este 1 dacă la experiența i se obține rezultatul 1 și are valoarea 0 în caz contrar. Prin urmare statistica ce estimează parametrul p este

$$\mathcal{P} = \frac{X_1 + X_2 + \dots + X_n}{n}$$

care este chiar media de selecție. La fel ca în cazul repartiției Poisson se arată că

$$\begin{aligned} D(\mathcal{P}) &= \frac{pq}{n} \\ &= \frac{1}{n \sum_{x=0}^1 \left(\frac{\partial \ln f(x;p)}{\partial p} \right)^2 f(x;p)} = \frac{1}{n \left(\left(\frac{\partial \ln(1-p)}{\partial p} \right)^2 (1-p) + \left(\frac{\partial \ln p}{\partial p} \right)^2 p \right)} \end{aligned}$$

deci estimarea lui p este absolut corectă (exercițiu).

Dacă avem de estimat mai mulți parametri $\theta_1, \theta_2, \dots, \theta_p$, știind că densitatea de probabilitate a variabilei aleatoare X este $f(x; \theta_1, \dots, \theta_p)$, atunci în mod analog cu cazul unui singur parametru, principiul verosimilității maxime spune că în urma a n experiențe independente care dau rezultatele x_1, x_2, \dots, x_n , se aleg pentru parametri acele valori care maximizează funcția de verosimilitate $L(x_1, \dots, x_n; \theta_1, \dots, \theta_p) = f(x_1; \theta_1, \theta_2, \dots, \theta_p) \cdot f(x_2; \theta_1, \dots, \theta_p) \cdot \dots \cdot f(x_n; \theta_1, \theta_2, \dots, \theta_p)$ sau ceea ce este același lucru acele valori care maximizează $\ln L(x_1, x_2, \dots, x_n; \theta_1, \dots, \theta_p)$. Aceasta implică:

$$\begin{cases} \frac{\partial \ln L(x_1, x_2, \dots, x_n; \theta_1, \dots, \theta_p)}{\partial \theta_1} = 0 \\ \frac{\partial \ln L(x_1, x_2, \dots, x_n; \theta_1, \dots, \theta_p)}{\partial \theta_2} = 0 \\ \dots \\ \frac{\partial \ln L(x_1, x_2, \dots, x_n; \theta_1, \dots, \theta_p)}{\partial \theta_p} = 0 \end{cases} \quad (9.11)$$

Exemplul 9.18 Să presupunem că v.a. X are o distribuție normală cu

$f(x; m; \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-m)^2}{2\sigma^2}}$. În acest caz avem $\ln L = -\frac{n}{2} \ln(2\pi) - n \ln \sigma - \frac{\sum (x_i - m)^2}{2\sigma^2}$. Sistemul (9.11) devine:

$$\begin{cases} \frac{1}{\sigma^2} \sum (x_i - m) = 0 \\ -\frac{n}{\sigma} + \frac{\sum (x_i - m)^2}{\sigma^3} = 0 \end{cases}$$

care are ca soluții $m = \frac{\sum x_i}{n} = m^*(x_1, \dots, x_n)$ și $\sigma = \sqrt{\frac{\sum (x_i - m^*)^2}{n}}$.

Știm din această lecție că $M(X_1, X_2, \dots, X_n) = \bar{X}_{(n)} = \frac{X_1 + X_2 + \dots + X_n}{n}$ și $D(X_1, \dots, X_n) = \frac{(X_1 - \bar{X})^2 + (X_2 - \bar{X})^2 + \dots + (X_n - \bar{X})^2}{n}$ sunt estimări corecte pentru media și dispersia unei variabile aleatoare, în cazul nostru pentru m și σ^2 .

9.2 Metoda momentelor (K. Pearson)

Dată selecția $\{x_1, \dots, x_n\}$ noi putem calcula momentul de ordin k al selecției: $m_k^* = \frac{\sum_i x_i^k}{n}$, pentru orice $k = 0, 1, 2, \dots$. Obținem astfel estimatori pentru medie, dispersie, momente de diferite ordine. Funcția caracteristică $X_c(t)$ are toate derivatele în $t = 0$ date de $X_c^{(k)}(0) = i^k M^k(X)$. Prin urmare în condiții foarte generale, care asigură că $X_c(t)$ este analitică (se poate dezvolta în serie convergentă de puteri în jurul oricărui punct), rezultă că momentele $M^k(X)$ determină pe $X_c(t)$ care la rândul ei determină repartiția lui X (vezi Lecția 3). Această observație a fost folosită de K. Pearson pentru a găsi estimatori pentru parametrii unei legi de probabilitate.

Fie $\rho(x; \theta_1, \theta_2, \dots, \theta_p)$ densitatea de probabilitate a v.a. X , unde parametrii $\theta_1, \dots, \theta_p$ sunt necunoscuți. Există relațiile: $m_k = \int_{-\infty}^{\infty} \rho(x; \theta_1, \theta_2, \dots, \theta_p) x^k dx$ pentru orice k . Am văzut că m_k este estimat de m_k^* . Egalând valoarea teoretică exactă cu estimarea practică, $m_k = m_k^*$, adică:

$$\left\{ \int_{-\infty}^{\infty} x^k \cdot \rho(x; \theta_1, \theta_2, \dots, \theta_p) dx = \frac{\sum_{i=1}^n x_i^k}{n} \right. \quad (9.12)$$

pentru $k=1, 2, \dots, p$, obținem un sistem care dă prin rezolvare $\theta_k = \Theta_n(x_1, x_2, \dots, x_n)$, pentru $k=1, 2, \dots, p$.

Ca estimatori pentru θ_k se iau v.a. $\Theta_k(X_1, X_2, \dots, X_n)$.

Dacă v.a. X este discretă atunci integrala din (9.12) devine sumă, la fel ca în cazul metodei verosimilității maxime.

Exemplul 9.19 Fie v.a. cu densitatea $\rho(x; \lambda) = \frac{1}{\Gamma(\lambda)} x^{\lambda-1} e^{-x}$, pentru $\lambda > 0$ $x > 0$. Aici parametrul este λ . Se cere o metodă de a estima pe λ prin selecții. Dacă folosim metoda momentelor găsim

$$\int_0^{\infty} x \cdot \frac{1}{\Gamma(\lambda)} x^{\lambda-1} e^{-x} dx = \frac{x_1 + x_2 + \dots + x_n}{n}$$

adică $\lambda = \frac{x_1 + x_2 + \dots + x_n}{n}$.

Exemplul 9.20 Densitatea de probabilitate a unei v.a. X are forma:

$$f(x; a; b; c) = \begin{cases} a + 2bx, & x \in [0, c] \\ 0, & \text{în rest} \end{cases}$$

Rezultatele unei selecții de volum $n=3$ dau pentru X valorile $\{x_1, x_2, x_3\} = \{-1, 0, 1\}$. Să se estimeze parametrii a, b, c prin metoda momentelor.

Mai întâi punem condiția ca $\int_{-\infty}^{\infty} f(x; a; b; c) dx = 1$, de unde găsim relația:

$$ac + bc^2 = 1 \quad (9.13)$$

Calculăm acum media v.a. X :

$$M(X) = \int_0^c x(a + 2bx)dx = \frac{ac^2}{2} + \frac{2bc^3}{3} \quad (9.14)$$

momentul de ordin 2:

$$M_2(X) = \int_0^c x^2(a + 2bx)dx = \frac{ac^3}{3} + \frac{2bc^4}{4} \quad (9.15)$$

Momentele de selectie $m_1^* = (-1+0+1)/3=0$ si $m_2^* = ((-1)^2+0^2+1^2)/3=2/3$ vor estima pe $M(X)$ si pe $M_2(X)$. Deci vom obtine sistemul nelinier de ecuatii:

$$\begin{cases} ac + bc^2 = 1 \\ \frac{ac^2}{2} + \frac{2bc^3}{3} = 0 \\ \frac{ac^3}{3} + \frac{2bc^4}{4} = \frac{2}{3} \end{cases} \quad (9.16)$$

In general rezolvarea acestor sisteme (care se obtin folosind metoda momentelor) este foarte complicată. Din ecuatia a doua găsim $c = -\frac{3a}{4b}$. înlocuim expresia lui c în prima si în ultima ecuatie si găsim:

$$\begin{cases} -\frac{3a^2}{4b} + \frac{9a^2}{16b} = 1 \\ -\frac{9a^4}{64b^3} + \frac{81a^4}{512b^3} = \frac{2}{3} \end{cases} \quad (9.17)$$

Din prima ecuatie a sistemului (9.17) găsim $b = -\frac{3a^2}{16}$. Inlocuim expresia lui b în ecuatia a doua si găsim că $2a^2 = -8$, lucru imposibil. Concluzia este că selectia $\{-1,0,1\}$ nu poate fi pentru variabila X . Sondajul respectiv este eronat sau X nu are densitatea de probabilitate propusă. De altfel, examinând cu atenție rezultatele sondajului vedem că valoarea -1 în principiu nu ar fi trebuit să se obțină deoarece X are densitatea pozitivă pe $[0,c]$.

9.3 Exerciții

1. Pentru o selectie de volum n dintr-o distributie exponentială ($\rho(t) = \lambda e^{-\lambda t}$, dacă $t \geq 0$ si 0 în rest) cu parametrul λ , să se găsească un estimator pentru λ folosind metoda verosimilității maxime. Presupunem că $\rho(t)$ este densitatea de probabilitate a duratei dintre doua sosiri succesive la o staie de benzină. Se cronometrează 11 sosiri și se găsesc următorii timpi între ele 4, 3, 6, 1, 1, 4, 2, 6, 1, 3. Calculati prin metoda verosimilității maxime pe λ .

2. Considerăm o selectie de volum n dintr-o populatie cu distributia gama ($f(t) = \lambda(\lambda t)^{r-1}e^{-\lambda t} \cdot \frac{1}{(r-1)!}$, dacă $t \geq 0$ si 0 în rest). Găsiți un estimator pentru λ prin metoda verosimilității maxime si un altul prin metoda momentelor.

3. Viata unui bec electric, măsurată în numărul de ore de functionare continuă până când se arde, se presupune uniform distribuită cu parametrii a și b :

$$f(x) = \begin{cases} 1/(b-a), & a \leq x \leq b \\ 0, & \text{în rest.} \end{cases}$$

Se face o selecție de n becuri și se notează cu x_1, \dots, x_n timpurile de functionare ai acestora până când se ard. Determinați estimatori pentru a și b prin metoda momentelor.

4. Funcția de probabilitate a v.a. X este dată de

$$f(x) = \begin{cases} \frac{2b(c-bx)}{c^2}, & \text{dacă } 0 \leq x \leq \frac{c}{b} \\ 0, & \text{în rest} \end{cases}$$

. Stim că media $M(X) = \frac{c}{3b}$ și $\sigma_X^2 = \frac{c^2}{18b^2}$.

i) dacă $c=3$, este oare media de selecție M a unui esanșon de volum n un estimator nedeplasat pentru parametrul b ?

ii) dacă $b=1/3$, este M un estimator pentru c ? ($P(|M-c| < \epsilon) \rightarrow 1$, când $n \rightarrow \infty$). Indicație: folosiți inegalitatea lui Cebîșev sau teoria din această lecție.

5. Fie statistica $g(X_1, X_2, \dots, X_n) = a_1X_1 + a_2X_2 + \dots + a_nX_n$, cu $a_1, \dots, a_n \in \mathbf{R}$. Cum trebuie să fie numerele a_1, \dots, a_n , astfel încât g să fie un estimator nedeplasat pentru media m a populației? Indicație: cereți $M(g)=m$.

6. Dacă $g(X_1, X_2, \dots, X_n)$ este un estimator nedeplasat pentru parametrul θ , este adevărat că și g^2 este un estimator nedeplasat pentru θ^2 ?

7. Greutatea unor utilaje produse de o firmă este distribuită normal cu dispersia cunoscută σ_X^2 , dar cu media m necunoscută. Fie statisticile

$$G = \frac{(X_1 + X_2 + X_3 + X_4 + \dots + X_n)}{n} \text{ și } H = (X_1 + 2X_2 + 3X_3 + 4X_4 + \dots + nX_n) \cdot \frac{2}{n(n+1)}, \quad n \in N.$$

a) Să se arate că G și H sunt estimatori nedeplasați pentru m .

b) Care estimator are dispersia mai mică?

Lecția 10

Intervale de încredere

Definiția 10.1 Fie P o populație, θ un parametru al ei și $g = g(X_1, \dots, X_n)$, $h = h(X_1, \dots, X_n)$ două statistici astfel încât $g(X_1, \dots, X_n) \leq h(X_1, \dots, X_n)$, adică oricare ar fi selecția $\{x_1, \dots, x_n\}$ să avem că $g(x_1, \dots, x_n) \leq h(x_1, \dots, x_n)$. Spunem că intervalul $[g, h]$ este un interval de încredere pentru parametrul θ , de nivel de încredere α dacă avem relația:

$$\text{Prob}\{g \leq \theta \leq h\} \geq \alpha \quad (10.1)$$

Numărul $\varepsilon = 1 - \alpha$ se mai numește *prag de încredere*. De obicei α se exprimă în procente, de exemplu pentru $\alpha = 0,95$ putem scrie $\alpha = 95\%$.

Cerința (10.1) trebuie înțeleasă astfel: dacă după un număr mare de selecții $\{x_1, \dots, x_n\}$, să zicem N , K dintre ele dau intervale $[g(x_1, x_2, \dots, x_n), h(x_1, x_2, \dots, x_n)]$ cu proprietatea că $\theta \in [g, h]$ (pentru fiecare selecție fixată, intervalul devine interval obișnuit, numeric), atunci $K/N \geq \alpha$. Altfel spus, intervalele $[g, h]$ acoperă pe θ în proporție de cel puțin α % (de exemplu, dacă $\alpha = 1/5 = 20/100$, $\alpha = 20\%$).

Definiția 10.2 Pentru un α , un interval de încredere $[g_\alpha, h_\alpha]$ de lungime minimă, astfel încât $\text{Prob}\{g_\alpha \leq \theta \leq h_\alpha\} = \alpha$, se zice interval de încredere eficient, relativ la încrederea α .

Pentru calculele următoare vom avea nevoie de teorema:

Teorema 10.3 Fie X o v.a. normală, de tip $N(m, \sigma)$. Fie X_1, X_2, \dots, X_n variabilele de selecție asociate cu X . Atunci avem:

- a) Variabila $\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$ este de tipul $N\left(m, \frac{\sigma}{\sqrt{n}}\right)$.
- b) Variabila $\left(\frac{X_1 - m}{\sigma}\right)^2 + \left(\frac{X_2 - m}{\sigma}\right)^2 + \dots + \left(\frac{X_n - m}{\sigma}\right)^2$ este de tip $H(n)$ adică este o variabilă χ^2 standard cu n grade de libertate.
- c) Variabila $\left(\frac{X_1 - \bar{X}}{\sigma}\right)^2 + \left(\frac{X_2 - \bar{X}}{\sigma}\right)^2 + \dots + \left(\frac{X_n - \bar{X}}{\sigma}\right)^2$ este de tip $H(n-1)$ adică este de tip χ^2 cu $n-1$ grade de libertate și este independentă față de variabila \bar{X} .

d) Variabila $\sqrt{(n-1)n} \frac{\bar{X}-m}{\sqrt{(X_1-\bar{X})^2+(X_2-\bar{X})^2+\dots+(X_n-\bar{X})^2}}$ este de tip Student cu $n-1$ grade de libertate.

Demonstrație. a) Această afirmație este demonstrată în lecția 4, secțiunea "Repartiția normală".

b) Deoarece $\left(\frac{X_i-m}{\sigma}\right)$ sunt normale de tip $N(0,1)$ și independente, afirmația de la acest punct rezultă din lecția 4, secțiunea "Distribuția χ^2 ".

c) Faptul că \bar{X} și $\left(\frac{X_1-\bar{X}}{\sigma}\right)^2 + \left(\frac{X_2-\bar{X}}{\sigma}\right)^2 + \dots + \left(\frac{X_n-\bar{X}}{\sigma}\right)^2$ sunt independente nu se demonstrează în acest curs. Acum scriem că $X_i-m = (X_i-\bar{X}) + (\bar{X}-m)$, de unde $\sum (X_i-m)^2 = \sum (X_i-\bar{X})^2 + \sum (\bar{X}-m)^2 + 2(\bar{X}-m) \sum (X_i-\bar{X})$. Dar $\sum (X_i-\bar{X}) = 0$, deoarece $\bar{X} = \frac{1}{n} \sum X_i$. Prin urmare

$$\sum \left(\frac{X_i-m}{\sigma}\right)^2 = \sum \left(\frac{X_i-\bar{X}}{\sigma}\right)^2 + \frac{n(\bar{X}-m)^2}{\sigma^2} = \sum \left(\frac{X_i-\bar{X}}{\sigma}\right)^2 + \left(\frac{\bar{X}-m}{\sigma/\sqrt{n}}\right)^2$$

Membrul stâng este de tip $H(n)$, iar în membrul doi avem o sumă de v.a. independente, dintre care a doua este de tip $H(1)$ fiind pătratul unei v.a. normale, de tip $N(0,1)$ (vezi lecția 4). Prin urmare am găsit $\chi_{(n)}^2 = ? + \chi_{(1)}^2$. Comparând această relație cu $\chi_{(p+q)}^2 = \chi_{(p)}^2 + \chi_{(q)}^2$, unde indicii de jos indică numărul de grade de libertate (vezi lecția 4), găsim că $\sum \left(\frac{X_i-\bar{X}}{\sigma}\right)^2$ este de tip $\chi_{(n-1)}^2$.

d) Conform cu lecția 6, secțiunea "Distribuția Student", variabila aleatoare

$$\frac{\frac{\bar{X}-m}{\sigma/\sqrt{n}}}{\sqrt{\frac{\sum \left(\frac{X_i-\bar{X}}{\sigma}\right)^2}{n-1}}} = \sqrt{(n-1)n} \frac{\bar{X}-m}{\sqrt{(X_1-\bar{X})^2+(X_2-\bar{X})^2+\dots+(X_n-\bar{X})^2}}$$

fiind de tipul $\frac{f}{\sqrt{\frac{g}{n-1}}}$, cu f de tipul $N(0,1)$ și g de tipul $H(n-1) = H(n-1,1)$, rezultă că are o distribuție Student cu $n-1$ grade de libertate.

QED.

10.1 Intervale de încredere pentru medie

Să considerăm o caracteristică numerică X care are o distribuție normală de medie m și dispersie σ^2 . Dacă în urma unei selecții de volum n s-au obținut rezultatele x_1, x_2, \dots, x_n pentru X , atunci, conform celor arătate în lecția trecută valoarea $\frac{x_1+x_2+\dots+x_n}{n}$ este o estimare bună pentru m iar $\frac{(x_1-\bar{x})^2+(x_2-\bar{x})^2+\dots+(x_n-\bar{x})^2}{n}$ este o estimare bună pentru σ^2 . Ce încredere putem avea în aceste estimări? În continuare vom da un răspuns la această întrebare?

10.1.1 Dispersia este cunoscută

Să considerăm cazul când dispersia σ^2 este cunoscută. Valorile x_1, x_2, \dots, x_n sunt valorile variabilelor aleatoare de selecție, independente, X_1, X_2, \dots, X_n , care au aceeași distribuție normală ca X . Deoarece variabila $\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$ este normală cu media m și dispersia $\frac{\sigma^2}{n}$ rezultă că variabila $Z = \frac{\bar{X} - m}{\sigma/\sqrt{n}}$ este normală cu media 0 și dispersia 1. Ca urmare:

$$P(-a \leq Z \leq a) = \Phi(a) - \Phi(-a) = 2\Phi(a)$$

Dacă înlocuim pe Z cu $\frac{\frac{X_1 + X_2 + \dots + X_n}{n} - m}{\sigma/\sqrt{n}}$ găsim:

$$P\left(\frac{X_1 + X_2 + \dots + X_n}{n} - a\frac{\sigma}{\sqrt{n}} \leq m \leq \frac{X_1 + X_2 + \dots + X_n}{n} + a\frac{\sigma}{\sqrt{n}}\right) = 2\Phi(a)$$

Prin urmare intervalul $\left[\frac{x_1 + x_2 + \dots + x_n}{n} - a\frac{\sigma}{\sqrt{n}}, \frac{x_1 + x_2 + \dots + x_n}{n} + a\frac{\sigma}{\sqrt{n}}\right]$ este un interval de încredere pentru m cu nivelul de încredere $2\Phi(a)$. Introduscând pragul de încredere ε , avem $1 - \varepsilon = 2\Phi(a)$ sau $\Phi(a) = \frac{1 - \varepsilon}{2}$. Am demonstrat deci:

Propoziția 10.4 Fie X o variabilă normală de dispersie cunoscută σ^2 și de medie m necunoscută. Dacă $\varepsilon \in (0, 1)$ și $a \in \mathbb{R}_+$, atunci, la selecții de volum n , o condiție suficientă ca intervalul $\left[\bar{X} - a\frac{\sigma}{\sqrt{n}}, \bar{X} + a\frac{\sigma}{\sqrt{n}}\right]$ să fie interval de încredere de nivel $1 - \varepsilon$ pentru media m , este ca a să verifice ecuația $\Phi(a) \geq \frac{1}{2}(1 - \varepsilon)$.

QED.

Observația 10.5 La același prag de încredere ε , creșterea volumului n de selecție conduce la un interval de încredere mai scurt.

Exemplul 10.6 Fie P o populație normală de variantă (dispersie) cunoscută σ^2 și de medie m necunoscută (de estimat). Considerăm selecții de volum fixat n . Vom găsi un interval de încredere, de nivel de încredere 95% pentru medie, dacă alegem astfel pe a încât $\Phi(a) \geq \frac{1}{2} \frac{95}{100}$. Din tabelul pentru Φ găsim $a \geq 1,96$. Deci un interval de încredere de nivel 95% va fi de forma:

$$\left[\bar{X} - 1,96\frac{\sigma}{\sqrt{n}}, \bar{X} + 1,96\frac{\sigma}{\sqrt{n}}\right].$$

Exemplul 10.7 O firmă produce piese cilindrice de diametru $\phi = 10$ mm. Abaterile de la acest diametru impus respectă o lege normală de variație (dispersie) egală cu 0,04 mm (practică a arătat acest lucru). Se face un sondaj pe 100 de piese și se găsește că media de selecție (empirică) este de 10,01 mm. Să se găsească un interval de estimatie pentru media reală cu nivelul de încredere de 90%.

Soluție Aici $n=100$, $\sigma=0,2$, $\bar{X}_{(100)}=0,01$, $(1-\varepsilon)=0,90$, deci $\varepsilon=0,10$. Din tabelul funcției Φ găsim $\Phi(a) \geq \frac{90}{2 \cdot 100} = 0,45$ pentru $\alpha \geq 1,65$. Deci, un interval de estimatie pentru media reală este: $[0,01 - 1,65 \cdot \frac{0,2}{10}, 0,01 + 1,65 \cdot \frac{0,2}{10}] = [9,977, 10,043]$.

Ce informație obține de aici producătorul? El este sigur în proporție de 90% că abaterea medie de la diametru real $\phi=10$ mm este de cel mult 0,043 mm.

10.1.2 Dispersia este necunoscută

Am văzut până acum că dacă dispersia unei populații normale este cunoscută putem estima prin intervale de încredere media populației cu ajutorul v.a. normale standard $Z = \frac{\bar{X}-m}{\sigma/\sqrt{n}}$, unde \bar{X} este media de selecție, iar m este media reală a populației.

Dacă media m nu este cunoscută atunci putem folosi punctul d) al teoremei precedente care spune că variabila

$$\begin{aligned} T &= \frac{\sqrt{(n-1)n} \frac{\bar{X}-m}{\sqrt{(X_1-\bar{X})^2 + (X_2-\bar{X})^2 + \dots + (X_n-\bar{X})^2}}}{\sqrt{n-1}} \\ &= \frac{\bar{X}-m}{S} \end{aligned}$$

are o distribuție Student cu $n-1$ grade de libertate. Aici utilizat notația (vezi lecția 9) $S^2 = \frac{(X_1-\bar{X})^2 + (X_2-\bar{X})^2 + \dots + (X_n-\bar{X})^2}{n}$. Așa cum se vede în lecția 4, densitatea de probabilitate este simetrică față de $x=0$, deci pentru funcția de repartiție $F(t)$ avem relația $F(-t) = 1 - F(t)$. Această observație ne ajută să folosim tabelul II pentru găsirea cuantilelor corespunzătoare acestei distribuții. Pe coloana din stânga a tabelului avem gradele de libertate $\nu = n - 1$ (n volumul selecției), pe prima linie orizontală avem valorile funcției $F(t)$ de la 0,60 până la 0,999. Fie de aflat la $\nu = n - 1 = 4$ valoarea lui a astfel ca $F(a)=0,40$. Avem $1-F(a)=F(-a)=0,60$ și pentru 0,60 avem cuantila în tabel: $-a=0,271$. Deci $a = -0,271$.

Să punem aceste rezultate în următoarea propoziție:

Propoziția 10.8 Fie P o populație normală cu media m și dispersia σ^2 necunoscute. Pentru orice n , pentru un prag $\varepsilon \in (0, 1)$ și $a \in \mathbb{R}_+$, o condiție suficientă ca intervalul

$$\left[\bar{X} - a \frac{S}{\sqrt{n-1}}; \bar{X} + a \frac{S}{\sqrt{n-1}} \right]$$

să fie interval de încredere de nivel $1 - \varepsilon$ (sau de prag ε) pentru media m , este ca a să fie cuantilă de ordin $1 - \varepsilon/2$ a distribuției Student cu $n - 1$ grade de libertate (adică $F(a) = 1 - \varepsilon/2$).

Demonstrație. Relația $P\left(\bar{X} - a\frac{S}{\sqrt{n-1}} \leq m \leq \bar{X} + a\frac{S}{\sqrt{n-1}}\right) = 1 - \varepsilon$ se mai poate scrie

$$P\left(-a \leq \sqrt{n-1} \frac{\bar{X} - m}{S} \leq a\right) = 1 - \varepsilon,$$

$$\text{sau } P(-a \leq T \leq a) = 1 - \varepsilon$$

Dar

$$P(-a \leq T \leq a) = F(a) - F(-a) =$$

$$F(a) - 1 + F(a) = 2F(a) - 1 = 2(1 - \varepsilon/2) - 1 = 1 - \varepsilon$$

QED.

Exemplul 10.9 Presupunem că în exemplul precedent nu cunoaștem dispersia 0,04 mm și că o estimăm cu formula $S^2 = \frac{(X_1 - \bar{X})^2 + (X_2 - \bar{X})^2 + \dots + (X_n - \bar{X})^2}{n}$ găsind-o ca fiind egală cu 0,09 mm. Avem $1 - (\varepsilon/2) = 0,95$, $n = 100$, $S = 0,3$ (în cazul nostru). Cuantila corespunzătoare lui 0,95 o găsim din tabelul cu distribuția Student. La $\nu = n - 1 = 99$ nu găsim date dar putem folosi linia lui $\nu = 120$, deoarece cuantilele vecine diferă puțin unele de altele (pentru același prag bineînțeles). Aici găsim $a = 1,658$. Cu acest a găsim, intervalul de încredere va fi: $\left[10,01 - 1,658\frac{0,3}{\sqrt{99}}; 10,01 + 1,658\frac{0,3}{\sqrt{99}}\right]$, adică $[9,9798, 10,0599]$. Să observăm că $a = 1,65$ pentru situația când am folosit v.a. Z și $a = 1,658$ pentru situația când am folosit v.a. T . Acest lucru se explică, deoarece pentru n mare (mai mare ca 40), în cazul nostru 100, cele două cuantile diferă foarte puțin.

10.2 Intervale de încredere pentru dispersie

Dacă media m a variabilei aleatoare X este cunoscută, atunci putem folosi punctul b) al teoremei precedente care spune că variabila $\left(\frac{X_1 - m}{\sigma}\right)^2 + \left(\frac{X_2 - m}{\sigma}\right)^2 + \dots + \left(\frac{X_n - m}{\sigma}\right)^2$ este de tip $H(n)$ iar dacă media m nu este cunoscută putem folosi punctul c) al teoremei, anume că variabila aleatoare $\left(\frac{X_1 - \bar{X}}{\sigma}\right)^2 + \left(\frac{X_2 - \bar{X}}{\sigma}\right)^2 + \dots + \left(\frac{X_n - \bar{X}}{\sigma}\right)^2$ este de tip $H(n - 1)$, în scopul de a determina intervale de încredere pentru dispersie.

$$\frac{nS^2}{\sigma^2} = \chi_{(n-1)}^2 \quad (10.2)$$

adică $\frac{nS^2}{\sigma^2}$ este v.a χ^2 cu $n - 1$ grade de libertate. Dacă în loc de S^2 se folosește estimatorul nedeplasat S'^2 se obține formula

$$\frac{(n-1)S'^2}{\sigma^2} = \chi_{(n-1)}^2 \quad (10.3)$$

Teorema 10.10 Fie $\varepsilon \in (0, 1)$. Un interval de încredere de nivel $100(1 - \varepsilon)$ procente pentru dispersia σ^2 a unei populații normale cu media cunoscută m , în cazul selecțiilor de volum n , este

$$\left[\frac{(n-1)S'^2}{b}, \frac{(n-1)S'^2}{a} \right] \quad (10.4)$$

unde a este cuantila de ordin $\varepsilon/2$ și b este cuantila de ordin $1 - (\varepsilon/2)$ a distribuției χ^2 cu $n - 1$ grade de libertate.

Demonstratie Notăm cu $F(t)$ funcția de repartiție a v.a. $\chi^2_{(n-1)}$. Avem $F(a) = \varepsilon/2$ și $F(b) = 1 - (\varepsilon/2)$. Atunci $P\left(a \leq \frac{(n-1)S'^2}{\sigma^2} \leq b\right) = F(b) - F(a) = 1 - (\varepsilon/2) - (\varepsilon/2) = 1 - \varepsilon$. Dar, din $a \leq \frac{(n-1)S'^2}{\sigma^2} \leq b$ găsim că $\frac{(n-1)S'^2}{b} \leq \sigma^2 \leq \frac{(n-1)S'^2}{a}$. Prin urmare, intervalul din (10.4) acoperă pe σ^2 cu probabilitatea $1 - \varepsilon$.

QED.

Exemplul 10.11 Media erorilor de măsurare a lungimilor unor baghete metalice este de 3 mm. Presupunem că aceste erori respectă legea normală cu media 3 mm și dispersia necunoscută. Se face o selecție de volum 4: $\{-1, 4, 4, 1\}$. Se cere un interval de estimatie pentru σ^2 cu pragul de încredere de 90%.

Soluție În cazul nostru aplicăm Teorema 10.9 cu $1 - \varepsilon = 0,90$, deci $\varepsilon = 0,10$. Căutăm cuantilele pentru $\varepsilon/2 = 0,05$ și $1 - (\varepsilon/2) = 0,95$, când $n-1=3$ (grade de libertate). Găsim $a=0,351846$ și $b=7,81473$ în Tabelul III. Calculăm acum $S'^2 = \frac{1}{3}((-1-3)^2 + (4-3)^2 + (4-3)^2 + (1-3)^2) = \frac{22}{3}$. Intervalul va fi deci $\left[\frac{22}{7,81}, \frac{22}{0,35}\right] = [2,81; 62,85]$. Se observă că intervalul este destul de mare, deci precizia pentru σ^2 este mică, chiar dacă apare cu probabilitate mare!

10.3 Intervale de încredere pentru cîtul a două dispersii

Fie acum două populații distincte, normal distribuite. Facem o selecție de volum n_1 din prima populație și o selecție de volum n_2 din a doua populație. Stim din formula (10.3) că

$$\frac{(n_1 - 1)S_1'^2}{\sigma_1^2} = \chi^2_{(n_1-1)} \quad (10.5)$$

$$\frac{(n_2 - 1)S_2'^2}{\sigma_2^2} = \chi^2_{(n_2-1)}$$

unde $\sigma_1, S_1'^2$ și $\sigma_2, S_2'^2$ sunt dispersiile și dispersiile de selecție modificate pentru cele două populații. Notăm cu $\nu_1 = n_1 - 1$ și cu $\nu_2 = n_2 - 1$. Notăm cu F (de la Fischer) v.a.

$$\frac{S_1'^2/\sigma_1^2}{S_2'^2/\sigma_2^2} = \frac{\left[\chi_{(\nu_1)}^2/\nu_1\right]}{\left[\chi_{(\nu_2)}^2/\nu_2\right]} \quad (10.6)$$

Această v.a. are o densitate de probabilitate ce depinde de doi parametri ν_1 și ν_2 iar formula ei este complicată din punct de vedere matematic (vezi lecția 4). Ea apare ca un cât de v.a. χ^2 , înmulțit cu un număr care depinde de ν_1 și ν_2 , adică ν_2/ν_1 . Vom mai nota o asemenea variabilă F_{ν_1, ν_2} pentru a pune în evidență cei doi parametri de care depinde. Tabelul IV ne furnizează fractilele acestei distribuții numai pentru ordinele 0,95; 0,975 și 0,99. Pe coloane apar valorile parametrului ν_1 și pe linii apar valorile parametrului ν_2 . De exemplu, pentru $n_1 = 10$, $n_2 = 6$, $\nu_1 = 9$, $\nu_2 = 5$, presupunem că după selecție am obținut $F=7$. Ne uităm la cuantila de ordin 0,95 și găsim valoarea 3,48. Valoarea selecției, 7, este mai mare decât 3,48, deci cade în partea opusă, adică în partea cu probabilitatea 5%. Prin urmare "inferența" noastră asupra câtului $\frac{S_1'^2}{S_2'^2}$ nu este adevărată cu 95% probabilitate. În practică este utilă relația:

$$P(F_{(\nu_1, \nu_2)} \geq c) = P\left(F_{(\nu_2, \nu_1)} \leq \frac{1}{c}\right) \quad (10.7)$$

Aici $F_{(\nu_2, \nu_1)}$ se numește *inversa* v.a. $F_{(\nu_1, \nu_2)}$.

Din aceste observații rezultă imediat:

Teorema 10.12 *Avem relația*

$$P\left(\frac{aS_2'^2}{S_1'^2} \leq \frac{\sigma_2^2}{\sigma_1^2} \leq \frac{bS_2'^2}{S_1'^2}\right) = P\left(a \leq \frac{S_1'^2/\sigma_1^2}{S_2'^2/\sigma_2^2} \leq b\right) = F_{\nu_1, \nu_2}(b) - F_{\nu_1, \nu_2}(a) = 1 - \varepsilon$$

dacă a și b sunt alese astfel ca $F_{\nu_1, \nu_2}(b) = 1 - \frac{\varepsilon}{2}$ și $F(a) = \frac{\varepsilon}{2}$. În aceste condiții un interval de încredere pentru $\frac{\sigma_1^2}{\sigma_2^2}$, cu nivelul de încredere $(1 - \varepsilon)$ este $\left[\frac{aS_2'^2}{S_1'^2}, \frac{bS_2'^2}{S_1'^2}\right]$.

QED.

10.4 Intervale de încredere în cazul unor selecții mari

Dacă $f(x, \theta)$ este densitatea de probabilitate a variabilei aleatoare X atunci din $\int_{-\infty}^{\infty} f(x, \theta) dx = 1$ rezultă prin derivare în raport cu θ că $\int_{-\infty}^{\infty} \frac{\partial f}{\partial \theta}(x, \theta) d\theta = 0$ sau $\int \frac{\partial \ln(f(x, \theta))}{\partial \theta} f(x, \theta) dx = 0$. Deci variabila aleatoare $\frac{\partial \ln(f(x, \theta))}{\partial \theta}$ are media 0 și dispersia $\int_{-\infty}^{\infty} \left(\frac{\partial \ln(f(x, \theta))}{\partial \theta}\right)^2 f(x, \theta) dx$. Presupunând că dispersia este finită, rezultă din legea limită centrală (vezi lecția 5) că pentru n

mare, variabila aleatoare

$$\frac{\frac{\partial \ln(f(X_1, \theta))}{\partial \theta} + \frac{\partial \ln(f(X_1, \theta))}{\partial \theta} + \dots + \frac{\partial \ln(f(X_1, \theta))}{\partial \theta}}{\sqrt{n \int_{-\infty}^{\infty} \left(\frac{\partial \ln(f(x, \theta))}{\partial \theta} \right)^2 f(x, \theta) dx}}$$

unde X_1, X_2, \dots, X_n sunt variabilele de selecție asociate cu X , are o distribuție aproximativ normală cu media 0 și dispersia 1. Avem deci:

$$\Pr ob \left(a < \frac{\frac{\partial \ln(f(X_1, \theta))}{\partial \theta} + \frac{\partial \ln(f(X_1, \theta))}{\partial \theta} + \dots + \frac{\partial \ln(f(X_1, \theta))}{\partial \theta}}{\sqrt{n \int_{-\infty}^{\infty} \left(\frac{\partial \ln(f(x, \theta))}{\partial \theta} \right)^2 f(x, \theta) dx}} < b \right) \approx \int_a^b \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx \quad (10.8)$$

Un interval de încredere pentru θ , cu nivelul de încredere α , se poate obține pentru n mare astfel:

- Se determină prin n experiențe independente valorile x_1, \dots, x_n pentru X_1, \dots, X_n .
- Se determină a și b astfel ca $\int_a^b \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx = \alpha$.
- Din formula 10.8 rezultă că mulțimea valorilor θ care verifică inegalitatea

$$a < \frac{\frac{\partial \ln(f(X_1, \theta))}{\partial \theta} + \frac{\partial \ln(f(X_1, \theta))}{\partial \theta} + \dots + \frac{\partial \ln(f(X_1, \theta))}{\partial \theta}}{n \cdot \sqrt{\int_{-\infty}^{\infty} \left(\frac{\partial \ln(f(x, \theta))}{\partial \theta} \right)^2 f(x, \theta) dx}} < b$$

este o mulțime de încredere pentru θ cu încrederea α . În unele cazuri această mulțime este un interval.

10.5 Rezumat

1. Fie $\varepsilon > 0$, $\varepsilon \in (0, 1)$. Vom numi interval de încredere de prag ε (sau de nivel de încredere $1 - \varepsilon$) pentru parametrul λ două statistici $\Lambda_1 = f_1(X_1, \dots, X_n)$ și $\Lambda_2 = f_2(X_1, \dots, X_n)$ astfel încât $P(\Lambda_1 \leq \lambda \leq \Lambda_2) \geq 1 - \varepsilon$.

Pentru selecții efective x_1, \dots, x_n , vom nota valoarea statisticii Λ_1 cu $\hat{\lambda}_1$ și a statisticii Λ_2 cu $\hat{\lambda}_2$. Intervalul numeric $[\hat{\lambda}_1, \hat{\lambda}_2]$ este considerat încă ca interval de încredere de nivel $1 - \varepsilon$ (se mai spune de prag ε) pentru parametrul estimat λ .

2. Dacă există două statistici $Y = f(X_1, \dots, X_n)$ și $Z = g(X_1, \dots, X_n)$ astfel încât v.a. $T = \frac{Y - \mu}{Z}$ să fie normală redusă sau Student cu d grade de libertate și t_ε un număr pozitiv astfel încât $P(|T| > t_\varepsilon) \leq \varepsilon$, atunci $[Y - t_\varepsilon |Z|, Y + t_\varepsilon |Z|]$ este un interval de încredere de prag ε pentru media μ , adică:

$$P(Y - t_\varepsilon |Z| \leq \mu \leq Y + t_\varepsilon |Z|) \geq 1 - \varepsilon$$

• Fie acum de estimat σ^2 a v.a. X . Presupunem că am găsit o statistică Y a.î. $T = \frac{Y}{\sigma^2}$ are distribuția χ^2 cu d grade de libertate și fie două numere t'_ε și t''_ε astfel încât $P(T \geq t'_\varepsilon) \leq \frac{\varepsilon}{2}$ și $P(T \leq t''_\varepsilon) \leq \frac{\varepsilon}{2}$. Atunci $\left[\frac{Y}{t'_\varepsilon}, \frac{Y}{t''_\varepsilon}\right]$ este un interval de încredere pentru σ^2 de prag ε : $P\left(\frac{Y}{t'_\varepsilon} \leq \sigma^2 \leq \frac{Y}{t''_\varepsilon}\right) \geq 1 - \varepsilon$.

Se alege Y a.î. să aibă cât mai multe grade de libertate.

FORMULE UTILIZATE FRECVENT

In formulele de mai jos nivelul de încredere este $1 - \varepsilon$, iar rezultatele a n măsurători independente ale unei caracteristici numerice cu distribuție normală sunt x_1, x_2, \dots, x_n .

1. Un interval de încredere pentru media m a unei variabile aleatoare normale, dacă se cunoaște dispersia σ^2 este:

$$\left[\frac{x_1 + x_2 + \dots + x_n}{n} - a \frac{\sigma}{\sqrt{n}}, \frac{x_1 + x_2 + \dots + x_n}{n} + a \frac{\sigma}{\sqrt{n}} \right]$$

unde a se alege astfel ca $\Phi(a) = 0,5 - \frac{\varepsilon}{2}$.

2. Un interval de încredere pentru media m a unei v.a. normale, dacă nu se cunoaște dispersia, este:

$$\left[m^* - a \frac{s^*}{\sqrt{n-1}}, m^* + a \frac{s^*}{\sqrt{n-1}} \right]$$

unde $m^* = \frac{x_1 + x_2 + \dots + x_n}{n}$, $s^* = \sqrt{\frac{(x_1 - m^*)^2 + (x_2 - m^*)^2 + \dots + (x_n - m^*)^2}{n}}$, iar a se alege astfel ca $F(a) = 1 - \frac{\varepsilon}{2}$, F fiind funcția de repartiție a unei variabile Student cu $n - 1$ grade de libertate.

3. Un interval de încredere pentru dispersia σ^2 a unei v.a. normale este:

$$\left[\frac{(x_1 - m^*)^2 + \dots + (x_n - m^*)^2}{b}, \frac{(x_1 - m^*)^2 + \dots + (x_n - m^*)^2}{a} \right]$$

unde $F(a) = \frac{\varepsilon}{2}$, iar $F(b) = 1 - \frac{\varepsilon}{2}$, F fiind funcția de repartiție a unei variabile χ^2 cu $n - 1$ grade de libertate.

4. Un interval de încredere pentru câtu dispersiilor $\frac{\sigma_2^2}{\sigma_1^2}$ a două v.a. independente este:

$$\left[a \cdot \frac{\frac{\sum_{i=1}^{n_2} (y_i - m'^*)^2}{n_2 - 1}}{\frac{\sum_{i=1}^{n_1} (x_i - m^*)^2}{n_1 - 1}}, b \cdot \frac{\frac{\sum_{i=1}^{n_2} (y_i - m'^*)^2}{n_2 - 1}}{\frac{\sum_{i=1}^{n_1} (x_i - m^*)^2}{n_1 - 1}} \right]$$

unde $m^* = \frac{x_1 + x_2 + \dots + x_{n_1}}{n_1}$, $m'^* = \frac{y_1 + y_2 + \dots + y_{n_2}}{n_2}$, n_1 și n_2 sunt volumele celor două selecții, $\nu_1 = n_1 - 1$, $\nu_2 = n_2 - 1$, iar a și b sunt alese astfel ca $1 - \frac{\varepsilon}{2} = F_{\nu_1, \nu_2}(b)$, $\frac{\varepsilon}{2} = F_{\nu_1, \nu_2}(a)$

10.6 Exerciții rezolvate

1. Atunci când se nasc 2 copii simultan (gemeni) probabilitatea ca ei să fie gemeni adevărați este λ . Se presupune că:

- a) 2 gemeni adevărați au întotdeauna același sex și probabilitatea ca ei să fie băieți este $\frac{1}{2}$;
 b) 2 gemeni falsi au sexe diferite și probabilitatea ca unul dintre ei să băiat este $\frac{1}{2}$;

i) În cursul nasterii a 2 gemeni se consideră evenimentele: $A = (2 \text{ băieți})$; $B = (2 \text{ fete})$; $C = (1 \text{ băiat și } 1 \text{ fată})$. Calculați în funcție de λ probabilitățile $p(A)$, $p(B)$, $p(C)$.

ii) În cursul a 1000 de nașteri se realizează evenimentul C de 328 de ori. Dați pentru λ un interval de încredere de prag $\varepsilon = 0,05$.

iii) Observăm acum n nașteri de gemeni și notăm cu Y_C numărul de realizări ale evenimentului C . Ce lege guvernează v.a. Y_C ? Definiți cu ajutorul lui Y_C un esantion nedeplasat Z pentru λ . Calculați varianta lui Z . Dați pentru n mare o condiție independentă de λ și suficientă pentru a putea defini cu ajutorul lui Z un interval de încredere de prag $\varepsilon = 0,05$ a cărui lungime să fie mai mică decât un $a \in \mathbb{R}$, dat. Caz particular $a = \frac{1}{100}$.

Soluție i) Notăm cu V evenimentul: \ll gemeni adevărați \gg și cu F : \ll gemeni falsi \gg . Atunci $A = (A \cap V) \cup (A \cap F)$ și $p(A) = p(V) \cdot p_V(A) + p(F) \cdot p_F(A) = \lambda \cdot \frac{1}{2} + (1 - \lambda) \cdot \frac{1}{4} = \frac{\lambda+1}{4}$. La fel $p(B) = \frac{\lambda+1}{4}$ și $p(C) = \frac{1-\lambda}{2}$.

ii) Fie \bar{X} v.a. care are valoarea 1, dacă se realizează evenimentul C , și 0 altfel. Este clar că $M(X) = \mu = \frac{1-\lambda}{2}$. $\bar{X} = \frac{1}{n} Y_C = \frac{328}{1000} = 0,328$. Cum $X = X_i$, avem că $X_i^2 = X_i$, deci $S^2 = \frac{1}{n} \sum_i X_i^2 - (\bar{X})^2 = \bar{X}(1 - \bar{X})$. Cum $T = \frac{\bar{X} - M(X)}{\sqrt{S^2/n}}$ este practic normală redusă, egalitatea $P(|T| \geq 1,96) = 0,05$ dă pentru μ intervalul de încredere cerut: $\bar{X} - 1,96 \frac{S}{\sqrt{n}} \leq \mu \leq \bar{X} + 1,96 \frac{S}{\sqrt{n}}$, sau $0,299 \leq \frac{1-\lambda}{2} \leq 0,357$. De aici rezultă intervalul de încredere căutat pentru λ : $0,286 \leq \lambda \leq 0,402$.

iii) Y_C este binomială cu $p = \frac{1-\lambda}{2}$. Avem deci $M(Y_C) = \frac{n(1-\lambda)}{2}$ și $D(Y_C) = \frac{n(1-\lambda^2)}{4}$. Egalitatea $\lambda = 1 - \frac{2}{n} M(Y_C)$ dă pentru λ estimatorul nedeplasat $Z = 1 - \frac{2}{n} Y_C$ de dispersie $\frac{4}{n^2} D(Y_C) = \frac{1-\lambda^2}{n}$. Când $n \rightarrow \infty$, Y_C este practic gaussiană (normală), deci și Z este la fel. Considerăm deci $T = \frac{Z - \lambda}{\sqrt{(1-\lambda^2)/n}}$ care este gaussiană redusă. Egalitatea $P(|T| \geq 1,96) = 0,05$ ne per-

mite să scriem $P(|Z - \lambda| \geq 1,96 \sqrt{(1-\lambda^2)/n}) = 0,05$, și cum $\sqrt{(1-\lambda^2)/n} \leq \sqrt{1/n}$ avem $P(|Z - \lambda| \geq 1,96/\sqrt{n}) < 0,05$, de unde $Z - \frac{1,96}{\sqrt{n}} \leq \lambda \leq Z + \frac{1,96}{\sqrt{n}}$. Lungimea lui va fi mai mică decât a atunci când $a \geq \frac{2 \times 1,96}{\sqrt{n}}$, sau $n \geq \left(\frac{3,92}{a}\right)^2$. Pentru $a = \frac{1}{100}$ găsim $n \geq 153664$.

2. Se măsoară forța de compresiune X (în Kg/cm^3) a cimentului din care sunt confecționați cilindri mici, limită de la care ei se sparg. Pentru $n=10$ cilindri se observă următoarele presiuni:

19,6	19,9	20,4	19,8	20,5
21,0	18,5	19,7	18,4	19,4

Presupunem că X are o lege gaussiană (normală).

i) Dati un interval de încredere de prag $\varepsilon = 0,1$ pentru $M(X)$.

ii) Dati o estimare nedeplasată $\hat{\sigma}^2$ pentru varianta σ^2 a v.a. X , găsiți apoi un interval de încredere de prag $0,1$ pentru σ^2 .

iii) Presupunem că $\sigma^2 = 0,69$. Găsiți pentru $M(X)$ un nou interval de încredere de prag $0,1$. Comparati cu rezultatul de la 1).

Solutie i) Calculăm $\bar{X} = 19,72$ și $nS^2 = 6,0960$. $T = \frac{\bar{X} - M(X)}{\sqrt{S^2/(n-1)}}$ este o v.a. Student cu $n - 1 = 9$ grade de libertate, avem $P(|T| > t_\varepsilon) = 0,1$ pentru $t_\varepsilon = 1,833$. Intervalul de încredere cerut este deci $\bar{X} - t_\varepsilon \frac{S}{\sqrt{n-1}} \leq M(X) \leq \bar{X} + t_\varepsilon \frac{S}{\sqrt{n-1}}$, sau $19,243 \leq M(X) \leq 20,197$.

ii) O estimatie nedeplasată a lui σ^2 este $\hat{\sigma}^2 = \frac{nS^2}{n-1} = 0,6773$. Pe de altă parte stim că $U = \frac{nS^2}{\sigma^2}$ are distributia χ^2 (Pearson) cu $n - 1 = 9$ grade de libertate. Avem deci $P(U > t'_\varepsilon) = 0,05$ pentru $t'_\varepsilon = 3,33$. De aici găsim pentru σ^2 intervalul de încredere: $\frac{nS^2}{t'_\varepsilon} \leq \sigma^2 \leq \frac{nS^2}{t''_\varepsilon}$, adică $0,36 \leq \sigma^2 \leq 1,83$.

iii) Dacă stim dispersia $\sigma^2 = 0,69$ putem folosi faptul că \bar{X} este gaussiană $N\left(\mu, \frac{\sigma^2}{n}\right)$ și deci $T' = \frac{\bar{X} - M(X)}{\sqrt{\sigma^2/n}}$ este gaussiană redusă. Egalitatea $P(|T'| > t'_\varepsilon) = 0,1$ ne conduce la $t'_\varepsilon = 1,6449$.

Prin urmare, găsim un interval de încredere de prag $\frac{1}{10}$ pentru $M(X)$: $\bar{X} - t'_\varepsilon \sqrt{\sigma^2/n} \leq M(X) \leq \bar{X} + t'_\varepsilon \sqrt{\sigma^2/n}$, sau $19,287 \leq M(X) \leq 20,153$. Acest interval este mai mic decât acela găsit la 1) deoarece acum avem dispersia dată.

10.7 Exerciții

1. Notăm cu X vârsta în ani la care un om devine bunic. Presupunem că X are distributia normală cu varianta 225. 9 persoane luate la întâmplare au declarat că au devenit bunici la: 42, 56, 68, 56, 48, 36, 45, 71 și 64 ani.

a) Calculati media și dispersia de selectie.

b) Găsiți un interval de încredere de 80% pentru medie.

c) Găsiți un interval de încredere de 95% pentru medie.

2. În cadrul unui proces de estimare a mediei unei populatii oarecare, un statistician vrea ca probabilitatea ca media de selectie să difere de media adevărată cu mai puțin de $0,2\sigma$ să fie mai mare de 0,95.

a) Ce volum de selectie trebuie să folosească?

b) Dar dacă volumul de selectie este 100, care este marja de aproximare (în σ unități) a mediei reale cu media de selectie, pentru ca să se obțină un prag de încredere de 0,95?

c) Dar dacă se știe că populația este normală care trebuie să fie volumul de selecție ca

$$\Pr ob(|\bar{X} - m| \leq 0,2 \cdot \sigma) \geq 0,95$$

?

3. Fie distributia student T cu 12 grade de libertate.

- a) găsiți fractile pentru 0,10; 0,60 și 0,95.
- b) găsiți media și dispersia.
- c) $P(T < -0,695)$.
- d) $P(-2,179 < T < 1,356)$.

4. Fie distributia T cu 6 grade de libertate.

- a) găsiți $P(T \geq 1)$.
- b) găsiți fractila pentru 0,20.

5. Fie P o populație distribuită normal. Se face selecția: -2, -1, 0, 0, 1, 2, 2, 3, 0. Să se găsească un interval de încredere pentru media populației cu pragul de încredere de 80%, folosind distributia T .

6. Dintr-un lot de 100 piese, 4 sunt găsite defecte. Presupunem că procesul de producție se comportă ca un proces de tip Bernoulli. Determinați intervale de încredere de 95% și de 99% pentru probabilitatea ca luând o piesă la întâmplare ea să fie defectă.

7. Se urmărește indicatorul Dow-Jones (indicator la bursa americană) la stocurile industriale de la o zi la alta. Se presupune că acesta variază după o distribuție normală. Se face un sondaj pe parcursul a 81 de zile și se obține o medie de selecție de 0,20 și o dispersie de selecție de 1,50. Găsiți un interval de încredere pentru medie de 90%.

8. Dintr-o populație normală ($m=24$, $\sigma^2=9$) se ia un esanșion de volum 8. Notăm cu W suma pătratelor celor 8 valori standardizate $\left(Z_i^2 = \left(\frac{X_i - m}{\sigma}\right)^2\right)$. Găsiți: a) $P(W \geq 20,09)$; b) $P(2,73 < W < 5,07)$; c) $P(W > 2,18)$.

11. Pentru distributia F cu 5 grade de libertate la numărător și cu 8 grade de libertate la numitor, găsiți : a) fractilele de ordin 0,025 și 0,99; b) $P(F \leq 3,69)$; c) $P(F \geq 1,22)$.

12. Durata de viață a unui bec electric este o variabilă aleatoare normală. Testul pe 16 becuri a arătat o valoare medie de viață de 3000 ore și o abatere standard $\sigma^* = 20$. Să se determine intervale de încredere pentru medie și abatere standard cu pragul de risc $\varepsilon = 0,1$.

13. Un producător de rulmenți pretinde că diametrul mediu al rulmenților, în mm este de 10 mm cu o dispersie de 10^{-4} . Admitem că diametrul este o v.a. normală. Pe un lot de 20 rulmenți măsurați s-a găsit că diametrul mediu 9,98 mm și o dispersie empirică 0,0002. Să se determine intervale de încredere pentru diametrul mediu și dispersie cu încrederea 0,99. Valorile pretinse de fabricant se află în aceste intervale?

Lecția 11

Ipoteze statistice. Teste statistice

11.1 Ipoteze și testarea lor

În continuare vom face ipoteze asupra parametrilor unor populații, știind în prealabil clasa de distribuții din care fac parte (de exemplu: normală, Bernoulli, Poisson, etc.). Vom folosi rezultatele obținute în Lectia 10 asupra estimării prin intervale de încredere a unor parametri remarcabili pentru distribuții cunoscute (media și dispersia pentru populații normale, de exemplu).

O ipoteză statistică este o ipoteză făcută asupra unor însusiri statistice ale unei populații P . Ea este *simplă*, dacă se referă la întreaga informație care determină distribuția populației, de exemplu ipoteza:

H: populația este normală de medie $m=10$ și dispersie $\sigma^2=225$, sau

H: populația este Bernoulli cu $p=0,3$.

Ipoteza poate fi *compusă* dacă se referă numai la o parte din informațiile ce ar putea determina distribuția populației. Iată un exemplu de distribuție compusă:

H: populația este normală de medie 40, sau

H: populația este Poisson (nu facem nici o ipoteză asupra mediei λ).

În cazul ipotezelor compuse, ceilalți parametri care împreună cu cei testați ar duce la determinarea completă a distribuției, se estimează dintr-o selecție (sau mai multe) făcută asupra populației.

O ipoteză poate fi *exactă*, de exemplu ipoteza H: media populației Poisson este $\lambda=3$, sau poate fi *inexactă*: H: media populației normale este $m \geq 5$.

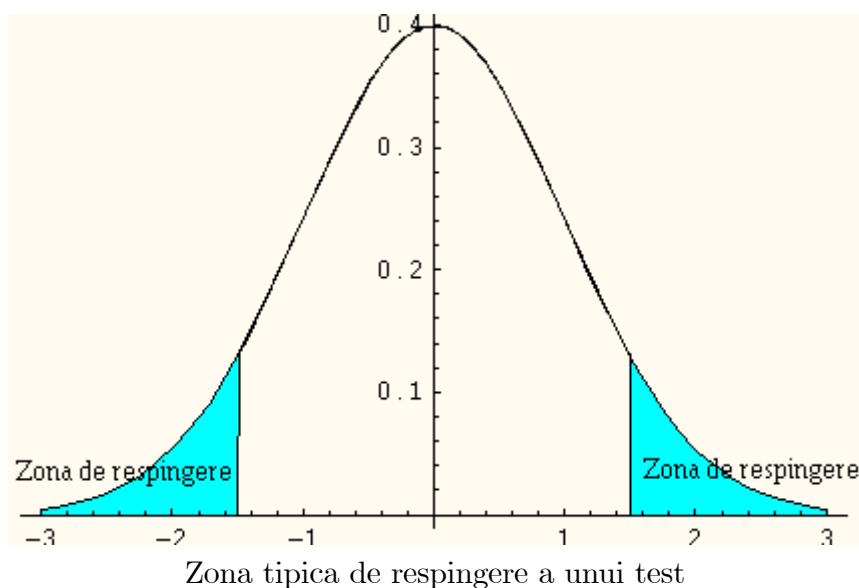
În aparență noi lucrăm cu o singură ipoteză H. De fapt lucrăm cu două ipoteze: $H=H_0$ și H_1 , ipoteza contrară ipotezei H_0 . În cele ce urmează vom considera două ipoteze alternative H_0 și H_1 . Nu întotdeauna ipoteza H_1 reprezintă negația logică obișnuită a ipotezei H_0 . De exemplu, H_0 : media populației este $m=30$, H_1 : media populației s-a micșorat, adică este $m<30$.

Operația de comparare a două ipoteze statistice în lumina informațiilor furnizate de selecție

se numeste *test statistic*. Dacă testul statistic se referă la unul sau mai multi parametri ce apar în legea ce definește populația spunem că testul este *parametric*. Dintre cele două ipoteze H_0 și H_1 , una dintre ele, notată cu H_0 , ocupă locul central: testăm pe H_0 "împotriva" alternativei H_1 . La finalul testării statisticianul, fie acceptă pe H_0 , fie că respinge pe H_0 în favoarea ipotezei H_1 . Oricum el trebuie să ia "o decizie". Fiecare test statistic implică o statistică de selecție, adică o funcție continuă de tipul $g(X_1, \dots, X_n)$, unde n este volumul selecției, iar X_1, \dots, X_n sunt variabilele de selecție. Anumite valori ale statisticii conduc la acceptarea ipotezei H_0 , alte valori ale ei conduc la respingerea acestei ipoteze. Vom vorbi deci de un *domeniu de respingere (de neacceptare)*. *Regula de decizie* este dată de fapt de specificarea domeniului de respingere al ipotezei H_0 , deoarece se consideră că domeniul complementar domeniului de respingere este exact domeniul de acceptare pentru ipoteza H_0 .

Mai exact, dacă ipoteza H_0 , numita și *ipoteza nula*, este adevărata atunci, în urma unei selecții concrete, este foarte probabil ca valoarea calculată pentru v.a. $g(X_1, \dots, X_n)$ să se găsească într-un interval "de probabilitate mare". Acest lucru se întâmplă deoarece statistica $g(X_1, \dots, X_n)$ are o repartitie bine determinată de ipoteza H_0 și eventual de unele estimări făcute în urma selecției concrete ce apare în problema. Deci noi trebuie să stabilim o *zona de acceptare*, adică o submulțime A din \mathbb{R} a.i. probabilitatea ca o valoare a v.a. $g(X_1, \dots, X_n)$ să aparțină mulțimii A să fie destul de mare (de obicei se ia ca fiind $\geq 0,9$). Mulțimea $\mathbb{R} \setminus A$ se zice *zona de respingere* și probabilitatea ca v.a. $g(X_1, \dots, X_n)$ să ia o valoare în $\mathbb{R} \setminus A$ este foarte mică ($\leq 0,1$). Numărul $\varepsilon = \text{prob}(g(X_1, \dots, X_n) \in \mathbb{R} \setminus A)$ se numeste *prag de semnificație* pentru testul pe care îl vom constitui, iar statistica $g(X_1, \dots, X_n)$ este o v.a. care depinde de v.a. de selecție X_1, \dots, X_n și este legată de ipoteza H_0 . De exemplu, dacă H_0 se referă la media populației, $g(X_1, \dots, X_n)$ va fi media de selecție standardizată, adică $Z = \frac{\bar{X} - m}{\sigma/\sqrt{n}}$, unde m este media ce rezultă din ipoteza H_0 , σ este deviația standard (presupusă cunoscută), n este volumul selecției, iar $\bar{X} = (X_1 + X_2 + \dots + X_n)/n$, este v.a. media de selecție. *Ipoteza alternativă*, H_1 este ipoteza ce rezultă *natural în urma negării ipotezei H_0* . De exemplu, dacă H_0 este " $m = 30$ " și noi știm sigur că media *nu poate să crească* în urma experimentului ce apare în problema, atunci H_1 va fi " $m < 30$ ". Dacă nu știm nimic despre modul în care se schimbă media, este natural să considerăm ipoteza alternativă ca fiind " $m \neq 30$ ", adică " $m < 30$ " sau " $m > 30$ ". Iată deci cum funcționează în general un *test parametric* referitor la parametrul θ al unei populații P :

- 1) construim ipotezele H_0 (*ipoteza nula*) și H_1 (*ipoteza alternativă*) asupra parametrului θ .
- 2) construim v.a. $g(X_1, \dots, X_n)$ care are o distribuție (repartitie) cunoscută dacă considerăm pe H_0 adevărata.
- 3) precizăm pragul de semnificație ε pentru v.a. $g(X_1, \dots, X_n)$ (ε este mic, $\leq 0,1$).
- 4) reprezentăm grafic (schematic și nu exact) zonele de respingere și respectiv de acceptare, $\rho(g)$ = densitatea de probabilitate a v.a. g .



Aria hasurata este $\varepsilon/2 + \varepsilon/2 = \varepsilon = \text{prob}(g \text{ sa ia valori in zona de respingere})$.

5) calculam $g_{calc.} = g(x_1, \dots, x_n)$ pentru valorile efective ale unei selectii furnizate de problema. Daca $g_{calc.} \in \{\text{zonei de acceptare}\}$ vom spune ca acceptam ipoteza H_0 cu pragul de semnificatie ε . Daca $g_{calc.} \in \{\text{zonei de respingere}\}$ acceptam ipoteza alternativa H_1 cu pragul de semnificatie ε .

11.1.1 Testul Z privind media unei populatii normale cu dispersia cunoscuta σ^2

Vrem sa testam ipoteza:

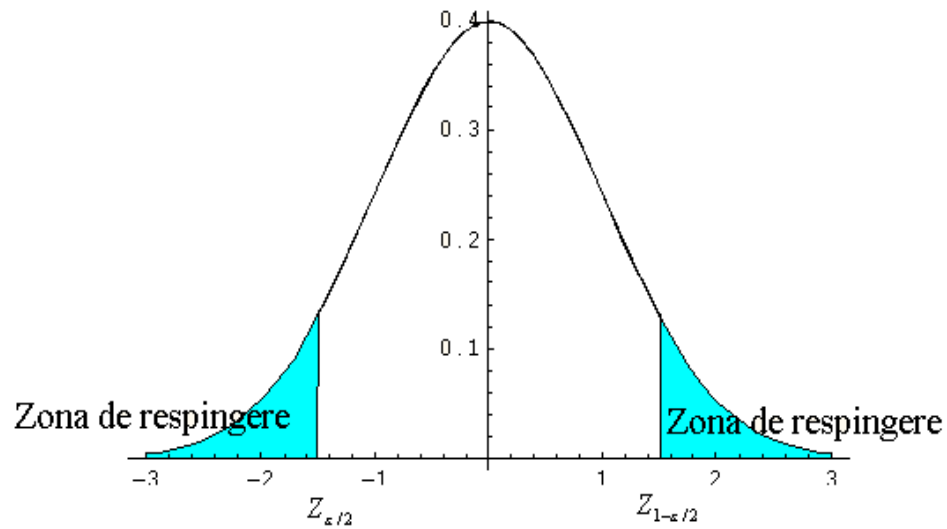
- $H_0: m = m_0$, m_0 specificat (m este media populatiei)

$H_1: m \neq m_0$

- Consideram statistica $Z = \frac{\bar{X} - m_0}{\sigma/\sqrt{n}}$. Stim ca Z este o v.a. normala redusa (are media 0 si dispersia 1), pentru n mare. Aici σ este precizat de problema.

- Fie $\varepsilon \in (0; 0, 1]$ pragul de semnificatie ales ($\varepsilon = 0, 05; 0, 01$, etc.)

-



Zona de respingere pentru un test bilateral

Calculam pe $Z_{\epsilon/2}$ ca fiind cuantila de ordin $\epsilon/2$, adica $F(Z_{\epsilon/2}) = 1 - \epsilon/2$ (vezi TABELUL I)

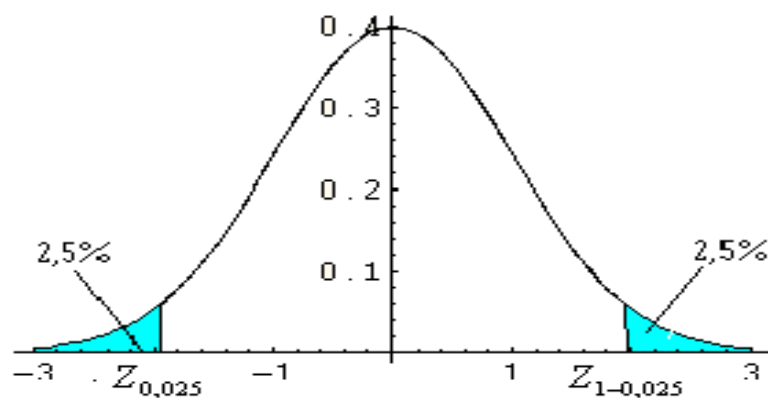
- calculam $Z_{calc} = \frac{\bar{X} - m_0}{\sigma/\sqrt{n}}$ pentru selectia din problema.

Daca $|Z_{calc}| \leq Z_{\alpha/2}$, acceptam ipoteza H_0 cu pragul de semnificatie ϵ .

Daca $|Z_{calc}| > Z_{\alpha/2}$, respingem ipoteza H_0 (acceptam H_1) cu pragul de semnificatie ϵ .

Exemplul 11.1 Testati cu un prag (nivel) de semnificatie de 5% daca o selectie de volum 1, $x_1 = 172$ provine dintr-o populatie normala cu media $m = 150$ si dispersia fixata (cunoscuta) $\sigma^2 = 100$.

Solutie $H_0: m = 150$; $H_1: m \neq 150$; $Z = \frac{\bar{X} - m}{\sigma/\sqrt{n}} = \frac{X - m}{\sigma}$, deoarece $n = 1$ si $\bar{X} = X_1 = X$.
 $F(Z_{0,025}) = 1 - 0,025 = 0,975$.



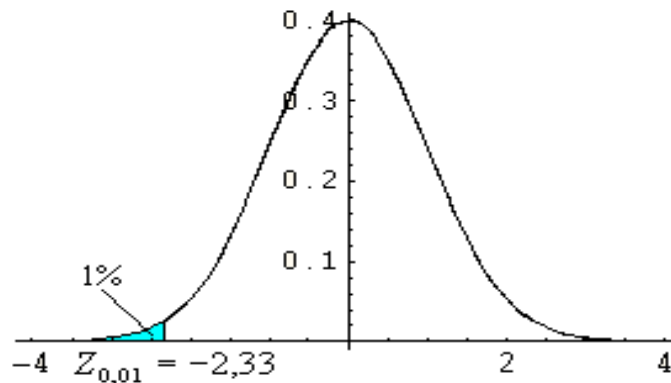
Deci, din TABELUL I gasim ca $Z_{0,025}$ (cuantila de ordin 0,025) este 1,96.

$Z_{calc} = \frac{172-150}{10} = 2,2$. Deoarece $|Z_{calc}| > 1,96$ respingem H_0 (acceptam H_1) cu pragul de semnificatie de 5%. Adica este putin probabil ca selectia sa provina dintr-o populatie normala cu media $m = 150$ si dispersia $\sigma^2 = 100$.

Aici am folosit un test bilateral (zona de respingere este simetrica fata de origine, adica are "2 cozi") deoarece pentru H_1 , m poate fi tot asa de bine < 150 sau > 150 .

Exemplul 11.2 Testati cu un prag de semnificatie $\varepsilon = 1\%$ daca selectia de volum 1, $x_1 = 54$, a fost facuta dintr-o populatie normala cu media $m = 65$ si dispersia $\sigma^2 = 30$, sau daca media este mai mica decat 65.

Solutie $H_0: m = 65$; $H_1: m < 65$. Vom avea deci un test unilateral (la stanga, cu o singura "coada"). $Z = \frac{\bar{X}-m}{\sigma/\sqrt{n}} = \frac{X-m}{\sigma}$.



Zona de respingere pentru un test unilateral

Deoarece in tabelele statistice se dau numai valorile functiei de repartitie normala, $F(z)$, pentru $z \geq 0$, va trebui sa folosim proprietatile de simetrie ale densitatii de probabilitate $\rho(z)$. Avem ca $P(Z < Z_{0,01}) = F(Z_{0,01}) = 0,01$. Deci $F(-Z_{0,01}) = P(Z < -Z_{0,01}) = 1 - P(Z < Z_{0,01}) = 1 - 0,01 = 0,99$, deci, din TABELUL I gasim ca $-Z_{0,01} = 2,33$, adica $Z_{0,01} = -2,33$.

Calculam $Z_{calc} = \frac{54-65}{\sqrt{30}} = -2,01$. Cum $-2,01 > -2,33$ vom accepta ipoteza H_0 cu pragul de semnificatie de 1%. Cum pragul este mic si Z_{calc} este foarte "aproape" de valoarea critica $-2,33$ statisticianul are dubii serioase asupra rezultatului si va trebui sa considere si alta selectie si sa foloseasca un test cu semnificatie mai mare, de exemplu de 5% pentru a fi "mai sigur" de concluzia pe care o da.

Exemplul 11.3 Din 100 de seminte plantate 83 au germinat. Folositi aproximarea distributiei binomiale cu o distributie normala pentru a testa pretentia comerciantului ca 90% din

seminte germineaza. Folositi doua teste: unul cu pragul de semnificatie de 5%, altul cu pragul de 1%.

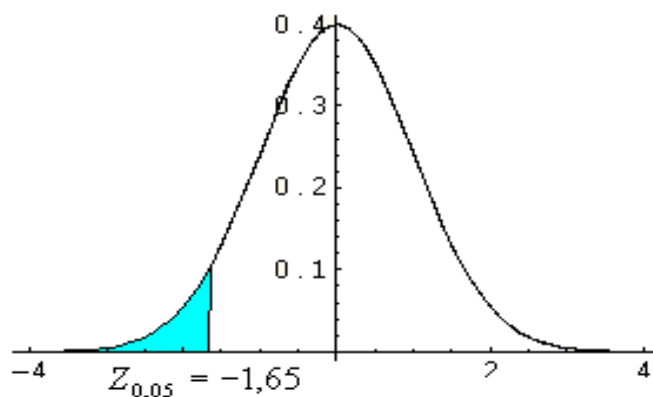
Solutie Fie X v.a. care numara cate seminte au germinat din cele n . $X \sim \text{Bin}(n, p)$, unde $n = 100$.

$H_0: p = 0,9$ (rata de germinare este de 90%).

$H_1: p < 0,9$ (rata de germinare este mai mica de 90%).

Vom avea deci un test unilateral, deoarece este putin probabil ca vanzatorul sa sustina o rata de germinare mai mica decat aceea reala.

Pentru pragul $\varepsilon = 0,05$ avem:



$P(Z < Z_{0,05}) = 0,05$; deci $P(Z < -Z_{0,05}) = 0,95 = F(-Z_{0,05})$.

Din TABELUL I gasim ca $-Z_{0,05} = 1,65$, deci $Z_{0,05} = -1,65$.

Cum $H_0: X \sim \text{Bin}(100; 0,9)$ avem ca $X \sim N(np, npq) = N(90, 9)$, deci $Z_{\text{calc}} = \frac{X - np}{\sqrt{npq}} = \frac{83 - 90}{3} = -2,33$.

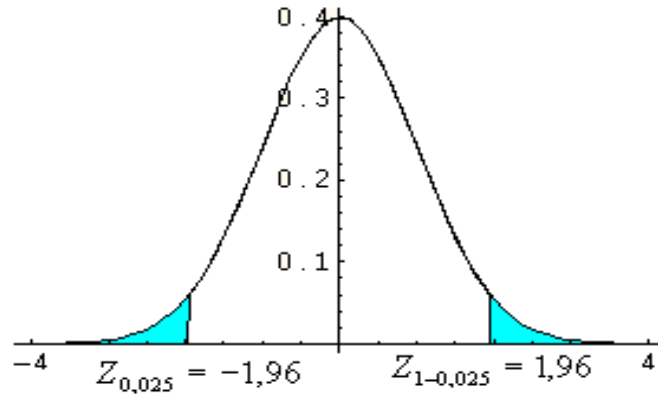
Dar $Z_{\text{calc}} = -2,33 < -1,65$ si deci va trebui sa resping H_0 cu pragul de 5%. Adica fabricantul de seminte... minte!

Pentru pragul $\varepsilon = 0,01$, $\phi(-Z_{0,01}) = 0,99$, deci $-Z_{0,01} = 2,32$, sau $Z_{0,01} = -2,32$. Cum $Z_{\text{calc}} = -2,33 < -2,32$, dar aproape insensibil mai mic, testul in acst caz **nu** poate fi concludent deoarece valoarea calculata Z_{calc} este prea aproape de valoarea critica $-2,32$. Prin urmare, fabricantul... minte, dar nu minte prea mult! Este nevoie de alte solutii pentru a capata o certitudine mai mare.

Exemplul 11.4 O masina produce benzi elastice cu tensiuni de rupere normal distribuite cu media $m = 45N$ si $\sigma = 4,36N$. Intr-o zi s-a facut o selectie de volum 50 si s-a gasit media selectiei $\bar{x} = 43,46N$. Testati cu un prag de semnificatie de 5% daca acest lucru indica sau nu o schimbare a mediei tensiunilor de rupere.

Solutie $H_0: m = 45$ (media nu s-a schimbat)

$H_1: m \neq 45$ (media s-a schimbat)-test bilateral!



$$\bar{X} \sim N\left(m, \frac{\sigma^2}{n}\right), m = 45; \sigma = 4,36; n = 50.$$

$Z_{\text{calc}} = \frac{\bar{x} - m}{\sigma/\sqrt{n}} = \frac{43,46 - 45}{43,36/\sqrt{50}} = -2,4975 < -1,96$. Prin urmare respingem ipoteza H_0 cu pragul de semnificatie de 5%.

Un interval de incredere de nivel 95% pentru medie este $\bar{x} \pm 1,96 \frac{\sigma}{\sqrt{n}} = (42,25; 44,67)$. Vedem ca $45 \notin (42,25; 44,67)$. Cea mai mica valoare a lui σ astfel incat 45 sa fie in intervalul de incredere $\bar{x} \pm 1,96 \frac{\sigma}{\sqrt{n}}$ este $\sigma = 5,56$ (vezi ecuatia $43,46 + 1,96 \frac{\sigma}{\sqrt{50}} = 45$).

Exemplul 11.5 Tensiunea de rupere a unor cabluri produse de o fabrica este normal distribuita cu media 6000N si deviatia standard $\sigma = 150$ N. Gasiti probabilitatea ca un cablu luat la intamplare se aiba tensiunea de rupere > 6200 N.

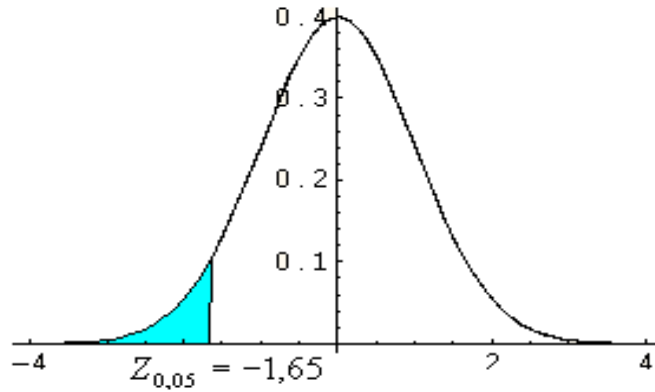
S-a modificat procesul de productie si media tensiunilor de rupere se modifica. Se aleg 6 cabluri la intamplare dupa aceasta modificare, se testeaza si se gaseste o medie de rupere $\bar{x} = 5920$ N. Testati cu un prag de 5% daca dupa modificare media tensiunilor s-a micorat. Gasiti o constanta C a.i. noi sa putem spune cu un nivel de incredere de 90% ca media de rupere este mai mare decat C .

Solutie $X \sim N(6000, 150^2)$;

$$\begin{aligned} P(X > 6200) &= P\left(\frac{X - m}{\sigma} > \frac{6200 - m}{\sigma}\right) = P(Z > 1,333) = 1 - P(Z \leq 1,333) \\ &= 1 - F(1,333) = 1 - 0,90 = 0,1. \end{aligned}$$

$$\bar{x} = 5920\text{N}; H_0: m = 6000\text{N}; H_1: m < 6000\text{N};$$

$$\bar{X} \sim N\left(m, \frac{\sigma^2}{n}\right) = N\left(6000, \frac{150^2}{6}\right);$$



$Z_{calc} = \frac{\bar{x}-m}{\sigma/\sqrt{n}} = \frac{5920-6000}{150/\sqrt{6}} = -1,306 > -1,65$, deci acceptam ipoteza H_0 cu pragul de semnificatie 5%.

• Trebuie sa gasim C a.i. $P(C < m < \infty) = 0,9$, sau $P(-C > -m) = 0,9$, sau inca $P\left(\frac{\bar{X}-C}{\sigma/\sqrt{6}} > \frac{\bar{X}-m}{\sigma/\sqrt{6}} = Z\right) = 0,9$. Deci $F\left(\frac{\bar{X}-C}{\sigma/\sqrt{6}} > Z\right) = 0,9$ si de aici gasim ca $\frac{\bar{X}-C}{\sigma/\sqrt{6}} = Z_{0,9} = 1,29$. Prin urmare $C = \bar{x} - 1,29 \frac{150}{\sqrt{6}}$.

Exemplul 11.6 O distributie normala se crede a avea media 50. Se face o selectie de volum 100 si se gaseste o medie de 52,6 si o deviatie standard de selectie de 14,5. Testati cu nivelul de 5% daca media populatiei a crescut.

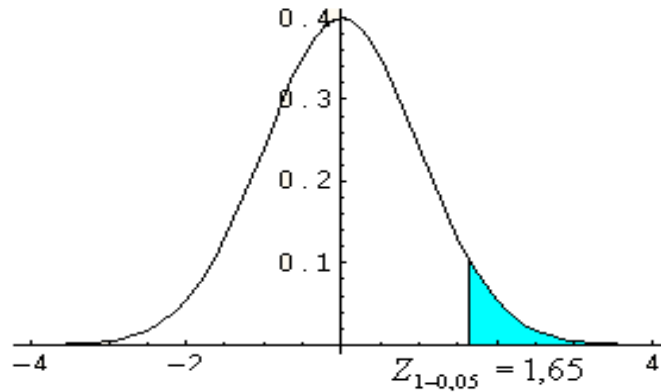
Solutie Fie m media reala si σ^2 dispersia reala a populatiei. $\bar{x} = 52,6$, iar $s = \sqrt{s^2} = 14,5$,

unde $s^2 = \frac{1}{n} \sum (x - \bar{x})^2$.

$H_0: m = 50$ (nu exista o schimbare a mediei)

$H_1: m > 50$ (media populatiei a crescut)

$\bar{X} \sim N\left(m, \frac{\sigma^2}{n}\right)$. Estimam pe σ^2 cu $\hat{\sigma}^2 = \frac{ns^2}{n-1} \approx s^2$, deoarece $n = 100$ este considerat mare ($100 > 30$). Folosim deci statistica $Z = \frac{\bar{X}-m}{\hat{\sigma}/\sqrt{n}}$ si calculam $Z_{calc} = \frac{52,6-50}{14,5/\sqrt{100}} = 1,793$.



Cum $Z_{calc} = 1,793 > 1,645$, vom respinge H_0 , adica acceptam H_1 cu pragul de semnificatie de 5%. Deci acceptam ca mesia a crescut cu pragul de semnificatie de 5%.

11.1.2 Testul T privind media unei populatii normale cu dispersia estimata prin estimatorul nedeplasat S'^2

Ipoteza nula $H_0: m = m_0$ se refera la o populatie normala careia nu ii cunoastem dispersia. Aceasta se va estima fie prin estimatorul deplasat $S^2 = \frac{1}{n} \sum (x - \bar{x})^2$, fie prin estimatorul nedeplasat $S'^2 = \frac{1}{n-1} \sum (x - \bar{x})^2$. Daca volumul selectiei n este mare ($n \geq 30$) atunci $S'^2 \approx S^2$ si putem folosi ca statistica pentru un test de semnificatie, statistica Z (vezi Exemplul 11.1.6). Daca insa volumul selectiei este mic *nu* mai putem folosi aceasta statistica ci este indicat sa folosim statistica Student cu $n - 1$ grade de libertate, $T = \frac{\bar{X} - m_0}{S'/\sqrt{n}}$, unde $S' = \sqrt{S'^2}$.

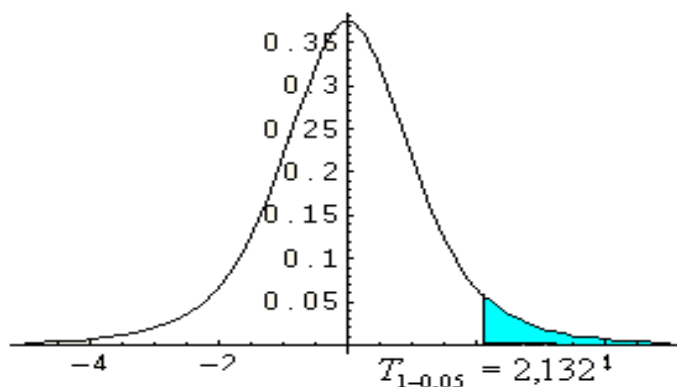
In rest, testul lucreaza exact ca testele de la 11.1.1.

Exemplul 11.7 Se testeaza rezistenta in ohmi pentru 5 bucati de cablu si se gasesc valorile: 1,51; 1,49; 1,54; 1,52; 1,54. Daca cablul ar fi din argint pur, rezistenta lui ar fi de 1,5 ohmi. Daca argintul nu este pur, rezistenta creste. Testati cu un nivel de semnificatie de 5% faptul ca argintul din cablu nu este pur.

Solutie $H_0: m = 1,5$ ohmi.

$H_1: m > 1,5$ ohmi (test unilateral)

Deoarece esantionul selectiei este mic ($n = 5$) vom folosi statistica student $T = \frac{\bar{X} - m}{S'/\sqrt{n}} = \frac{\bar{X} - m}{S/\sqrt{n-1}}$, cu $n - 1 = 4$ grade de libertate.



$F(T_{0.05}) = P(T < T_{0.05}) = 0.95$. Din TABELUL II, pentru $\nu = 4$ grade de libertate gasim ca $T_{0.05} = 2.132$.

$T_{calc} = \frac{\bar{x} - m}{S/\sqrt{n-1}} = \frac{1.52 - 1.50}{0.019/2} = 2.105$, deoarece $\bar{x} = \frac{7.6}{5} = 1.52$, iar $S^2 = \frac{1}{n} \sum (x - \bar{x})^2 = \frac{0.0018}{5} = 0.00036$, sau $S = 0.019$.

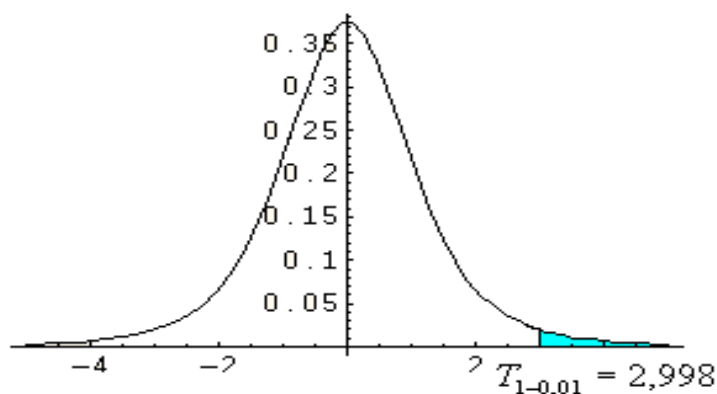
Deoarece $T_{calc} < 2.132$, valoarea critica a testului, atunci trebuie sa acceptam ipoteza H_0 cu nivelul de semnificatie 5%.

Exemplul 11.8 Se fac 8 observatii dintr-o populatie normala si gasim $\bar{x} = 4.65$ si $\sum (x - \bar{x})^2 = 0.74$. Testati cu nivelul de semnificatie de 2% daca media distributiei este 4.3.

Solutie $H_0: m = 4.3$

$H_1: m \neq 4.3$ (test bilateral)

$T = \frac{\bar{X} - m}{S/\sqrt{n-1}}; T_{calc} = \frac{4.65 - 4.3}{\frac{0.74}{8\sqrt{7}}} = 3.05$.



Cum $T_{calc} > 2,998$, resping H_0 cu nivelul de semnificatie de 2%.

11.1.3 Test pentru proportia de succese

Sa notam cu $P_s = (\text{numarul de realizari ale evenimentului A din } n \text{ incercari})/n$, adica proportia de realizari a evenimentului A, cu X v.a. $\begin{pmatrix} 1 & 0 \\ p & q \end{pmatrix}$, unde $q = 1 - p$, iar p este probabilitatea teoretica pentru a se realiza evenimentul A la o incercare. Daca punem $X_1 = X_2 = \dots = X_n = X$, este clar ca P_s este chiar $\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$. Pentru n mare, stim ca $\bar{X} \sim N(M(P_s), D(P_s))$, dar $M(P_s) = M(\bar{X}) = p$, $D(P_s) = D(\bar{X}) = \frac{D(X)}{n} = \frac{p-p^2}{n} = \frac{pq}{n}$. Deci

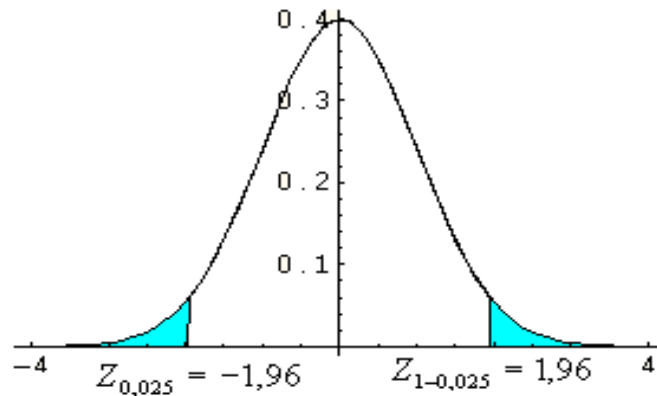
$P_s = \text{"proportia de succese"} \sim N(p, \frac{pq}{n})$. Fie $Z = \frac{P_s - p}{\sqrt{\frac{pq}{n}}} \sim N(0, 1)$. Testul proportiei de succese se realizeaza cu ajutorul statisticii Z dupa cum se va vedea in exemplul urmator.

Exemplul 11.9 *La o universitate americana senatul sustine ca nu se face discriminare sexuala la admitere. Se aleg 500 studenti si se gasesc 267 baieti. Testati cu nivelul de semnificatie 5% daca senatul universitatii spune adevarul sau nu.*

Solutie $H_0: p = 0,5 = \text{probabilitatea ca un student sa fie baiat};$

$H_1: p \neq 0,5$ (test bilateral)

$$Z = \frac{P_s - p}{\sqrt{\frac{pq}{n}}}; z_{calc} = \frac{\frac{267}{500} - 0,5}{\sqrt{\frac{0,5 \cdot 0,5}{500}}} = 1,52.$$



Cum $z_{calc} = 1,52 < 1,96$, acceptam ipoteza H_0 cu pragul de semnificatie de 5%. Prin urmare este foarte probabil ca senatul sa spuna adevarul.

Exemplul 11.10 ([Hays], pag. 447) Un producător de frigidere pretinde că temperatura medie în compartimentul de congelare este de 10 grade Fahrenheit (aproximativ -12,3 grade Celsius). Vrem să vedem adevărul acestei afirmații și facem ipotezele: $H_0 : m \leq 10$ (ipoteza nulă), versus $H_1 : m = 10$ (ipoteza alternativă). Facem un sondaj pe un lot de 16 frigidere alese la întâmplare și măsurăm temperaturile în congelatoarele acestora. Găsim că media selecției este 10,24 grade, iar dispersia de selecție modificată (nedeplasată) este de $S'^2 = 0,36$. Presupunem că distribuția temperaturilor este normală (deoarece apare un fenomen de repartiție de "erori"). Cum dispersia este calculată tot din selecție folosim v.a. de selecție T . Calculăm $t = \frac{10,24-10}{0,6/4} = 1,6$. Vrem să folosim statistica T cu pragul de semnificație de 5% pentru a vedea dacă producătorul are dreptate sau nu. Testul va fi unilateral deoarece zona de respingere este data de $m > 10$. $T_{0,95} = 1,753$ pentru 15 grade de libertate, după cum se poate constata în TABELUL II. $T_{calc} = 1,6 < 1,753$. Prin urmare acceptăm ipoteza H_0 cu pragul de 5%.

Observația 11.11 a) În exemplul de mai sus regiunea de respingere este: $m > 10$, concentrată într-o singură direcție, adică la dreapta lui 10. Un astfel de test se zice *directional* (sau *unilateral*).

b) Dacă vrem să testăm o ipoteză despre dispersia unei populații: $H_0 : \sigma^2 = \sigma_0^2$, unde σ_0 este dată, va trebui să utilizăm testul χ^2 (vezi Lectia 12). Aici $H_1 : \sigma^2 \neq \sigma_0^2$, și α este dat. Cele două (α și H_1) descriu regiunea de respingere. Statistica folosită este $\chi_{(n-1)}^2 = \frac{(n-1)S'^2}{\sigma_0^2}$, unde $(n-1)$ reprezintă numărul gradelor de libertate. Se urmează apoi aceleași cale ca și în cazul mediei (vezi Rezumatul și exemplele date acolo).

11.1.4 Testul T pentru compararea a două esantioane

Fie $\{x_1, \dots, x_n\}$ și $\{x'_1, \dots, x'_m\}$ două esantioane. $H_0 : \ll \text{cele două esantioane provin din aceeași populație} \gg$.

Fie $\alpha \in (0,1)$ un prag de semnificație. Notăm cu $T_\nu = \frac{\bar{X} - \bar{X}'}{S'/\sqrt{\frac{nm}{n+m}}}$, unde $\bar{X} = (\sum x_i)/n$, $\bar{X}' = (\sum x'_i)/m$, $S'^2 = \frac{1}{\nu} \left(\sum (x_i - \bar{X})^2 + \sum (x'_i - \bar{X}')^2 \right)$, iar $\nu = n + m - 2$. R. A. Fisher a arătat că T_ν tinde către o repartiție Student cu ν grade de libertate. Calculăm T_ν . Găsim cuantila de ordin $1-\alpha/2$ în Tabelul corespunzător lui T pentru ν grade de libertate și o notăm cu $T_{\alpha/2}$. Testul funcționează astfel:

- Dacă $|T_\nu| > T_{\alpha/2}$ atunci vom respinge ipoteza H_0 cu pragul de semnificație α .
- Dacă $|T_\nu| \leq T_{\alpha/2}$ atunci vom accepta ipoteza H_0 cu pragul de semnificație α .

Exemplul 11.12 (R. A. Fisher) 8 ghivece cu fire de orez au fost supuse la descărcări electrice. Altele 9 au fost ferite de descărcări. Rezultatul recoltei a fost (număr de spice):

Izolate: 17, 27, 18, 25, 27, 29, 27, 23, 17;

Electrizate: 16, 16, 20, 16, 20, 17, 15, 21.

Să se testeze ipoteza $H_0 : \ll \text{electricitatea influențează creșterea orezului} \gg$.

11.2 Tipuri de erori. Reguli de decizie

Vom începe cu un exemplu.

Exemplul 11.13 ([Hays], pag. 404) *Un economist are două ipoteze asupra implicațiilor ce derivă din creșterea impozitelor la un anumit moment dat. Prima ipoteză este că după această creștere 80% din populație va trebui să-și reducă economiile. A doua ipoteză este că numai 40% din populație va trebui să facă acest lucru. Cum s-ar putea afla care ipoteză este adevărată?*

Soluție. S-ar putea ca nici una dintre cele două ipoteze să nu fie adevărată. Totuși, aici, noi vom considera că *sigur* una dintre ele este adevărată. Vom nota:

$H_0: p=0,8$

$H_1: p=0,4$, unde p este proporția de consumatori care urmează să-și reducă economiile datorită creșterii impozitelor. Iată că impozitele au crescut și economistul nostru vrea să vadă care dintre cele două ipoteze ale sale (fiecare are în spatele ei rationamente și teorii economice sofisticate) este adevărată. Pentru aceasta face un sondaj pe un esantion de n consumatori. Deoarece fiecare din consumatori spune DA sau NU (și-a redus sau nu și-a redus economiile) avem un proces de tip Bernoulli cu n dat și p dat. Pentru H_0 trebuie să considerăm $p=0,8$, iar pentru H_1 trebuie să considerăm $p=0,4$. Presupunem, pentru ușurință, că $n=10$. Notăm cu r numărul acelor consumatori care au răspuns cu DA (dintre cei n chestionați). Statistica de selecție care poate fi comparată cu p este R/n , unde R este v.a. ce poate lua valorile $r: 0, 1, \dots, n$. Valorile teoretice pe care le poate lua R/n și probabilitățile lor le găsim în următorul tabel (am folosit distribuția binomială, vezi Tabelul V):

r	r/n	$P(r/n \mid p = 0,8)$ aprox. 4 zecimale	$P(r/n \mid p = 0,4)$ aprox. 4 zecimale
0	0	0	0,006
1	0,1	0	0,040
2	0,2	0	0,121
3	0,3	0,001	0,215
4	0,4	0,006	0,251
5	0,5	0,026	0,2
6	0,6	0,088	0,111
7	0,7	0,201	0,042
8	0,8	0,302	0,011
9	0,9	0,268	0,002
10	1	0,107	0,0001

(11.1)

Să presupunem că statisticianul are datele de selecție ale sondajului, are deci raportul R/n calculat din sondaj. El are nevoie de o *REGULĂ DE DECIZIE* pentru a putea alege H_0 sau H_1 . Există foarte multe posibilități de a construi astfel de reguli de decizie. Unele

sunt mai "bune", altele nu sunt asa de "bune". Vom alege acum următoarea regulă: <<Dacă $R/n < 0,8$, alegem H_1 ; dacă $R/n \geq 0,8$, alegem H_0 >>(REGULA 1).

Ce se poate întâmpla după ce statisticianul a folosit această regulă?

El poate gresi sau nu. Să calculăm probabilitățile în toate cele patru cazuri care pot să apară:

		<i>Situatia reală</i>	
		H_0	H_1
<i>Decizia luată</i>	H_0	corect	eroare II
	H_1	eroare I	corect

(11.2)

De exemplu, să presupunem că din selectie am obtinut $R/n \geq 0,8$ și totusi în realitate $p=0,4$. Deci am ales H_0 și totusi H_1 este adevărată. Apare al doilea tip de eroare (eroare II). Să calculăm probabilitatea acestei erori (folosind tabelul de mai sus):

$$\begin{aligned}
 & P(R/n \geq 0,8 \mid p=0,4) \\
 = & P(R/n=0,8 \mid p=0,4) + P(R/n=0,9 \mid p=0,4) + P(R/n=1 \mid p=0,4) \\
 = & 0,011 + 0,002 + 0,0001 \approx 0,013.
 \end{aligned}$$

Primul tip de eroare (eroare I) apare dacă alegem H_1 și totusi H_0 este adevărată. Probabilitatea acesteia este:

$$\begin{aligned}
 & P(R/n < 0,8 \mid p=0,8) \\
 = & +0 + 0 + 0,001 + 0,006 + 0,026 + 0,088 + 0,201 = 0,322.
 \end{aligned}$$

Cele două situatii corecte au următoarele probabilități:

$$P(R/n \geq 0,8 \mid p=0,8) = 0,677,$$

$$P(R/n < 0,8 \mid p=0,4) = 0,987.$$

Punem aceste rezultate în următorul tabel:

		<i>Situatia reală</i>	
		H_0	H_1
<i>Decizia luată</i>	H_0	0,677	0,013
	H_1	0,323	0,987

(11.3)

Ce spune acest tabel? Dacă după selectie $R/n < 0,8$, este foarte probabil ca $p=0,4$. Oricum este mai probabil acest lucru decât faptul că $R/n \geq 0,8$ și $p=0,8$. Este foarte puțin probabil "să gresim" cu această regulă de decizie, deoarece $0,323+0,013 < 0,677+0,987$.

Iată o altă regulă de decizie: <<Dacă $R/n < 0,6$, alegem H_1 ; dacă $R/n \geq 0,6$, alegem H_0 >>(REGULA 2). Tabelul corespunzător acestei reguli este următorul:

		<i>Situația reală</i>		
		H_0	H_1	
<i>Decizia luată</i>	H_0	0,966	0,116	(11.4)
	H_1	0,034	0,834	

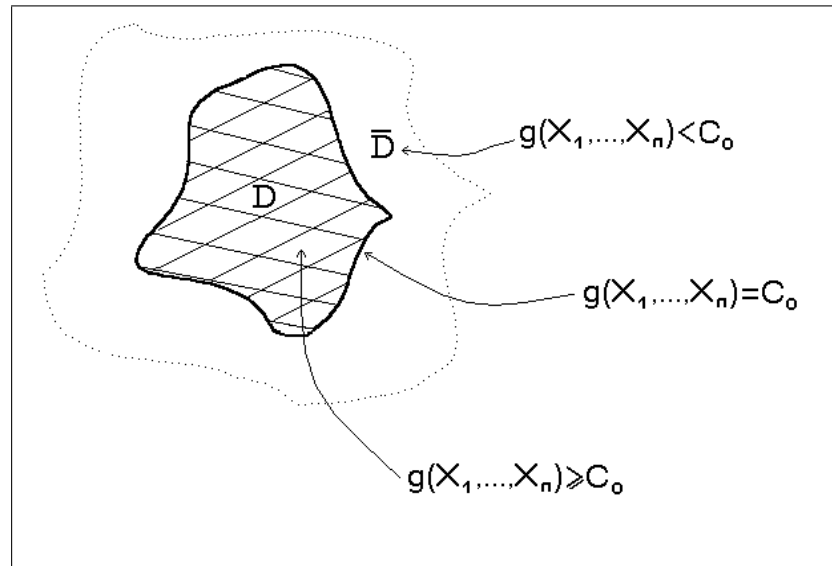
Este clar că dintre cele două reguli de decizie este "mai bună" a doua regulă deoarece probabilitățile de eroare sunt mici.

În prima regulă de decizie valoarea lui $R/n=0,8$ (care face trecerea de la zona ipotezei H_0 la zona ipotezei H_1) se numește *valoarea critică a lui R/n* (adică a statisticii R/n).

Pentru a doua regulă de decizie valoarea critică a statisticii R/n este 0,6.

• În general, frontiera dintre domeniul de acceptare D și domeniul de respingere \bar{D} , formează mulțimea de puncte în care statistica de testare $g(X_1, \dots, X_n)$ ia valoarea critică C_0 .

În figura următoare avem reprezentarea grafică a *domeniului de acceptare* D pentru ipoteza H_0 , a *domeniului de respingere* \bar{D} pentru ipoteza H_0 și a mulțimii punctelor de valoare critică pentru următoarea regulă de decizie: <<Dacă $g(X_1, \dots, X_n) \geq C_0$, acceptă H_0 ; dacă $g(X_1, \dots, X_n) < C_0$, acceptă H_1 >>. Aici C_0 este valoarea critică a testului.



Regiunea de acceptare și regiunea de respingere a unui test

Desigur că statistica $g(X_1, \dots, X_n)$ este aleasă astfel încât să fie o legătură naturală între ea și ipotezele H_0 și H_1 . Se consideră, de asemenea, că H_0 și H_1 se exclud reciproc ($D \cap \bar{D} = \emptyset$).

Vom presupune în continuare că H_0 este *ipoteza care se testează*. Statisticianul poate face două tipuri de erori:

- Eroare de tipul I dacă respinge H_0 , ea fiind adevărată;
- Eroare de tipul II dacă acceptă H_0 , ea nefiind neadevărată.

Notăm cu $\alpha = P(\text{eroare de tipul I}) = P(\text{respinge } H_0 \mid H_0 \text{ este adevărată})$;

$\beta = P(\text{eroare de tipul II}) = P(\text{acceptă } H_0 \mid H_0 \text{ este falsă})$.

Tabelele (11.3) și (11.4) din Exemplul 11.1 se generalizează la următorul tabel:

		Situatia reală	
		H_0	H_1
Decizia luată	Accept H_0	$1-\alpha$	β
	Resping H_0	α	$1-\beta$

(11.5)

Orice regulă de decizie are un cuplu de numere (α, β) .

Idealul ar fi ca α și β să fie 0, sau foarte mici.

Dintre două reguli de decizie cu (α_1, β_1) și (α_2, β_2) astfel încât $\alpha_1 \leq \alpha_2$, $\beta_1 \leq \beta_2$, vom elimina pe cea de-a doua. Spunem că regula de decizie cu (α_1, β_1) *domină* (este mai tare) regula de decizie cu (α_2, β_2) . În cazul Exemplului 11.1 cele două reguli nu se pot compara, ele nu se domină una pe alta.

O regulă de decizie dominată de o altă regulă de decizie se zice *inadmisibilă*.

Vom da acum un exemplu de regulă de decizie inadmisibilă: <<Dacă $R/n \in (-\infty; 0, 2) \cup (0, 8; \infty)$, alegem H_1 ; dacă $R/n \in [0, 2; 0, 8]$, alegem H_0 >>(REGULA III).

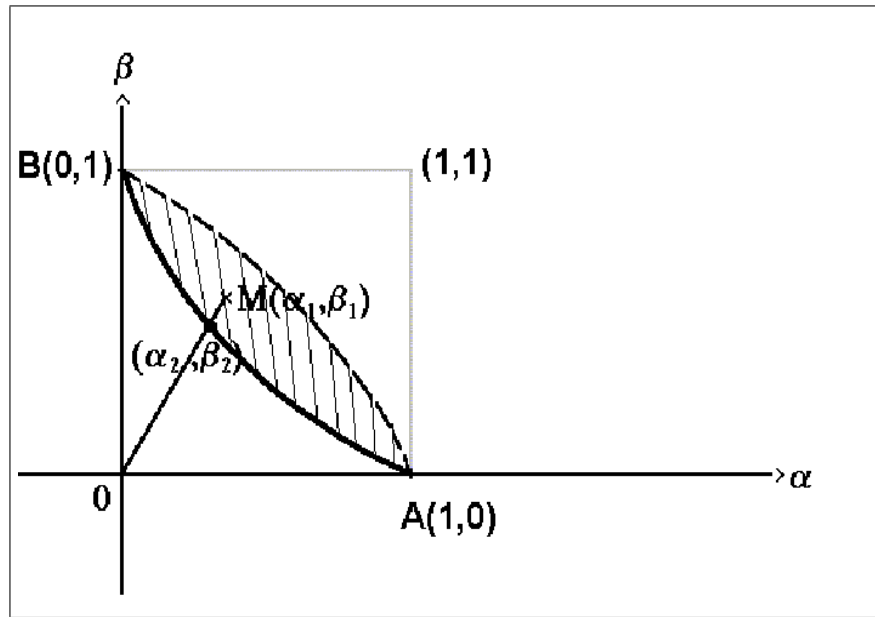
Tabelul corespunzător ei este:

		Situatia reală	
		H_0	H_1
Decizia luată	Accept H_0	0,625	0,952
	Resping H_0	0,375	0,048

(11.6)

Aici $\alpha_2=0,375$, $\beta_2=0,952$. Cum $\alpha_1=0,323$ și $\beta_1=0,013$ erau probabilitatile de eroare pentru prima regulă, rezultă că regula <<Dacă $R/n \in (-\infty; 0, 2) \cup (0, 8; \infty)$, alegem H_1 ; dacă $R/n \in [0, 2; 0, 8]$, alegem H_0 >>(REGULA III) este inadmisibilă deoarece este dominată de prima regulă de decizie. Nu vom lucra niciodată cu reguli de decizie despre care stim că sunt inadmisibile. Dacă se analizează îndeaproape regula <<Dacă $R/n \in (-\infty; 0, 2) \cup (0, 8; \infty)$, alegem H_1 ; dacă $R/n \in [0, 2; 0, 8]$, alegem H_0 >>(REGULA III) se constată că ea contrazice chiar "bunul simț" probabilistic (De ce?).

În general statisticianul este interesat de modul în care variaza probabilitatile de eroare α și β atunci când el schimbă legea de decizie. Evident că pentru orice statistică de testare $g(X_1, \dots, X_n)$ și pentru orice lege de decizie fixată avem un α și un β bine determinați. Păstrăm statistica de testare fixată și variem legile de decizie. Se obține un domeniu *al probabilităților de eroare* (un *domeniu de risc*) dacă α se consideră abscisa și β ordonata punctului (α, β) — vezi figura următoare:



Domeniul probabilităților de eroare

Curba groasă \widehat{AB} din figura anterioară este curba corespunzătoare tuturor cuplurilor (α, β) care derivă din decizii "bune" sau admisibile. Oricare alt punct (α_1, β_1) din domeniul de risc care *nu* se află pe curba \widehat{AB} provine dintr-o decizie inadmisibilă, deoarece aceasta este dominată de decizia corespunzătoare punctului (α_2, β_2) de pe curba \widehat{AB} și aflat la intersecția dintre această curbă și segmentul OM.

Experiența arată că dacă volumul de selecție n crește, atunci curba \widehat{AB} "se apropie" de origine, adică riscul devine mai mic cel puțin pentru deciziile admisibile (deoarece precizia de predicție asupra populației crește odată cu n).

Singura problemă care rămâne pentru statistician este aceea legată de alegerea regulii de decizie admisibile. Aici intervine "negocierea" între cazurile " α mare, β mic", sau invers, " α mic, β mare", $(\alpha, \beta) \in \widehat{AB}$.

Situația neutră este aceea când $\alpha = \beta$. În practică această negociere depinde de factori subiectivi sau obiectivi dar.

- Din punct de vedere istoric ipoteza H_0 se numește *ipoteza nulă* (nevinovăția prezumtivă în cazul unui acuzat!), iar ipoteza H_1 se zice *ipoteză alternativă* (vinovăția acuzatului). În practică, numărul α (notat uneori și cu ε) se dă ca fiind 0,05 sau 0,01 (rareori se utilizează altă valoare). α reprezintă probabilitatea de a respinge H_0 cu toate că H_0 este adevărată.

Presupunem că H_0 este adevărată și alegem așa numita *regiune (domeniu) de respingere* în concordanță cu ipoteza H_1 . Să notăm domeniul de respingere cu $Res = \{\text{acele valori ale statisticii de testare } g(X_1, \dots, X_n) \text{ pentru care } H_1 \text{ este adevărată}\}$. Găsim mulțimea Res din ipoteza:

$$P(g(X_1, \dots, X_n) \in Res \mid H_0 = \text{adevărată}) = \alpha(dat)$$

Dacă la o selecție de volum n : X_1, \dots, X_n , valoarea v.a. $g(X_1, \dots, X_n) \in Res$, spunem că ipoteza H_1 este adevărată cu pragul de semnificație α (cu condiția ca α să fie mic, ca mai sus).

Vom da acum câteva exemple de testare a unor ipoteze statistice în lumina celor spuse mai sus.

Exemplul 11.14 ([Hays], pag. 415) *Un muncitor poate să realizeze X piese pe oră. După îndelungate cercetări statistice s-a stabilit că X este o v.a. normală cu media $m=138$ și deviația standard $\sigma=20$. Un inginer pretinde că poate aduce o inovație în procesul de producție astfel încât să ridice media la $m=142$ piese pe oră, fără a perturba normalitatea v.a. X și pe σ . Este chemat un statistician să testeze pretentia inginerului.*

Soluție Introducem două ipoteze:

$$H_0: m=138$$

$$H_1: m=142$$

Facem o selecție de 100 muncitori care lucrează câte o oră fiecare. Alegem pragul de semnificație $\alpha=0,05=P(\text{respingem } H_0 \mid H_0 \text{ este adevărată})$. Vrem acum să găsim regiunea de respingere în acest caz concret. Stim că un estimator bun pentru media m este media de selecție $\bar{X} = (X_1 + X_2 + \dots + X_n)/n$, care este normal distribuită cu media m și deviația standard (pentru selecția noastră $n=100$) $\sigma_{\bar{X}} = 20/\sqrt{100} = 2$. Chiar dacă ipoteza de normalitate a v.a. X nu este întrutotul adevărată (sau chiar falsă!), deoarece $n=100$ este mare, din teorema limită centrală, rezultă că \bar{X} este distribuită normal. Dacă H_0 este adevărată, atunci \bar{X} este normal distribuită cu media 138 și deviația standard 2, pe când, dacă H_1 este adevărată, $m=142$ și deviația standard tot egală cu 2. Vedem de aici că valorile mari ale v.a. \bar{X} favorizează ipoteza H_1 , iar valorile mici ale v.a. \bar{X} favorizează ipoteza H_0 .

Regiunea de respingere va fi deci de forma: $\langle\langle \text{Respinge } H_0 \text{ dacă } \bar{X} \geq C \rangle\rangle$. Aici C reprezintă valoarea critică a regulii de decizie: $\langle\langle \text{Dacă } \bar{X} < C, \text{ alegem } H_0; \text{ dacă } \bar{X} \geq C, \text{ alegem } H_1 \rangle\rangle$.

Vrem acum să determinăm valoarea critică C a statisticii \bar{X} dacă stim pe $\alpha=0,05$. Scriem că $\alpha = P(\text{respinge } H_0 \mid H_0 = \text{adevărată}) = P(\bar{X} \geq C \mid m=138) = 0,05$. Dar $P(\bar{X} \geq C \mid m=138) = P(Z \geq \frac{C-138}{2})$, unde $Z = \frac{\bar{X}-138}{2}$ este v.a. standard normală cu media 0 și cu deviația standard 1. Vrem ca $P(Z \geq \frac{C-138}{2}) = 0,05 = 1 - F(\frac{C-138}{2})$, unde $F(z)$ este funcția de repartiție (funcția de distribuție cumulată) a v.a. normale standard Z . Din Tabelul I găsim că $\frac{C-138}{2}$ este 1,65, adică cuantila de ordin $1-0,05=0,95$ a v.a. Z . De aici rezultă că punctul critic $C=141,30$. Acum putem calcula și probabilitatea

$$\begin{aligned} \beta &= P(\text{accept } H_0 \mid H_1 = \text{adevărată}) = P(\bar{X} < 141,30 \mid m = 142) \\ &= P\left(Z < \frac{141,30 - 142}{2}\right) = P(Z < -0,35) = F(-0,35) = 1 - F(0,35) = 0,36. \end{aligned}$$

Tabelul corespunzător regulii de decizie de mai sus este:

		<i>Situația reală</i>	
		H_0	H_1
<i>Decizia luată</i>	Accept H_0	0,95	0,36
	Resping H_0	$\alpha=0,05$	0,64

(11.7)

Dacă din calcule $\bar{X} < 141,30$, atunci alegem H_0 , adică rămânem la vechiul procedeu de lucru.

Să comentăm puțin rezultatele din tabelul (11.7). $\alpha=0,05$, chiar dacă $\bar{X} \geq 141,30$, adică aleg H_1 și resping H_0 , probabilitatea de eroare în cazul când H_0 este adevărată, este foarte mică. Prin urmare, dacă alegem noul procedeu și-l înlocuiesc cu primul (primul fiind mai bun) riscul este mai mic, aproape zero. Totuși, dacă alegem H_0 ($\bar{X} < 141,30$)—rămânem la vechiul procedeu în timp ce noul procedeu este mai bun, riscul este mai mare: $\beta=0,36$.

Statisticianul vrea să micsoreze și acest risc β fără însă a mări pe $\alpha=0,05$ (putem micsora pe β dacă reducem valoarea critică la 140, de exemplu; dar, în acest caz crește și α). Teoretic stim că α și β se micsorează dacă mărim volumul selecției.

Vom mări pe n la 400. În acest caz deviația standard a v.a. \bar{X} devine $\sigma_{\bar{X}} = 20/\sqrt{400} = 1$. Determinăm valoarea critică ca mai sus și găsim $C = 139,65$. În acest caz

$$\begin{aligned}\beta &= P(\text{acceptă } H_0 \mid H_1 = \text{adevărată}) = P(\bar{X} < 139,65 \mid m = 142) \\ &= P\left(Z < \frac{139,65 - 142}{1}\right) = P(Z < -2,35) = 0,01.\end{aligned}$$

Iată deci că am redus pe β de la 0,36 (când $n=100$) la 0,01 (când $n=400$).

Acum statisticianul poate să răspundă cu riscuri foarte mici (α, β mici) dacă este bine să schimbăm procedeul de producere al pieselor (accept H_1 și resping H_0) sau, dimpotrivă, să lucrăm după același procedeu (accept H_0 și resping H_1).

În primul caz pot greși cu 5%, iar în al doilea caz cu 1%. Facem deci un sondaj de 400 muncitori/oră. Calculăm \bar{X} și vedem dacă $\bar{X} < 139,65$ sau $\bar{X} \geq 139,65$.

Observația 11.15 Pentru valoarea critică $C=141,30$, regiunea de respingere este $[141,30;\infty)$, pe când pentru valoarea critică $C=139,65$, regiunea de respingere este $[139,65;\infty)$.

11.3 Puterea unui test statistic

Fie θ un parametru al unei populații și ipoteza: $H_0 : \theta = \theta_0$

Pentru o regiune de respingere dată și pentru o valoare particulară a lui θ , de exemplu θ_1 , putem calcula $P(\text{resping } H_0 \mid \theta = \theta_1)$. Dar, dacă $\theta_1 = \theta_0$, această ultimă probabilitate este chiar α . Dacă $\theta_1 \neq \theta_0$ notăm cu $H_1 : \theta = \theta_1$ și probabilitatea de mai sus devine $P(\text{resping } H_0 \mid H_1 = \text{adevărată}) = 1 - P(\text{acceptă } H_0 \mid H_1 = \text{adevărată}) = 1 - \beta$.

Definiția 11.16 Numărul $P(\text{resping } H_0 \mid \theta = \theta_1)$, calculat mai sus, se numeste puterea testului ipotezei H_0 contra alternativei H_1 .

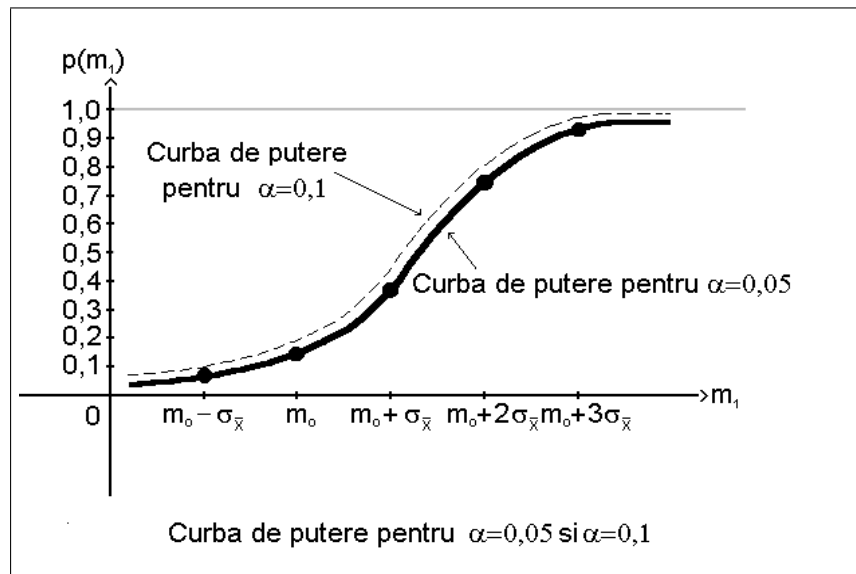
În Exemplul 11.2 puterea testului $H_0 : m = 138$ contra alternativei $H_1 : m = 142$ este egală cu $1-\beta=1-0,36=0,64$, dacă regiunea de respingere este $[141,30;\infty)$. Un test este cu atât mai puternic cu cât ipoteza $H_0 : \theta = \theta_0$ este "mai departe" decât valoarea reală a lui θ , $H_1 : \theta = \theta_1$. Mai mult, $1-\beta$ mare înseamnă β mic, adică un test este cu atât mai puternic cu cât este mai puțin probabil să accepte H_0 când ea de fapt nu este adevărată. Să amintim că noi alegem α foarte mic. Acest α determină zona de respingere (regiunea de respingere). Pentru această zonă de respingere determinăm pe β dacă stim pe θ_1 . Prin urmare puterea unui test $1-\beta$ depinde de θ_1 . Ea va fi maximă pentru acel θ_1 pentru care riscul de a accepta pe H_0 când ea este falsă, comparativ cu H_1 , este minim (aici H_1 este considerată adevărată!).

Exemplul 11.17 Să reluăm exemplul anterior într-un cadru mai general. Presupunem că vrem să studiem media unei populații și facem ipoteza nulă (initială, sau de bază): $H_0 : m = m_0$, și ipoteza alternativă: $H_1 : m = m_1$. Presupunem că $m_1 > m_0$ și că v.a. media de selecție $\bar{X} = (X_1 + X_2 + \dots + X_n) / n$ este normală (pentru n mare \bar{X} poate fi considerată normală) cu deviatia standard $\sigma_{\bar{X}}$ cunoscută. Deoarece valorile mari ale v.a. \bar{X} favorizează H_1 și valorile mici ale v.a. \bar{X} favorizează H_0 , este natural să căutăm zona de respingere de forma $\bar{X} \geq C$, unde C este valoarea critică a regulii de decizie: << respinge H_0 , dacă $\bar{X} \geq C$; acceptă H_0 , dacă $\bar{X} < C$ >>. Ca și în exemplul 11.2 găsim că pentru $\alpha=0,05$ valoarea critică $C = m_0 + 1,65\sigma_{\bar{X}}$. Acum putem calcula puterea acestui test pentru orice valoare dată lui m_1 . De exemplu, dacă $m_1 = m_0 + \sigma_{\bar{X}}$ puterea testului este $P(\text{respinge } H_0 \mid m = m_0 + \sigma_{\bar{X}}) = P(\bar{X} \geq m_0 + 1,65\sigma_{\bar{X}} \mid m = m_0 + \sigma_{\bar{X}}) = P\left(Z \geq \frac{(m_0 + 1,65\sigma_{\bar{X}}) - (m_0 + \sigma_{\bar{X}})}{\sigma_{\bar{X}}}\right) = P(Z \geq 0,65) = 0,26$.

Dacă $m_1 = m_0 + 3\sigma_{\bar{X}}$, puterea testului crește: $P(\text{respinge } H_0 \mid m = m_0 + 3\sigma_{\bar{X}}) = P(Z \geq -1,35) = 0,91$. Această ultimă valoare este mare, deci, dacă cumva media $m = m_0 + 3\sigma_{\bar{X}}$, atunci testul de mai sus poate detecta acest lucru cu o probabilitate de 0,91.

• Numim curbă de putere a testului de mai sus graficul funcției de putere $P(\bar{X} \geq m_0 + 1,65\sigma_{\bar{X}} \mid m = m_1)$, ca funcție de variabila m_1 . Să o notăm cu $p(m_1)$.

Două teste se pot compara și putem spune care dintre ele este mai tare. Mai mult, în anumite condiții putem alege testul "cel mai puternic". Noi însă nu ne ocupăm aici de aceste lucruri. În cadrul rezumatului vom lucra cu funcția $p(m_1)$ și vom construi "cel mai puternic test" în sensul precizat acolo. Ne limităm la analiza funcției de putere din figura următoare:



Curba de putere depinde de alegerea pragului de semnificație α . Este natural ca pentru α mai mare puterea testului să crească, după cum se vede în figură. De asemenea, dacă mărim volumul de selecție n , puterea testului crește pentru valori $m_1 > m_0$ și scade pentru valori $m_1 < m_0$ (de ce?).

• În exemplele 2 și 3 s-a testat media în cazul când dispersia populației era cunoscută. Putem observa că dacă micșorăm dispersia, puterea testului crește. Pe baza teoriei intervalor de încredere (Lecția 10) se pot construi teste parametrice în care și dispersia este necunoscută.

11.4 Rezumat

• Partea teoretică

Fie X un model statistic pentru o populație P , model care depinde de un parametru $\theta \in \Theta$. Fie H o submulțime specificată în Θ . Ea va fi numită ipoteză. Fie $K = \Theta \setminus H$. Spunem că:

- H se realizează dacă $\theta \in H$ (prin natura lucrurilor)
- H *nu* se realizează dacă $\theta \in K$.

În Lecția 11 am notat H cu H_0 și K cu H_1 .

Să notăm cu Ω spațiul tuturor observațiilor pe care le vom face asupra populației P (spațiul de selecție). Pentru fiecare $\theta \in \Theta$ avem o distribuție p_θ pentru X . Dacă pentru un sondaj $\omega \in \Omega$, găsim că $\theta \in H$, spunem că acceptăm ipoteza H , dacă $\theta \in K$ spunem că respingem ipoteza H . Spațiul de selecție se împarte în două părți distincte: $\Omega = A \cup R$, unde A este *zona de acceptare* a ipotezei H , adică $A = \{\omega \in \Omega \mid \theta \in H\}$, iar $R = \{\omega \in \Omega \mid \theta \in K\}$ este *zona de respingere* a ipotezei H .

Să notăm cu $p_\theta(R) = \beta_\theta$ probabilitatea de a respinge ipoteza H când de fapt θ este bine ales de natură. Se zice *eroare de primul tip* respingerea ipotezei H când ea este adevărată:

$\theta \in H$ și $\omega \in R$. Ea are probabilitatea β_θ .

Se zice *eroare de al doilea tip* acceptarea ipotezei H când ea de fapt nu este adevărată: $\theta \in H$ și $\omega \in A$. Testul are *pragul* ε dacă $\beta_\theta \leq \varepsilon$, $(\forall) \theta \in H$. Funcția $\theta \rightsquigarrow \beta_\theta$, $\theta \in K$ se zice *puterea testului*.

Pentru a face un test se alege o statistică Y care este "mică" când H se realizează și este "mai mare" când H nu se realizează. Se alege un număr y . Regiunea de respingere este de forma: $\text{Resping } H \iff Y \geq y$. Testul este de prag ε dacă $\theta \in H$ implică $p_\theta(Y \geq y) \leq \varepsilon$. Dacă nu se dă dinainte pragul ε , se observă valorile y ale statisticii Y și maximul probabilității $p_\theta(Y \geq y)$ când $\theta \in H$. Dacă acest maxim este mic se respinge ipoteza H . Dacă acest maxim este mare se acceptă ipoteza H .

Partea practica

1) *Test asupra mediei m a unei populatii normale cu dispersia cunoscuta σ^2*

$Z = \frac{\bar{X}-m}{\sigma/\sqrt{n}} \sim N(0,1)$; $H_0 : m = m_0$; $H_1 : m \neq m_0$ (test bilateral), sau $H_1 : m > m_0$ (sau $m < m_0$) (test unilateral)

2) *Test asupra mediei m cu dispersia σ^2 necunoscuta ($n \geq 30$)*

Se estimează σ^2 cu $S^2 = \frac{1}{n} \sum (x - \bar{x})^2$ și se folosește statistica $Z = \frac{\bar{X}-m}{S/\sqrt{n}} \sim N(0,1)$. H_0 și H_1 apar exact ca la **1**).

3) *Test asupra mediei m cu dispersia σ^2 necunoscuta ($n \leq 30$)*

Se estimează σ^2 cu $S'^2 = \frac{1}{n-1} \sum (x - \bar{x})^2 = \frac{nS^2}{n-1}$ și se folosește statistica

$$T = \frac{\bar{X} - m}{S/\sqrt{n-1}} \left(= \frac{\bar{X} - m}{S'/\sqrt{n}} \right) \sim t(n-1)$$

adică distribuția Student cu $n-1$ grade de libertate. H_0 și H_1 apar exact ca la **1**).

4) *Test pentru proporția de succese*

$P_s \sim N\left(p, \frac{pq}{n}\right)$, $q = 1 - p$. $Z = \frac{P_s - p}{\sqrt{\frac{pq}{n}}} \sim N(0,1)$. H_0 și H_1 apar exact ca la **1**).

5) *Erori de tipul I și erori de tipul II*

	Situatia reala	Concluzia noastra		
(1)	H_0 este adevarata	Acceptam H_0	Decizie corecta	
(2)	H_0 este adevarata	Respingem H_0	Decizie falsa	Eroare I
(3)	H_0 este falsa	Acceptam H_0	Decizie falsa	Eroare II
(4)	H_0 este falsa	Respingem H_0	Decizie corecta	

$\alpha = P(\text{Eroare de tipul I}) = P(\text{resping } H_0 \mid H_0 = \text{adevarata})$

$\beta = P(\text{Eroare de tipul II}) = P(\text{accept } H_0 \mid H_1 = \text{adevarata})$

11.5 Exerciții rezolvate (mai dificile)

1. Fie X o v.a. Poisson de medie λ . Considerăm ipoteza $H: \lambda \leq 0,5$. Fie X_1, \dots, X_{20} un esantion din X și punem $Y = X_1 + X_2 + \dots + X_{20}$. Definiți cu ajutorul lui Y (care are tot o distribuție Poisson) cel mai puternic test posibil de prag $0,05$. Pentru acest test determinați:

- 1) Probabilitatea erorii de primul tip pentru $\lambda = 0,3$.
- 2) Probabilitatea erorii de al doilea tip pentru $\lambda = 0,8$.

Presupunem că găsim $Y=20$. Cu ce prag se poate accepta H ?

Soluție. Cum $Y=20\bar{X}$ este natural să considerăm un test de forma: Respinge $H \iff Y \geq y$. Când H se verifică, cea mai mare probabilitate de a respinge H se obține pentru $\lambda = 0,5$. În acest caz Y este Poisson de parametru $20 \times 0,5 = 10$. Tabelul arată ca $P(Y \geq y) \leq 0,05$ pentru $y \geq 16$. Testul cerut se obține pentru cel mai mic y posibil (pentru a avea puterea maximă). Deci testul este: Respinge $H \iff Y \geq 16$. Pentru $\lambda = 0,3$, Y este Poisson de parametru $20 \times 0,3 = 6$ și probabilitatea erorii de primul tip este $P(Y \geq 16) = 0,0005$. Pentru $\lambda = 0,8$, Y este Poisson de parametru $20 \times 0,8 = 16$ și probabilitatea erorii de al doilea tip este $P(Y \leq 15) = 0,4667$.

Presupunem că $Y=20$. Maximul lui $P(Y \geq 20)$, când H este verificată, se obține pentru $\lambda = 0,5$. În acest caz Y este Poisson de parametru 10. Se găsește în acest caz $P(Y \geq 20) = 0,0035$ și trebuie să acceptăm ipoteza pentru un prag $\leq 0,0035$.

2. Fie X o v.a. gaussiană de dispersie 1 și de medie 0 sau 1. Vrem să testăm ipoteza $H: M(X)=0$ contra alternativei $M(X)=1$ cu un prag $\varepsilon = 0,05$. Se ia pentru acesta un esantion de volum n .

- 1) Definiți un test (cel mai puternic posibil).
- 2) Plecând de la ce valoare a lui n testul obținut va avea puterea $\geq 0,9$?

Soluție. Considerăm un test de forma: Respinge $H \iff \bar{X} \geq y$. Fie $\lambda = M(X)$. Variabila aleatoare $T = \sqrt{n}(\bar{X} - \lambda)$ este gaussiană redusă. În tabele găsim $P(T \geq t) = 0,05$ pentru $t = 1,6449$. $P(T \geq u) = 0,05$ pentru $u = -1,2816$.

1) n este considerat fix. Să găsim y astfel încât testul să fie de prag ε . Dacă H se verifică, $\lambda = 0$ și $T = \sqrt{n} \cdot \bar{X}$. Testul va avea pragul ε dacă $P(\bar{X} \geq y) \leq \varepsilon$, sau $P(T \geq \sqrt{n}y) \leq \varepsilon$, adică $\sqrt{n}y \geq t$. Testul va fi cel mai puternic dacă $y = \frac{t}{\sqrt{n}}$, deci avem testul:

$$\text{Respinge } H \iff \bar{X} \geq \frac{1,6449}{\sqrt{n}}.$$

2) Presupunem că H nu se verifică, adică $\lambda = 1$ și $T = \sqrt{n}\bar{X} - \sqrt{n}$. Se comite o eroare de al doilea tip pentru $\bar{X} \leq y$ sau $T \leq \sqrt{n}y - \sqrt{n} = t - \sqrt{n}$. Puterea testului va fi deci $\geq 0,9$ dacă $P(T \leq t - \sqrt{n})$, sau $t - \sqrt{n} < u$, deci $\sqrt{n} > t - u = 2,9265$. Cel mai mic n pentru care puterea testului este $\geq 0,9$ este $n = 9$.

3. Fie X pretul aceleiasi articol luat la întâmplare din 15 magazine. Găsim tabelul acestor

preturi în \$:

42,7	42,6	43,0	43,5	42,8
43,1	43,6	42,9	41,6	42,8
42,9	43,2	42,6	43,1	43,1

a) Se poate admite ipoteza $M(X)=43,0$?

b) Se poate admite ipoteza $D(X)=0,1$? În ambele cazuri se ia pragul $\varepsilon = 0,05$ și se consideră legea gaussiană pentru X .

Soluție. 1) Avem că $n = 15$, $\bar{X} = 42,9$, $S^2 = 0,2$. De aici deducem că $|T_0| = \frac{|\bar{X}-43|}{\sqrt{S^2/(n-1)}} = 0,8366$. O variabilă aleatoare cu $n-1=14$ grade de libertate are o probabilitate de cel puțin 0,40 pentru a lua o astfel de valoare. Ipoteza $M(X)=43$ este perfect acceptabilă.

2) O estimare nedeplasată pentru $\sigma^2 = D(X)$ este $\frac{n}{n-1}S^2 = 0,21$. Această valoare este dublul valorii testate 0,1. Vrem să vedem dacă nu cumva ea este prea mare. Pentru aceasta folosim faptul că $\frac{nS^2}{\sigma^2}$ are o distribuție χ^2 (Pearson) cu $n-1$ grade de libertate. Găsim că $\frac{nS^2}{0,1} = 30$. Dar χ^2 cu 14 grade de libertate are o probabilitate $< 0,01$ pentru a lua o astfel de valoare ridicată. Cu pragul 0,05 va trebui deci să respingem ipoteza $\sigma^2 = 0,1$.

4. Într-un oras A, 300 de locuitori din 1500 interogați declară că nu s-au uitat niciodată la TV. Într-un alt oras B, 320 din 1800 declară același lucru. Ce credeți despre ipoteza H: proporția de locuitori care nu se uită la TV este aceeași în ambele orase. (pragul 0,05).

Soluție Vom accepta pentru cele două v.a. X și Y ce reprezintă numărul acelor care nu se uită deloc la TV că au o distribuție uniformă. Presupunem că X_1, \dots, X_{1500} și Y_1, \dots, Y_{1800} sunt independente. Avem că $m = 1500$, $n = 1800$, $\bar{X} = \frac{300}{1500} = 0,2$, $\bar{Y} = \frac{320}{1800} = 0,172$; $S_X^2 = \frac{1}{m} \sum \bar{X}_i^2 - \bar{X}^2 = \bar{X}(1 - \bar{X})$, $S_Y^2 = \bar{Y}(1 - \bar{Y})$; $|T_0| = \frac{|\bar{Y}-\bar{X}|}{\sqrt{\frac{1}{m}S_X^2 + \frac{1}{n}S_Y^2}} = 1,61$.

O variabilă gaussiană redusă are o probabilitate mai mare decât 0,1 pentru a lua o valoare atât de mare. Prin urmare, la pragul de 0,05 ipoteza H se acceptă.

5. Pentru a măsura o masă μ putem folosi două procedee A și B. Rezultatul măsurătorii lui μ prin procedeul A este o v.a. gaussiană X cu media μ și cu dispersia σ_X^2 . Prin procedeul B, avem o variabilă gaussiană Y cu media μ și dispersia σ_Y^2 . S-au făcut 8 măsuri independente pentru μ prin procedeul A care au condus, la o dispersie de selecție empirică egală cu 0,24. S-au făcut apoi 12 măsurători independente pentru aceeași masă μ prin procedeul B și s-a obținut dispersia de selecție 0,08. Ce credeți despre ipoteza H: cele două procedee A și B au aceeași precizie? (prag $\varepsilon = 0,05$)

Soluție. Testul se referă la ipoteza H: $\sigma_X^2 = \sigma_Y^2$, cu $m = 8$, $n = 12$, $S_X^2 = 0,24$, $S_Y^2 = 0,08$. De aici s-ar putea crede că $\sigma_X^2 > \sigma_Y^2$. Vrem să vedem dacă nu cumva raportul $\frac{S_X^2}{S_Y^2}$ nu este prea mare (relativ la pragul ε) pentru ca ipoteza H să fie verificată.

Stim că variabilele aleatoare $\frac{mS_X^2}{\sigma_X^2}$ și $\frac{nS_Y^2}{\sigma_Y^2}$ sunt χ^2 cu $m-1 = 7$ și respectiv $n-1 = 11$ grade de libertate. Dacă cumva ipoteza H ar fi verificată v.a. $U = \frac{mS_X^2/(m-1)}{nS_Y^2/(n-1)}$ ar avea o

distributie F (Fisher-Snedecor) cu $\left\{ \begin{matrix} m-1 \\ n-1 \end{matrix} \right\} = \left\{ \begin{matrix} 7 \\ 11 \end{matrix} \right\}$ grade de libertate. În cazul nostru

$$U = \frac{8 \times 0,24/7}{12 \times 0,08/11} = 3,14.$$

Se vede din cercetarea directă a tabelului pentru distributia F că probabilitatea ca v.a. acesta să ia o astfel de valoare este ceva mai mică decât 0,05. Vom fi deci obligați să respingem ipoteza că σ_X^2 este cu mult mai mare decât σ_Y^2 . Totuși nu putem avea mare certitudine că $\sigma_X^2 = \sigma_Y^2$.

11.6 Exerciții propuse

1. Se arunca o moneda de 64 de ori. Testați la un nivel de semnificație de 5% dacă moneda este corectă, sau dacă moneda este contrafacută în favoarea "capului", sau dacă

a) apare "capul" de 38 ori;

b) apare "capul" de 42 ori.

R: a) $z=1,5$, corectă; b) $z=2,5$, contrafacută.

2. Un producător de hrană pentru câini pretinde că 8 din 10 câini preferă hrana produsă de el hranei produse de alți producători. Se aleg 120 de câini și se găsesc 88 care să prefere această hrană. Testați cu un nivel de 5% dacă pretentia producătorului este corectă.

R: $z=-1,826$, afirmația producătorului se respinge.

3. O mașină produce componente a căror greutate variază după o lege normală cu greutatea medie de 15,4 g și cu deviația standard $\sigma = 2,3$ g. Mașina s-a uzat cu timpul și se bănuiește că greutatea medie a pieselor produse de ea s-a micșorat. Se face o selecție aleatoare de 81 de piese și se găsește că masa medie a lor $\bar{x} = 15$ g. Oare acest lucru ne îndreptățește pe noi să credem cu un nivel de semnificație de 5% că greutatea medie a pieselor s-a micșorat? Presupunem că deviația standard σ nu s-a schimbat.

R: $z=-1,565$, Nu.

4. Un producător de castet audio pretinde că o casetă care durează 90 min. de obicei, durează de fapt 92 min. în medie cu deviația standard $\sigma = 1,8$ min. Se selectează 36 de benzi și se încearcă. Cel care verifică casetele respinge pretentia producătorului la un nivel de 5% și spune că media de timp al unei casete este mai mică decât 92 min. Ce puteți spune despre valoarea mediei de selecție pe care a cercetat-o statisticianul, valoare care l-a condus la decizia luată?

R: $\bar{x} < 91,51$ min.

5. După un sondaj făcut asupra a 300 de studenți s-a găsit că 35% fumează. Se poate respinge ipoteza $H_0 : p = 0,4$ în favoarea ipotezei $H_1 : p \neq 0,4$ la un prag de semnificație de 0,05?

6. Un esantion de volum 25 a fost extras dintr-o populatie normala, $X \sim N(m, 4)$. Media de selectie \bar{x} este 10,72. Fie ipoteza nula $H_0 : m = 10$ si ipoteza alternativa: a) $H_1 : m > 10$; b) $H_1 : m \neq 10$. Gasiti in ambele cazuri nivelul de semnificatie la care trebuie sa respingem ipoteza nula in favoarea ipotezei alternative.

R: a) $\varepsilon \geq 3,59\%$; b) $\varepsilon \geq 7,18\%$.

7. Un producator de becuri sustine ca media de viata a unui bec electric produs de el este de 2000 de ore. Se iau 64 de becuri la intamplare si se testeaza viata lor in ore, x .

Se obtine $\sum x = 127808$, $\sum (x - \bar{x})^2 = 9694,6$. Oare la un nivel de semnificatie de 2% putem spune ca producatorul si-a supraestimat produsul? Presupunem ca durata de viata a unui bec are distributia normala.

R: $n = 64 > 30$ este mai mare; deci estimam σ cu $\hat{\sigma} = S = \sqrt{\sum (x - \bar{x})^2 / n}$.

Folosim statistica $Z = \frac{\bar{X} - m}{\hat{\sigma} / \sqrt{n}}$ si gasim $z_{calc} = -1,95$ (test unilateral). Producatorul **nu** si-a supraestimat produsul la nivelul de 2%.

8. Dintr-o populatie normala X se extrage un esantion de volum 40 cu $\sum x = 24$ si $\sum x^2 = 596$. Testati (test bilateral) la nivelul de semnificatie de 5% daca media populatiei este 0 (estimati σ cu S).

R: $z = 0,995$. Se accepta afirmatia.

9. La un examen se analizeaza punctajul x obtinut de fiecare candidat in parte. Se aleg 250 de candidati si gasim ca $\sum x = 11872$ si $\sum x^2 = 646193$. Gasiti un interval de incredere de nivel 90% pentru media m . Testati ipoteza $H_0 : m = 49,5$ impotriva ipotezei $H_1 : m < 49,5$ cu nivelul de semnificatie $\alpha\%$ si gasiti α a.i. H_0 sa fie respinsa.

R: $(45,6; 49,4)$; $\alpha > 4$.

10. Un esantion de volum 8 dintr-o populatie normala are $\bar{x} = 4,65$, $\sum (x - \bar{x})^2 = 0,74$. Testati la un nivel de semnificatie de 2% daca media distributiei este 4,3.

R: Folositi statistica $T = \frac{\bar{X} - m}{S / \sqrt{n-1}} \in t(n-1)$, deoarece esantionul este mic ($8 < 30$); $t_{calc} = 3,05$ (test bilateral). Respingem afirmatia.

11. O masina produce ace de otel de lungume 2 cm. Se face o selectie de 10 ace si se gaseste: 1,98; 1,96; 1,99; 2,00; 2,01; 1,95; 1,97; 1,96; 1,97; 1,99. Presupunand ca lungimile acelor sunt normal distribuite, testati cu 1% daca masina functioneaza bine sau nu.

R: $n = 10 < 30$; folosim deci statistica T ; $t_{calc} = -3,601$. Nu functioneaza bine.

12. Se aleg 8 femei la intamplare si li se masoara nivelul de colesterol: 3,1; 2,8; 1,5; 1,7; 2,4; 1,9; 3,3; 1,6.

Vrem sa testam daca nivelul mediu de colesterol este 3,1.

a) Presupunand ca selectia face parte dintr-o populatie normala aratati de ce testul T este cel mai potrivit.

b) Aplicati testul asupra mediei 3,1 si gasiti cu nivelul de semnificatie de 2% daca media 3,1 este corecta sau nu.

c) Gasiti un interval de incredere de 90% pentru nivelul mediu al colesterolului.

R: b) Resping H_0 ($m=3,1$); c) (1,81;2,72).

13. Teoria prezice realizarea unui eveniment A cu probabilitatea $p=0,4$. Se experimenteaza realizarea evenimentului A de 400 de ori. Din cele 400 de incercari A se realizeaza de 140 de ori. Testati cu nivelul de semnificatie de 1% daca p este mai mica sau nu decat 0,4.

R: Se modeleaza cu proportia de succese: $P_s \sim N\left(p, \frac{pq}{n}\right)$; $Z = \frac{P_s - p}{\sqrt{\frac{pq}{n}}}$; $z_{calc} = -2,04$. Nu.

14. Se cerceteaza proportia de studenti care au un calculator personal. Din 200 de studenti, 143 au un calculator personal. Testati cu un nivel de 5% ipoteza $H_0: p=75\%$ (probabilitatea ca un student sa aiba un calculator) contra ipotezei H_1 : probabilitatea p este mai mica decat 75%.

R: $z_{calc} = -1,143$. Se accepta.

15. Gasiti probabilitatile erorilor de tipul I si celor de tipul II in testarea urmatoarelor ipoteze. Se stie ca o cutie poate sa contina fie (H_0) 10 jetoane albe si 90 negre, fie (H_1) 50 jetoane albe si 50 negre. Pentru a testa ipoteza H_0 impotriva ipotezei H_1 se aleg 4 jetoane din cutie fara a le pune inapoi. Daca toate 4 jetoane sunt negre, accept H_0 . Altfel, o resping (regula de decizie a testului).

Indicatie H_0 : cutia contine 10 jetoane albe si 90 negre.

H_1 : cutia contine 50 jetoane albe si 50 negre.

$P(\text{Eroare I}) = P(\text{resping } H_0 \mid H_0 \text{ este adevarata}) = P(\text{cel putin un jeton este alb} \mid \text{in cutie sunt 10 albe si 90 negre}) = 1 - P(\text{toate 4 sunt negre} \mid \text{in cutie sunt 10 albe si 90 negre}) = 1 - \frac{90}{100} \cdot \frac{89}{99} \cdot \frac{88}{98} \cdot \frac{87}{97} = 1 - 0,652 = 0,348$.

$P(\text{Eroare II}) = 0,059$.

16. Pentru datarea specimenelor arheologice se foloseste faptul ca aceste specimene emit particule radioactive. Numarul de particule emise in n minute are o distributie Poisson cu parametrul $n\lambda$, unde λ este un parametru ce depinde de varsta specimenului.

Se fac doua ipoteze asupra varstei unui specimen

H_0 : specimenul este de 7000 de ani ($\lambda = 0,1$)

H_1 : specimenul este de 15000 de ani ($\lambda = 4,0$)

S-a decis sa se contorizeze numarul X al particulelor radioactive emise in n minute de specimen si

Acceptam H_0 , daca $X \leq 1$ (respingem H_1)

Acceptam H_1 , daca $X \geq 2$ (respingem H_0)

Daca $n = 1$, se cere: a) $P(\text{resping } H_0 \mid H_0 = \text{adevarata})$ si b) $P(\text{resping } H_1 \mid H_1 = \text{adevarata})$.

Presupunem acum ca $P(\text{resping } H_1 \mid H_1 = \text{adevarata}) \leq 0,001$; aratati ca numarul minim de minute complete necesare inregistrarii este de 3 minute.

Pentru acest numar de 3 minute sa se calculeze $P(\text{resping } H_0 \mid H_0 = \text{adevarata})$.

Indicatie $X \sim Po(n\lambda)$. Daca $n = 1$, $X \sim Po(\lambda)$ si a) $P(\text{resping } H_0 \mid H_0 = \text{adevarata}) = P(X \geq 2 \mid X \sim Po(1, 0))$. $P(X \geq 2) = 1 - P(X = 0) - P(X = 1) = 1 - e^{-1} - e^{-1} = 1 - 2e^{-1} = 1 - 0,736 = 0,264$. Deci $P(\text{resping } H_0 \mid H_0 = \text{adevarata}) = 0,246$.

b) $P(\text{resping } H_1 \mid H_1 = \text{adevarata}) = P(X \leq 1 \mid X \sim Po(4)) = 5e^{-4} = 0,092$. Daca $P(\text{resping } H_1 \mid H_1 = \text{adevarata}) \leq 0,001$, atunci $P(X \leq 1 \mid X \sim Po(4n)) \leq 0,01$. Daca $X \sim Po(4n)$, atunci $P(X \leq 1) = e^{-4n} (1 + 4n) \leq 0,001$ si gasim $n = 3$ si $P(\text{resping } H_0 \mid H_0 = \text{adevarata}) = P(X \geq 2 \mid X \sim Po(3 \cdot 1)) = 1 - 4e^{-3} = 0,801$.

17. Se fac doua ipoteze asupra functiei densitate de probabilitate pentru o v.a. X .

$$H_0 : f(x) = \begin{cases} \frac{1}{4}(x+1), & 0 < x < 2 \\ 0, & \text{altfel} \end{cases}$$

$$H_1 : f(x) = \begin{cases} \frac{1}{4}x^3, & 0 < x < 2 \\ 0, & \text{altfel} \end{cases}$$

Se construiesc urmatorul test. Se face o singura observatie asupra v.a. X si daca $X < k$, cu k dat, $0 < k < 2$, atunci H_0 se accepta. Altfel H_1 se accepta: a) Gasiti k a.i. $P(\text{Eroare I}) = 0,1$. b) Cu valoarea lui k de la a) gasiti $P(\text{Eroare II})$

Indicatie Faceti graficele pentru $f(x)$ in fiecare din cazurile H_0 si H_1 . $P(\text{accept } H_1 \mid H_0 = \text{adevarata}) = 0,1$ implica $P(X \geq k \mid f(x) = \frac{1}{4}(x+1)) = 0,1$ sau $\int_k^2 \frac{1}{4}(x+1) dx = 0,1$, de unde $k = 1,86$.

$$P(\text{Eroare II}) = P(X < 1,86 \mid f(x) = \frac{1}{4}x^3) = \int_0^{1,86} \frac{1}{4}x^3 dx = 0,748.$$

18. Dintr-o populatie normala $N(m, 36)$ se ia un esantion de volum 100. Un cercetator vrea sa testeze ipotezele $H_0 : m = 65$, $H_1 : m > 65$. El decide sa foloseasca urmatoarea regula de decizie:

accept H_0 daca media de selectie $\bar{x} \leq 66,5$

resping H_0 daca $\bar{x} > 66,5$.

a) Gasiti $P(\text{Eroare I})$;

b) Daca el foloseste alternativa $H_1 : m = 67,9$, gasiti $P(\text{Eroare II})$.

c) Ce valoare critica trebuie sa considere el pentru media de selectie daca vrea ca

$$P(\text{Eroare I}) = P(\text{Eroare II})$$

Indicatie a) Sub H_0 avem: $\bar{X} \sim N(65, \frac{36}{100})$. El respinge H_0 daca $\bar{x} > 66,5$, adica $P(\bar{X} > 66,5) = P(Z > 2,5) = 0,00621$. Deci $P(\text{Eroare I}) = 0,00621$.

b) Daca el considera $H_1 : m = 67,9$, atunci sub H_1 avem ca $\bar{X} \sim N(67,9, \frac{36}{100})$ si

$$P(\text{Eroare II}) = P(\text{accept } H_0 \mid H_1 \text{ este adevarata}) = P(\bar{X} \leq 66,5 \mid m = 67,9)$$

$$P(\bar{X} \leq 66,5 \mid m = 67,9) = P(Z \leq -2,333) = 0,00982.$$

c) $P(\bar{X} > \bar{x} \mid H_0 \text{ este adevarata}) = P(\bar{X} \leq \bar{x} \mid H_1 \text{ este adevarata})$, deci \bar{x} se afla la mijlocul segmentului $[65; 67,9]$, adica $\bar{x} = 66,45$.

19. Ingredientele care se amesteca pentru a forma betonul au asemenea proportii incat rezistenta medie la rupere sa fie de 2000N. Daca rezistenta medie la rupere cade sub 1800N atunci compozitia trebuie schimbata. Distributia rezistentei de rupere este normal distribuita cu deviatia standard de 200N.

Se iau esantioane pentru a se cerceta ipotezele:

$$H_0 : m = 2000N$$

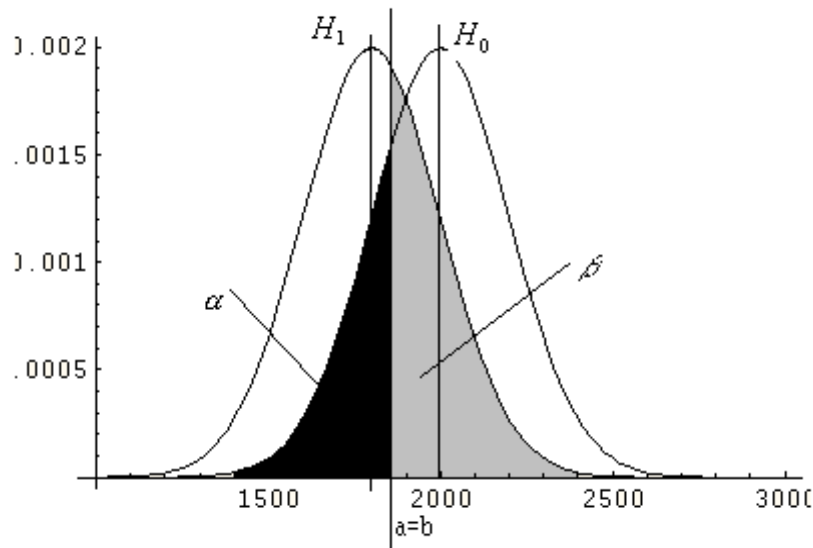
$$H_1 : m = 1800N$$

Cate esantioane trebuie testate pentru ca sa avem:

$$P(\text{Eroare I}) = \alpha = 0,05 \text{ si}$$

$$P(\text{Eroare II}) = \beta = 0,1.$$

Indicatie Sub H_0 avem ca $X \sim N(2000, 200^2)$, deci $\bar{X} \sim N\left(2000, \frac{200^2}{n}\right)$



Daca Z_α este cuantila ce corespunde lui $\alpha = 0,05$, avem ca $Z_\alpha = -1,645$ si valoarea corespunzatoare a lui \bar{X} este $a = 2000 - 1,645 \left(\frac{200}{\sqrt{n}}\right)$. Sub H_1 , $\bar{X} \sim N\left(1800, \frac{200^2}{n}\right)$, deci, pentru $Z_\beta = 1,282$ ($\beta = 0,1$) avem ca $b = 1800 + 1,282 \left(\frac{200}{\sqrt{n}}\right)$. Din $a = b$ gasim $n = 8,57$, deci vom lua 9 probe.

20. Se aleg aleator esantioane de 400 seminte dintr-un anumit sortiment. Probabilitatea ca o seminta sa germineze este egala cu α . Sa notam cu X v.a. ce reprezinta numarul de seminte care au germinat din totalul semintelor dintr-un esantion.

Folositi o aproximare convenabila a modelului probabilistic cu un model normal pentru a determina:

a) $P(X \leq 340 \mid \alpha = 0,9)$

b) $P(X \geq 340 \mid \alpha = 0,8)$

c) Controlorul de calitate al semintelor stie ca α este fie 0,8, fie 0,9. Sa presupunem ca, de fapt, din totalul de 400 de seminte dintr-un esantion au germinat numai x . Controlorul decide ca valoarea lui α este 0,8 daca

$$Z = P(X \geq x \mid \alpha = 0,8) - P(X \leq x \mid \alpha = 0,9)$$

este pozitiv. Altfel el decide ca $\alpha = 0,9$. Gasiti care este decizia controlorului pentru fiecare dintre cazurile $x = 330$, $x = 340$, $x = 350$.

R: a) 0,000577; b) 0,00738; c) 0,8; 0,8; 0,9.

Lecția 12

Testul neparametric χ^2

12.1 Principiul testului χ^2

Vom introduce acest test prin analiza atentă a unui exemplu.

Exemplul 12.1 (*fictiv*) S-a făcut un sondaj într-un oras din România pe un esantion de 200 de tineri bărbați de 27 de ani, în 1985, asupra situației pregătirii lor școlare. Ei au fost împărțiți în 6 categorii după cum urmează:

1. cu studii superioare terminate;
2. cu studii superioare neterminate dar începute;
3. cu liceul de 12 ani terminat;
4. cu liceul de 12 ani neterminat dar început;
5. cu școala generală de 8 ani terminată;
6. cu școala generală de 8 ani începută și neterminată.

Același sondaj se repetă (în aceleași condiții) în 1995. Vom nota cu $f_{o,j}$, $j = 1, 2, \dots, 6$ frecvența (absolută) observată în 1995 pentru categoria "j" din cele 6 expuse mai sus. Vom nota cu $f_{s,j}$ frecvența (sperată, la care ne așteptăm și în 1995!) găsită în 1985. Iată tabelul cu cele două sondaje:

categoria "j"	Frecvența observată în 1995, $f_{o,j}$	Frecvența observată în 1985, $f_{s,j}$
1	35	36
2	40	34
3	83	64
4	16	26
5	26	34
6	0	6
	— — — 200	— — — 200

Facem următoarea ipoteză statistică (inferență statistică) $H_0: \langle\langle \text{distributia populatiei în 1995 este aceeași cu distributia populatiei din 1985} \rangle\rangle$.

Este "natural" (de ce?) ca "discrepanța" sau deviația dintre cele două situații să o măsurăm prin suma

$$\sum_{j=1}^6 \frac{(f_{o,j} - f_{s,j})^2}{f_{s,j}} \quad (12.1)$$

Se arată (Teorema lui K. Pearson) că pentru n mare (cel puțin 20-25), în cazul nostru $n=200$, această sumă tinde către valoarea distribuției χ^2 cu 6-1 grade de libertate în ipoteza că H_0 este adevărată: avem aceeași distribuție, adică distribuțiile celor două sondaje concordă. În general, dacă am repartizat elementele din selecție în J clase, pentru un volum mare ($n \geq 25$) suma (12.1) cu J în loc de 6, este aproximativ egală cu valoarea distribuției χ^2 cu $J - 1$ grade de libertate.

Ipoteza H_0 are șanse să fie adevărată dacă valoarea sumei (12.1) este mică. Găsim

$$\chi^2 = \frac{(35 - 36)^2}{36} + \frac{(40 - 34)^2}{34} + \dots + \frac{(0 - 6)^2}{6} = 18,46$$

Căutăm în Tabelul III pe linia corespunzătoare $\nu = 5$ (grade de libertate) și găsim că $16,7496 \leq 18,46 \leq 20,515$.

Dar

$$\begin{aligned} P(16,7496 \leq \chi^2 \leq 20,515) \\ = F(20,515) - F(16,7496) = 0,999 - 0,995 = 0,004 \end{aligned}$$

deci extrem de mică. Prin urmare, este foarte puțin probabil ca cele două distribuții să concorde. Deci se respinge ipoteza H_0 cu probabilitatea $1 - 0,004 = 0,996$.

Retinem din Exemplul 1 că ori de câte ori vrem să facem o ipoteză asupra probabilității ca două distribuții să concorde putem folosi "testul χ^2 ", adică metodologia din 12.1, cu condiția ca numărul observațiilor să fie mare ($n \geq 25$), grupele de divizare ale sondajului să fie disjuncte și numărul lor k să respecte următoarele reguli:

$$\begin{aligned} 25 &< n \leq 100, k \in (10, 15) \\ 100 &< n < 200, k \in (15, 18) \\ 200 &\leq n < 400, k \in (18, 20) \\ 400 &\leq n < 1000, k \in (25, 30) \\ 1000 &\leq n < 2000, k \in (35, 40) \end{aligned}$$

Acste reguli au fost deduse din practică și nu teoretic.

Prezentăm acum situația generală în care putem utiliza testul de concordanță χ^2 . Facem următoarea ipoteză: $H_0: \ll \text{Presupunem că populația } P \text{ are o funcție de probabilitate } g(x), \text{ specificată (de exemplu normală cu } m = 0 \text{ și } \sigma = 2) \gg \gg$.

Vrem să folosim testul χ^2 pentru a "măsura" câtă dreptate avem să presupunem acest lucru pornind de la un sondaj de volum n : x_1, x_2, \dots, x_n din populația P . După regulile de mai sus împărțim datele x_1, x_2, \dots, x_n în J grupe disjuncte. Notăm cu ν_j frecvența absolută în grupa " j " (numărul acelor x_i care se află în grupa " j "). Notăm cu p_j probabilitatea teoretică (obținută cu ajutorul funcției de probabilitate $g(x)$, sau cu ajutorul funcției de repartiție corespunzătoare $G(x)$) ca un element x din populația P să se afle în grupa j . Atunci frecvența teoretică este np_j , și suma care măsoară deviația din formula (12.1) devine:

$$d = \sum_{j=1}^J \frac{(\nu_j - np_j)^2}{np_j} \quad (12.2)$$

Calculăm acest număr d . Privim Tabelul III și încercăm să "estimăm" probabilitatea ca χ^2 să aibă valoarea d . De obicei se fixează un prag de semnificație $\alpha \in (0, 1)$. Se caută cuantila χ_α^2 corespunzătoare acestui prag, adică acea valoare a lui χ^2 astfel încât funcția de repartiție a lui χ^2 să aibă valoarea α . Cum $P(\chi^2 \leq \chi_\alpha^2) = \alpha$, din definiția funcției de repartiție, vom face următorul raționament:

- dacă $d \geq \chi_\alpha^2$, vom accepta ipoteza H_0 (distribuția populației este de forma prescrisă) cu pragul de semnificație α .
- dacă $d < \chi_\alpha^2$, vom respinge ipoteza H_0 cu pragul de semnificație α .

Acesta este testul de concordanță χ^2 . De obicei α se ia mic: $\alpha = 0,05; 0,01; 0,001$.

Exemplul 12.2 Se fac 500 de măsurători asupra erorilor date de un aparat de măsură de la bordul unui avion. Ele se împart în 8 intervale consecutive. Frecvențele absolute cu care apar aceste erori pe un interval sunt date în următorul tabel:

$I_j :$	$[-4; -3)$	$[-3; -2)$	$[-2; -1)$	$[-1; 0)$	$[0; 1)$	$[1; 2)$	$[2; 3)$	$[3; 4)$
$\nu_j :$	6	25	72	133	120	88	46	10

Utilizând testul χ^2 să se verifice cu pragul de semnificație $\alpha = 0,95$ dacă distribuția erorilor este normală cu media estimată la $m = 0,168$ și dispersia estimată la $\sigma^2 = 1448^2$.

Soluție Aici avem un exemplu mai complicat deoarece cei doi parametri au deja estimări date. Vom avea mereu două relații de legătură: $0,168 = \frac{1}{500} \sum x_i$ și $1448^2 = \frac{1}{500} \sum (x_i - 0,168)^2$.

Deci "numărul gradelor de libertate" va scădea cu doi. Prin urmare χ^2 va avea $8-1-2=5$ grade de libertate.

Calculăm probabilitatea teoretică pe fiecare interval $[x_i, x_{i+1})$:

$$p_j = \Phi\left(\frac{x_{j+1} - m}{\sigma}\right) - \Phi\left(\frac{x_j - m}{\sigma}\right)$$

, unde Φ este funcția lui Laplace. Găsim pentru np_j următoarele valori (corespunzătoare intervalelor precizate deja în tabelul de mai sus): 6,2; 26,2; 71,2; 122,2; 131,8; 90,5; 32,8; 10,5. Calculăm $d = \sum_{j=1}^8 \frac{(v_j - np_j)^2}{np_j} = 3,94$.

Cuantila $\chi_{0,95}^2$ pentru 5 grade de libertate este 11,0705 (vezi Tabelul III). Deci $P(\chi^2 \leq 11,0705) = 0,95$, adică $P(\chi^2 > 11,0705) = 0,05$.

Cum $d = 3,94 < 11,0705$ acceptăm ipoteza de normalitate cu probabilitatea 0,95 (adică cu 95%), sau cu pragul $\alpha = 0,05$.

Vom prezenta acum testul χ^2 (se citește chi pătrat sau chi doi) dintr-un punct de vedere mai general.

Fie evenimentele $\{A_1, A_2, \dots, A_r\}$ cu probabilitatea $p_i = P(A_i)$, $i = \overline{1, r}$, cunoscute mai mult sau mai puțin. Putem cere ca $\sum p_i = 1$, adică $\{A_1, A_2, \dots, A_r\}$ să fie o descompunere a evenimentului singur. În esență testul χ^2 își propune să testeze o ipoteză oarecare H privind probabilitățile p_1, \dots, p_r . Se fac sondaje de volum mare pentru fiecare A_i , $i = \overline{1, r}$. Se notează cu Y_i numărul acelor probe care sunt favorabile evenimentului A_i .

Cazul I Se dau numerele $p'_1, p'_2, \dots, p'_r \geq 0$ cu $\sum p'_i = 1$ și se consideră ipoteza H: $p_1 = p'_1, p_2 = p'_2, \dots, p_r = p'_r$ (problemă de concordantă). Dacă H este adevărată statistica (Helmert-Pearson)

$$T = \sum_{1 \leq i \leq r} \frac{(Y_i - np'_i)^2}{np'_i} = \left(\sum_i \frac{Y_i^2}{np'_i} \right) - n \quad (12.3)$$

are practic (adică pentru n mare) distribuția χ^2 cu $r - 1$ grade de libertate. Pentru a arăta ultima egalitate în (12.3) am folosit egalitățile $\sum Y_i = n$ și $\sum p_i = 1$. Testul constă în a respinge ipoteza H dacă T ia o valoare semnificativ prea mare relativ la χ^2 (vezi exemplul 12.1). De exemplu dacă se dau v.a. X și legea Q cunoscută, pentru a testa dacă H: X are distribuția Q, descompunem dreapta reală \mathbb{R} în subintervale disjuncte: $\mathbb{R} = A_1 \cup A_2 \cup \dots \cup A_r$, considerăm evenimentele $(X \in A_i)_{i=\overline{1, r}}$ și punem $p'_i = Q(A_i)$.

Cazul II Fie acum r funcții pozitive $f_1(x_1, \dots, x_s), \dots, f_r(x_1, \dots, x_s)$ de variabile x_1, \dots, x_s cu $s < r$, astfel încât $f_1(x_1, \dots, x_s) + \dots + f_r(x_1, \dots, x_s) = 1$. Considerăm ipoteza H: Există numerele $\lambda_1, \dots, \lambda_s$ (estimate printr-un procedeu oarecare) astfel încât $p_i = f_i(\lambda_1, \dots, \lambda_s)$ pentru $1 \leq i \leq r$. Dacă H este adevărată se pot defini (plecând de la valorile Y_1, \dots, Y_r) estimatori convenabili pentru $\lambda_1, \dots, \lambda_s$. Fie $\hat{\lambda}_1 = g_1(Y_1, \dots, Y_r)$, $1 \leq i \leq s$, un estimator pentru λ_i . De aici găsim estimatori convenabili ai probabilităților p_i : $\hat{p}_i = f_i(\hat{\lambda}_1, \dots, \hat{\lambda}_s)$ pentru $1 \leq i \leq r$. Statistica

$$T = \sum_{1 \leq i \leq r} \frac{(Y_i - n\hat{p}_i)^2}{n\hat{p}_i} = \left(\sum_i \frac{Y_i^2}{n\hat{p}_i} \right) - n \quad (12.4)$$

are practic o distribuție χ^2 cu $r - s - 1$ grade de libertate (apar încă s legături datorită relațiilor de estimare). Testul constă în a respinge ipoteza H dacă T ia o valoare semnificativ

prea mare relativ la χ^2 (vezi Exemplul 2). Se pune acum problema de a găsi estimatori ”buni” pentru $\lambda_1, \dots, \lambda_n$. Iată câteva reguli bazate pe însăși definiția ”*formulei de deviatie*” (Helmert-Pearson) dată în (12.1) și (12.2) și pe alte noțiuni ce apar în Lectia 9.

Regula I (χ^2 minim) Se iau pentru $\hat{\lambda}_1, \dots, \hat{\lambda}_s$ acele funcții de Y_1, \dots, Y_r care minimizează expresia (deviația):

$$d = \sum_i \frac{(Y_i - n\hat{p}_i)^2}{n\hat{p}_i} = \min \quad (12.5)$$

adică acele $\hat{\lambda}_1, \dots, \hat{\lambda}_s$ care verifică ecuațiile diferențiale.

$$\sum_{1 \leq i \leq r} \frac{Y_i^2}{\hat{p}_i^2} \cdot \frac{\partial \hat{p}_i}{\partial \hat{\lambda}_j} = 0, 0 \leq j \leq s \quad (12.6)$$

Se pune condiția grad $d(\hat{\lambda}_1, \dots, \hat{\lambda}_s) = 0$ și se folosește faptul că $\sum_i \frac{\partial \hat{p}_i}{\partial \hat{\lambda}_j} = 0$, deoarece $\sum_j \hat{p}_j = 1$.

Regula II (verosimilitatea maximă) După un raționament asemănător cu acela din Lectia 9 (Principiul verosimilității maxime) se deduce că $\hat{\lambda}_1, \dots, \hat{\lambda}_n$ trebuie să verifice ecuațiile diferențiale:

$$\sum_{1 \leq i \leq r} \frac{Y_i}{\hat{p}_i} \cdot \frac{\partial \hat{p}_i}{\partial \hat{\lambda}_j} = 0 \quad (12.7)$$

pentru $1 \leq j \leq s$.

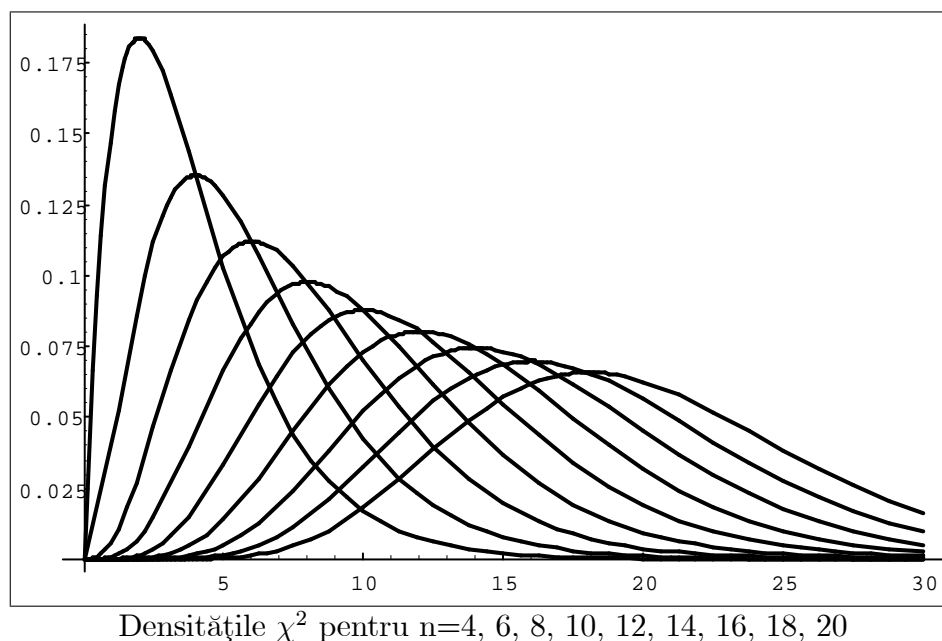
Regula III În locul formulelor (12.6) se pot folosi formulele:

$$\sum_{1 \leq i \leq r} \frac{\hat{p}_i}{Y_i} \cdot \frac{\partial \hat{p}_i}{\partial \hat{\lambda}_j} = 0 \quad (12.8)$$

pentru $1 \leq j \leq s$

(χ^2 minim modificat).

Această distribuție atât de folosită are pentru diverse valori ale lui n , densitățile ca în graficul următor:



Să aplicăm acum cele de mai sus la câteva tipuri de probleme.

12.1.1 Teste asupra formei unei distributii

Fie $P=Q(\lambda_1, \dots, \lambda_s)$ o familie de legi pe \mathbb{R} care depinde de s parametri $\lambda_1, \dots, \lambda_s$ și X o v.a. aleatoare. Pentru a testa ipoteza: $H: X$ se supune unei anume legi din familia P , împărțim pe \mathbb{R} în r subintervale: $\mathbb{R} = A_1 \cup \dots \cup A_r$, și punem $f_i(\lambda_1, \dots, \lambda_s) = Q(\lambda_1, \dots, \lambda_s)(A_i)$.

Exemplul 12.3 Fie X o v.a. ce poate lua valorile $0, 1, 2, \dots, k, \dots$. Testăm ipoteza: $H: X$ se supune unei legi Poisson.

Familia legilor Poisson depinde de un singur parametru $\lambda > 0$. Ca estimator "bun" pentru λ se alege $\hat{\lambda} =$ media de selecție (eventual după grupaje convenabile). Avem $s = 1$ în acest caz.

Exemplul 12.4 Vrem să testăm pentru o v.a. X continuă următoarea ipoteză $H: X$ se supune unei legi gaussiene. Aici P este familia $N(\mu, \sigma^2)$ care depinde de doi parametri: μ și σ^2 . După cum stim este convenabil să estimăm media μ cu media de selecție și pe σ^2 cu S^2 . Aici $s = 2$.

12.1.2 Teste de independentă

Fie X și Y două v.a. definite pe aceeași categorie de probe. Pentru a testa ipoteza $H: X$ și Y sunt independente, se descompune \mathbb{R} în două partitii: $\mathbb{R} = E_1 \cup \dots \cup E_a = F_1 \cup \dots \cup F_b$ și considerăm cele $r = a \cdot b$ evenimente: $X \in E_i$ și $Y \in F_j$, pentru $1 \leq i \leq a, 1 \leq j \leq b$.

Se face un număr mare de sondaje independente de tipul $(X_1, Y_1), \dots, (X_n, Y_n)$. Notăm cu N_{ij} numărul realizărilor $X \in E_i$ și $Y \in F_j$. Se introduc statisticile

$$\tilde{N}_{i\bullet} = \sum_{1 \leq j \leq b} N_{ij} \quad \tilde{N}_{\bullet j} = \sum_{1 \leq i \leq a} N_{ij} \quad (12.9)$$

$$T = n \sum_{i,j} \frac{\left(N_{ij} - \frac{1}{n} \tilde{N}_{i\bullet} \tilde{N}_{\bullet j}\right)^2}{\tilde{N}_{i\bullet} \tilde{N}_{\bullet j}} = n \left[\left(\sum_{i,j} \frac{N_{ij}^2}{\tilde{N}_{i\bullet} \tilde{N}_{\bullet j}} \right) - 1 \right] \quad (12.10)$$

Dacă H este adevărată, T se supune practic unei legi χ^2 cu $(a-1)(b-1)$ grade de libertate. Testul constă în a respinge ipoteza H dacă T ia o valoare semnificativ prea mare relativ la χ^2 . Dacă ipoteza H este verificată numărul $\hat{p}_{ij} = \frac{1}{n^2} \tilde{N}_{i\bullet} \tilde{N}_{\bullet j}$ este un estimator bun pentru probabilitatea $p_{ij} = P(X \in E_i \text{ și } Y \in F_j)$.

12.1.3 Teste de omogenitate

Se consideră t variabile aleatoare X_1, \dots, X_t . Pentru a testa ipoteza H : X_1, \dots, X_t satisfac aceeași lege, se face partiția $\mathbb{R} = A_1 \cup A_2 \cup \dots \cup A_r$, pentru fiecare i se ia un esantion $(X_{i1}, \dots, X_{in_i})$ de volum mare, $x_i \in X_i$, $n = n_1 + \dots + n_i$, v.a. X_{i1}, \dots, X_{in_i} fiind considerate independente. Pentru $1 \leq i \leq t$ se notează Y_{ij} numărul de indici k astfel încât $X_{ik} \in A_j$. Se introduc statisticile

$$\tilde{Y}_{\bullet j} = \sum_{1 \leq i \leq t} \tilde{Y}_{ij} \quad T = \sum_{\substack{1 \leq i \leq t \\ 1 \leq j \leq r}} \frac{\left(Y_{ij} - \frac{n_i}{n} \tilde{Y}_{\bullet j}\right)^2}{\frac{n_i}{n} \tilde{Y}_{\bullet j}} = n \left[\sum_{i,j} \frac{Y_{ij}^2}{n_i \tilde{Y}_{\bullet j}} - 1 \right]. \quad (12.11)$$

Dacă H este adevărată, T se supune practic unei legi χ^2 cu $(t-1)(r-1)$ grade de libertate. Testul constă în a respinge H atunci când T ia o valoare semnificativ mare relativ la χ^2 . Dacă H este adevărată $\hat{p}_j = \frac{1}{n} \tilde{Y}_{\bullet j}$ este un estimator bun pentru $P(X_i \in A_j)$ și nici nu depinde de i . Dacă $r = 2$ (12.10) devine

$$T = \frac{n^2}{\tilde{Y}_{\bullet 1} \tilde{Y}_{\bullet 2}} \sum_i \frac{Y_{i1}^2}{n_i} - n \frac{\tilde{Y}_{\bullet 1}}{\tilde{Y}_{\bullet 2}} \quad (12.12)$$

Exemplul 12.5 Teoria lui Mendel asupra eredității ne previne că dacă creștem 2 tipuri de plante va trebui să obținem produse de tipul A, B, C, D în proporție de 9, 3, 3 și 1. După experiențe se observă că s-au obținut 154 produse de tipul A , 44 de tip B , 63 de tip C și 21 de tip D . Ce părere aveți în acest caz de teoria lui Mendel (prag $\varepsilon = 0,05$)?

Soluție Aici evenimentele A_1, A_2, A_3, A_4 sunt A, B, C și D . Teoria lui Mendel prevede că $p_1 = \frac{9}{16}$, $p_2 = \frac{3}{16}$, $p_3 = \frac{3}{16}$, $p_4 = \frac{1}{16}$ (deoarece $9+3+3+1=16$). Suntem în *Cazul I* al testului χ^2

cu: $p'_1 = \frac{9}{16}$, $p'_2 = \frac{3}{16}$, $p'_3 = \frac{3}{16}$, $p'_4 = \frac{1}{16}$. Experiențele conduc la , $Y_1 = 154$, $Y_2 = 44$, $Y_3 = 63$, $Y_4 = 21$, $n = 154 + 44 + 63 + 21 = 282$.

Dacă U este variabila χ^2 cu $r-1 = 3$ grade de libertate, avem că $P(U \geq 7,81) = 0,05$. Testul de prag 0,05 se scrie: Respinge $H \iff \chi^2 \geq 7,81$. În cazul nostru avem np'_1 , $np'_2 = np'_3 = 53$, $np'_4 = 18$ și

$$\chi^2 = \frac{(154 - 159)^2}{159} + \frac{(44 - 53)^2}{53} + \frac{(63 - 53)^2}{53} + \frac{(21 - 18)^2}{18} = 4,06 < 7,81$$

. Se acceptă deci teoria lui Mendel ca fiind adevărată cu 0,95=1-0,05.

Exemplul 12.6 Fie p probabilitatea pentru ca o piesă din echipamentul fabricat de o anumită masină să aibă defecte. Vrem să testăm ipoteza $H: p = 0,2$. Pentru aceasta luăm 100 de piese și constatăm că 22 dintre ele au defecte. Care este probabilitatea, dacă ipoteza este adevărată, pentru ca χ^2 să aibă o valoare mare? Cum interpretați acest lucru?

Soluție Aici A_1 =succes și A_2 =eșec. Fie $p_1 = p = P(A_1)$ și $p = 1 - p = q = P(A_2)$. Vom avea $p'_1 = p'$ și $p'_2 = q' = 1 - p'$. Fie Y_1 =numărul de succese în cursul a n încercări(=Y) și $Y_2 = n - Y$. Avem

$$T = \frac{(Y - np')^2}{np'} + \frac{[n - Y - n(1 - p')]^2}{n(1 - p')} = \frac{(Y - np')^2}{np'q'} = Z^2$$

unde $Z = \frac{Y - np'}{\sqrt{np'q'}}$. Testul χ^2 capătă deci următoarea formă: Respinge $H \iff (T \geq y)$, sau dacă $|Z| \geq \sqrt{y}$. Se știe că dacă H este adevărată Z este practic gaussiană redusă (teorema limită centrală).

Aplicatie numerică: $p' = 0,2$, $n = 100$, $Y = 22$. Valoarea lui Z^2 este $Z^2 = \frac{(22-20)^2}{100 \times 0,2 \times 0,8} = \frac{1}{4}$.

Probabilitatea ca variabila gaussiană redusă să ia o valoare absolută $> \sqrt{\frac{1}{4}} = 0,5$ este 0,616. Se acceptă deci această ipoteză H .

Exemplul 12.7 O v.a. X ia numai valorile 0, 1, 2, 3, 4. Vrem să testăm dacă această v.a. se spune legii binomiale cu $p = \frac{1}{3}$ și $n = 4$ (numărul de probe). S-au făcut pentru aceasta 324 încercări independente care au condus la următoarele rezultate, f_i fiind frecvența absolută a valorii i :

i	0	1	2	3	4
f_i	67	122	94	38	3

Ce concluzie trageți?

Soluție Fie A_i evenimentul: $X=i$; $P(A_i) = C_4^i \cdot \left(\frac{1}{3}\right)^i \cdot \left(\frac{2}{3}\right)^{4-i}$, pentru $0 \leq i \leq 4$. Avem deci $p'_0 = \frac{16}{81}$; $p'_1 = \frac{32}{81}$; $p'_2 = \frac{24}{81}$; $p'_3 = \frac{8}{81}$; $p'_4 = \frac{1}{81}$; $Y_0 = 67$; $Y_1 = 122$; $Y_2 = 94$; $Y_3 = 38$; $Y_4 = 3$. Deoarece Y_4 este mic se grupează A_3 cu A_4 și găsim A_0 , A_1 , A_2 și $B = (A_3 \cup A_4)$ cu p'_0 , p'_1 , p'_2 și $q' = p'_3 + p'_4$. Experiența a condus la Y_0 , Y_1 , Y_2 și $Z = Y_3 + Y_4$. Se obține $\chi^2 = 1,15$. Probabilitatea ca χ^2 cu 4-1=3 grade de libertate să ia o valoare de această mărime este $> 0,1$. Deci ipoteza se acceptă.

Exemplul 12.8 Gazul de esapament al unui motor contine particule solide. Se consideră ipoteza H : numărul X al acestor particule continut într-un volum mic V de gaz se supune unei legi Poisson. Pentru a testa această ipoteză luăm 400 esantioane de același volum V și se găsesc 1872 de particule reprezentate după următorul tabel (n_i reprezintă numărul de esantioane care au continut i particule):

i	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
n_i	0	20	43	53	86	70	54	37	18	10	5	2	2	0	0

Aceste rezultate permit oare să acceptăm ipoteza H ?

Soluție $\bar{X} = \frac{1}{n} \sum in_i = \frac{1872}{400} = 4,68$. Se grupează clasele corespunzătoare lui $i = 0$ și $i = 1$, cât și clasele corespunzătoare lui $i \geq 10$. Se obțin în final 10 clase numerotate de la 1 la 10 care au probabilitățile estimate prin \hat{p}_i :

$\hat{p}_i = e^{-4,68} (1 + 4,68)$, $\hat{p}_i = e^{-4,68} \times \frac{(4,68)^i}{i!}$, pentru $2 \leq i \leq 9$, $\hat{p}_i = \frac{4,68}{i} \hat{p}_{i-1}$, pentru $3 \leq i \leq 9$.

și $\hat{p}_{10} = 1 - \sum_{1 \leq i \leq 9} \hat{p}_i$. Avem următorul tabel:

i	Y_i	\hat{p}_i	$n\hat{p}_i$	$\frac{(Y_i - n\hat{p}_i)^2}{n\hat{p}_i}$
0 sau 1	20	0,0527	21,1	0,0552
2	43	0,1016	40,6	0,1372
3	53	0,1585	63,4	1,7060
4	86	0,1855	74,2	1,8764
5	70	0,1736	69,4	0,0044
6	54	0,1354	54,1	0,0004
7	37	0,0905	36,2	0,0176
8	18	0,0529	21,1	0,4720
9	10	0,0275	11,0	0,0908
≥ 10	9	0,0218	8,7	0,0152
				4,3772

Deci suma termenilor de pe ultima coloană este chiar $\chi^2 \approx 4,38$. Probabilitatea pentru ca o variabilă χ^2 cu $r - s - 1 = 10 - 1 - 1 = 8$ grade de libertate să fie $\geq 4,38$ este $> 0,1$. Prin urmare H se va accepta.

Exemplul 12.9 Vrem să examinăm dacă aptitudinile manuale ale unui individ sunt independente de vedere. Pentru aceasta se definesc 2 caractere X și Y : X ia valorile 1, 2 sau 3 care corespund faptului că individul este mai abil cu mâna stângă, la fel de abil cu ambele mâini, sau mai abil cu mâna dreaptă. Y ia valorile 1, 2 sau 3 după cum individul vede mai bine cu ochiul stâng, cu ambii ochi sau vede mai bine cu ochiul drept. Facem deci ipoteza că X și Y sunt v.a. independente. În tabelul următor avem rezultatele observațiilor făcute asupra a 413

persoane. De exemplu, am găsit 20 persoane cu $X=2$ și $Y=3$, etc.

	Y	1	2	3
X				
1		34	62	28
2		27	28	20
3		57	105	52

.

Solutie Cu notatiile din *Aplicatia 2* din această lectie avem

$$\tilde{N}_{1\bullet} = 34 + 62 + 28 = 124$$

$$\tilde{N}_{2\bullet} = 75, \tilde{N}_{3\bullet} = 214$$

$$\tilde{N}_{\bullet 1} = 34 + 27 + 57 = 118$$

$$\tilde{N}_{\bullet 2} = 196, \tilde{N}_{\bullet 3} = 52 \text{ in } = 34 + 62 + \dots + 52 = 413.$$

Pentru numerele $\frac{1}{n}\tilde{N}_{i\bullet}\tilde{N}_{\bullet j}$ avem tabelul:

	j	1	2	3
i				
1		35	59	30
2		21	35	18
3		61	101	52

și

$\chi^2 = \frac{(34-35)^2}{35} + \dots + \frac{(52-52)^2}{52} = 3,5$. Cum v.a. χ^2 are $(a-1)(b-1) = 2 \cdot 2 = 4$ grade de libertate, ipoteza de independență se acceptă (vezi Tabelul ?).

Exemplul 12.10 *Inteligența unui copil se clasează în 6 nivele: de la A—foarte slabă, până la F—foarte bună. S-au clasat 1725 copii la întâmplare alesi din 8 școli numerotate de la 1 la 8 și s-au obținut rezultate următoare:*

	nivelul	A	B	C	D	E	F
scoala							
1		6	18	36	43	39	4
2		14	25	52	87	54	13
3		—	—	1	8	37	—
4		14	19	60	94	73	12
5		33	69	132	187	85	14
6		45	50	69	72	66	15
7		—	6	20	8	9	1
8		18	32	37	36	12	—

Aceste rezultate ne permit oare să acceptăm ipoteza după care repartitia nivelelor de inteligența este aceeași în oricare dintre școli?

Soluție Este vorba despre un test de omogenitate (vezi *Aplicatia 3*). Aici $t = 8$, $r = 6$, numerele Y_{ij} se află în tabel la intersecția liniei i cu coloana j , de exemplu $Y_{23} = 52$. Avem $n_1 = 6 + 18 + 36 + 43 + 39 + 4 = 146$, $n_2 = 243$, $n_3 = 46$, $n_4 = 272$, $n_5 = 520$, $n_6 = 317$, $n_7 = 44$, $n_8 = 135$. La fel $\tilde{Y}_{\bullet 1} = 6 + 14 + 14 + 33 + 44 + 18 = 130$, $\tilde{Y}_{\bullet 2} = 219$, $\tilde{Y}_{\bullet 3} = 407$, $\tilde{Y}_{\bullet 4} = 535$, $\tilde{Y}_{\bullet 5} = 375$, $\tilde{Y}_{\bullet 6} = 59$. Calculăm acum numerele $\frac{Y_{11}^2}{n_1 \tilde{Y}_{\bullet 1}} = 0,001897, \dots, \frac{Y_{85}^2}{n_8 \tilde{Y}_{\bullet 5}} = 0,002844$ și suma lor $\sum_{i,j} \frac{Y_{ij}^2}{n_i \tilde{Y}_{\bullet j}} = 1,0784$. De aici găsim că $\chi^2 = 1725 \times (1,0784 - 1) = 135$. Variabila aleatoare χ^2 cu $(t-1)(r-1) = 35$ grade de libertate are o probabilitate mai mică decât 0,001 ca să fie ≥ 135 . Deci ipoteza va trebui categoric respinsă.

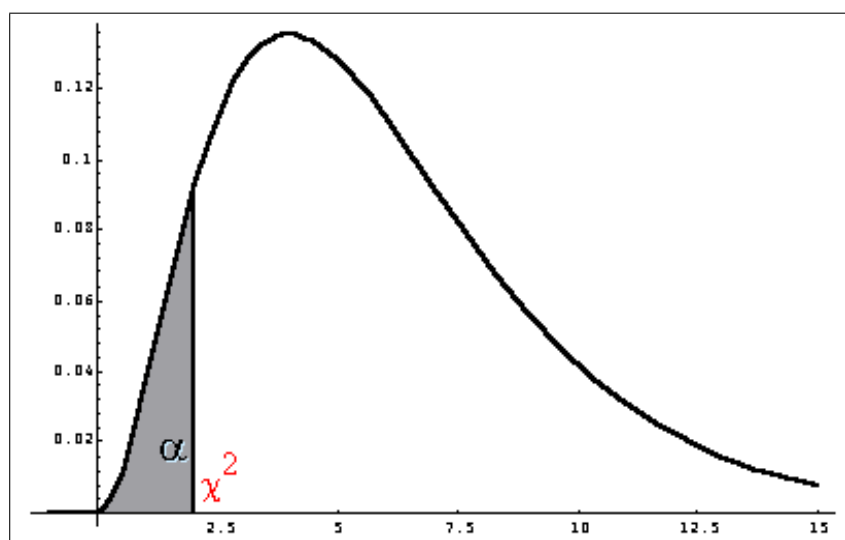
12.2 Rezumat

Fie X o v.a. ce guvernează populația P . Nu cunoaștem legea statistică a v.a. X și vrem să testăm dacă această lege este o lege cunoscută Q (dată). Pentru aceasta considerăm un esantion de volum n ($n \geq 25$): x_1, x_2, \dots, x_n . Grupăm aceste date în intervale J_1, J_2, \dots, J_r astfel încât $J_i \cap J_j = \emptyset$, pentru $i \neq j$ și $J_1 \cup J_2 \cup \dots \cup J_r = \mathbb{R}$ să ia valoarea în intervalul J_i . Este clar că nu cunoaștem aceste p_i -uri. Stim doar că $\sum p_i = 1$. Fie p'_i probabilitatea ca v.a. X să ia valoarea în intervalul J_i dacă am presupune că X are legea Q . Aceste numere se pot calcula folosind tabelele distribuției cunoscute Q . Notăm cu Y_i numărul acelor x_j -uri care se află în intervalul J_i . Y_i este de fapt frecvența absolută empirică de selecție pe intervalul J_i . Frecvența absolută teoretică (potrivit legii Q) pe intervalul J_i este egală cu np'_i . Ipoteza pe care o facem este următoarea: $H: p_i = p'_i, (\forall) i = \overline{1, r}$.

Construim acum statistica Helmert-Pearson: $T = \sum_{i=1}^r \frac{(Y_i - np'_i)^2}{np'_i}$. Ea măsoară "deviația" legii reale a v.a. X de la legea presupusă Q . Pentru n mare T tinde să aibă distribuția χ^2 (Pearson) cu $r - 1$ grade de libertate dacă cumva X ar avea într-adevăr legea de distribuție Q .

- Calculăm numărul T în cazul nostru și îl notăm cu χ^* .

- Fie acum $\alpha \in (0, 1)$ (de obicei α se ia mic $\alpha = 0,05; 0,01; 0,001$), un număr real dat pe care îl vom numi *prag de semnificație*.
- Fie χ_α^2 cuantila de ordin α pentru v.a. χ^2 , adică aceea valoare pentru care $P(\chi^2 < \chi_\alpha^2) = \alpha$ (vezi figura)



- Testul χ^2 funcționează astfel:
 - Dacă valoarea calculată din selecție $\chi^* > \chi_\alpha^2$ vom accepta ipoteza H cu pragul de semnificație α , adică este foarte probabil ca ipoteza să fie adevărată.
 - Dacă valoarea calculată $\chi^* \leq \chi_\alpha^2$ vom respinge ipoteza H cu pragul de semnificație α .

Aceasta este testul χ^2 simplu (Cazul I)

- Dacă însă legea Q are s parametri de estimat tot din selecție, atunci va trebui să micșorăm numărul gradelor de libertate cu s , adică să avem $r - s - 1$ grade de libertate, deoarece apar $s + 1$ legături: $\sum p_i = 1$ și cele s legături de estimare.

12.3 Exerciții rezolvate

1. Se arunca 4 monede simultan de 160 de ori. Se observa de fiecare data de câte ori a apărut "capul".

x = de câte ori poate să apară capul, când aruncăm o dată cele 4 monede	0	1	2	3	4
f = de câte ori s-a realizat x în cele 160 de aruncări x în cele 160 de aruncări	5	35	67	41	12

Să se testeze ipoteza H_0 : monedele nu sunt falsificate, cu pragul de 5%.

Soluție Fie X v.a. care ia valorile x = de câte ori apare "capul" când aruncăm o dată cele 4 monede (sau de 4 ori aceeași monedă). Avem ca:

$$P(X = x) = C_4^x \cdot \left(\frac{1}{2}\right)^x \cdot \left(\frac{1}{2}\right)^{4-x} = \left(\frac{1}{2}\right)^4 \cdot C_4^x.$$

Prin urmare frecventa absoluta sperata (daca presupunem H_0 adevarata) de aruncari pentru x este $160 \cdot P(X = x)$. Obtinem deci urmatorul tabel:

x	0	1	2	3	4
$P(X = x)$	1/16	4/16	6/16	4/16	1/16
E_i = fr. abs. sperata cand $x = i$ sperata cand $x = i$	10	40	60	40	10
O_i = fr. abs. observata cand $x = i$	5	35	67	41	12
$\frac{(O_i - E_i)^2}{E_i}$	2,5	0,625	0,817	0,025	0,4

Folosim testul χ^2 cu $\nu = 5$ clase-1 restrictie ($\sum O_i = 160$) = 4 grade de libertate.

Avem deci $\chi_{calc}^2 = \sum_{i=1}^5 \frac{(O_i - E_i)^2}{E_i} = 4,365 < \chi_{0,05}^2(4) = 9,49$. Prin urmare acceptam H_0 cu pragul de 95%.

2. La o fabrica de caramizi se aleg 500 de pachete de cate 5 caramizi la intervale regulate de-a lungul unei saptamani si se numara de fiecare data cate caramizi defecte sunt in fiecare pachet. S-au obtinut urmatoarele rezultate:

x = nr. de caramizi defecte intr-un singur pachet	0	1	2	3	4	5
f = nr. de pachete care au avut x caramizi defecte	170	180	120	20	8	2

Testati cu ajutorul testului χ^2 cu nivelul de semnificatie 5% daca numarul caramizilor defecte urmeaza legea binomiala.

Solutie Deoarece in legea binomiala $\text{Bin}(n = 5; p)$ nu stim pe p , trebuie sa-l estimam din esantion:

$$np = \bar{x} = \frac{\sum f \cdot x}{\sum f} = 1,044, \text{ deci } p = 0,2088.$$

Am presupus ca H_0 : "legea este binomiala" este adevarata. Atunci v.a. X = nr. caramizilor defecte intr-un esantion de 5 caramizi va fi binomiala cu $n = 5$ si $p = 0,2088 \equiv$ probabilitatea ca o caramida luata la intamplare sa fie defecta. Avem ca frecventa absoluta sperata de " i " caramizi defecte este $E_i = 500 \cdot P(X = i) = 500 \cdot C_5^i (0,2088)^i (0,7912)^{5-i}$, unde $0,7912 = 1 - 0,2088$, iar $i = 0, 1, 2, 3, 4, 5$. Cum $E_4 = 4$ si $E_5 = 0$, clasele E_3, E_4 , si E_5 le cumulam intr-o singura clasa renotata cu $E_3 (x \geq 3)$. La fel, cumulam pe $O_3 = 20, O_4 = 8$ si $O_5 = 2$ intr-o noua clasa $O_3 (x \geq 3)$. Obtinem deci urmatorul tabel:

x	0	1	2	$x \geq 3$
E_i = fr. abs. sperata cand $x = i$, daca acceptam legea binomiala	155	205	108	32
O_i = fr. abs. observata cand $x = i$, data in datele problemei	170	180	120	30
$\frac{(O_i - E_i)^2}{E_i}$	1,452	3,049	1,333	0,125

Folosim testul χ^2 cu $\nu = 4$ clase-2 restrictii ($\sum O_i = 500$ si $\bar{x} = 1,044$ este impus prin estimare)= 2 grade de libertate.

Avem deci $\chi^2_{calc} = 5,959 = \sum_{i=1}^4 \frac{(O_i - E_i)^2}{E_i} < \chi^2_{0,05}(2) = 5,99$. Acceptam deci ipoteza H_0 (repartitia este binomiala) cu pragul de semnificatie 5%. Avem o oarecare neincredere deoarece 5,959 este prea aproape de valoarea critica 5,99. Se indica sa se repete experienta pentru mai multa siguranta.

3. In 100 de meciuri o echipa de fotbal a inregistrat goluri dupa cum urmeaza:

x =nr. de goluri inregistrate intr-un meci	0	1	2	3	4	5	6	7
f =nr. de meciuri in care echipa a inregistrat x goluri	14	18	29	18	10	7	3	1

Testati cu χ^2 cu 5% daca golurile inregistrate se repartizeaza dupa o lege Poisson.

Solutie Deoarece nu cunoastem pe λ in legea lui Poisson $Po(\lambda)$, il vom estima din media de selectie:

$$\bar{x} = \frac{\sum f \cdot x}{\sum f} = \frac{230}{100} = 2,3, \text{ deci } \lambda \approx \bar{x} = 2,3.$$

$$P(X = x) = \frac{e^{-2,3} \cdot (2,3)^x}{x!}, \quad x = 0, 1, 2, \dots$$

Calculam $E_i = 100 \cdot P(X = i) \equiv$ frecventa absoluta sperata (teoretica) a "i" goluri inregistrate de echipa in cele 100 de meciuri. Aici $P(X = i)$ este probabilitatea ca echipa sa inregistreze "i" goluri intr-un singur meci. Folosim testul χ^2 cu $\nu = 6$ clase - 2 restrictii ($\sum O_i = 100$ si $\bar{x} = 2,3$ este impus)=4 grade de libertate. Teoretic avem 9 clase: $E_0, E_1, \dots, E_7, E_8$ ($i \geq 8$). Cum $E_5 = 5,4$; $E_6 = 2,1$; $E_7 = 0,7$, $E_8 = 0,2$, cumulam clasele E_5, E_6, E_7 si E_8 intr-una singura, desemnata tot cu $E_5 = 8,4$. Vom obtine tabelul:

x	0	1	2	3	4	5
E_i	10,0	23,1	26,5	20,3	11,7	8,4
O_i	14	18	29	18	10	11
$\frac{(O_i - E_i)^2}{E_i}$	1,6	1,126	0,236	0,261	0,247	0,805

Folosim deci testul χ^2 cu 4 grade de libertate. Cum $\chi^2_{calc} = \sum_{i=0}^5 \frac{(O_i - E_i)^2}{E_i} = 4,275 < \chi^2_{0,05}(4) = 9,49$ vom accepta ipoteza H_0 : "distributia golurilor urmeaza o lege Poisson" cu nivelul de incredere de 95%.

12.4 Exerciții

1. Se testează greutatea unui grup de 50 de bărbați și se obțin valorile (în Kg):

66	78	82	75	94	77	69	74	68	60
96	78	89	61	75	95	60	79	83	71
79	62	67	97	78	85	76	65	71	75
86	84	75	81	68	63	62	75	76	77
73	65	88	87	60	62	71	78	85	72

Folositi testul χ^2 pentru a decide cu pragul de $\alpha=0,90$ dacă populația este normală cu media 75 și dispersia 625. Faceti același lucru dar cu media și dispersia estimate din selecție. Împărțiți esanșionul în zece intervale.

2. Se urmăresc accidentele mortale pe o porțiune dintr-un drum național timp de 100 săptămâni. Timp de 45 de săptămâni nu a fost nici un accident mortal. În 29 de săptămâni s-a produs un accident, în 17 săptămâni 2 accidente, iar în 9 săptămâni au avut loc 3 accidente. Folositi testele χ^2 și K-S pentru a vedea, cu pragul $\alpha=0,90$, dacă distribuția accidentelor urmează modelul Poisson sau nu. Parametrul se estimează din selecție.

3. Se consideră o prismă care are ca baze două triunghiuri echilaterale B_1 și B_2 și ca fete laterale A_1 , A_2 și A_3 . Se aruncă prisma de 500 de ori și se constată că prisma a căzut de :

$Y_1 = 111$ ori pe fata A_1

$Y_2 = 113$ ori pe fata A_2

$Y_3 = 118$ ori pe fata A_3

$Z_1 = 81$ ori pe fata B_1

$Z_2 = 77$ ori pe fata B_2 . Testati ipoteza după care cele 3 fete laterale și cele 2 baze au aceeași probabilitate $\frac{1}{5}$ (pragul $\varepsilon = 0,05$ și $\varepsilon = 0,01$).

Indicatie. Avem Cazul I: $r = 5$ și $p'_1 = p'_2 = \dots = p'_5 = \frac{1}{5}$. Găsim $\chi^2 = 15,04$. Analizând tabelul vedem că trebuie să respingem ipoteza pentru cele două praguri.

4. Vrem să testăm ipoteza H după care o anumită v.a. X este gaussiană cu media 1,1 și dispersia 0,2. S-au făcut 1000 de probe independente care au condus la rezultatele următoare:

0,6	0,7	0,8	0,9	1	1,1	1,2	1,3	1,4	1,5
26	51	107	168	200	193	138	80	29	8

adică X a luat de 26 de ori o valoare mai mică decât 0,6 și de 51 de ori o valoare cuprinsă în intervalul $[0,6; 0,7)$, etc. Testati ipoteza H cu pragurile $\varepsilon = 0,05$, $\varepsilon = 0,01$ și $\varepsilon = 0,001$.

Indicatie Se partitionează \mathbb{R} după cum indică problema. Se găsește $\chi^2 = 13,97$. Folosim tabelul distribuției χ^2 cu $10 - 1 = 9$ grade de libertate. Ipoteza trebuie respinsă cu pragul 0,05 și 0,01, dar trebuie să o acceptăm cu pragul 0,001.

5. Se testează ipoteza după care culoarea ochilor este independentă de culoarea părului. Pentru aceasta se introduc 2 caractere X și Y :

X ia valorile 1, 2, 3, 4 după cum ochii sunt albastri, gri sau bruni.

Y ia valorile 1, 2, 3, 4 după cum părul este blond, brun, negru sau roscat. Se testează un număr de persoane și rezultatele le găsim în tabelul următor:

	Y	1	2	3	4
X					
1		1768	807	189	47
2		946	1387	746	53
3		1125	438	288	16

.
Indicatie Aici $(a - 1)(b - 1) = 2 \times 3 = 6$ și $\chi^2 = 1075$. Chiar cu pragul de $\varepsilon = 0,001$ ipoteza nu este acceptabilă.

6. Un articol de calitățile A, B sau C poate fi fabricat după 2 metode numerotate cu 1 și 2. Se examinează 100 de articole și se găsește tabelul:

	A	B	C
1	20	19	11
2	12	31	7

. Putem accepta oare ipoteza că, calitatea articolului nu depinde de modul său de fabricare (se ia $\varepsilon = 0,1$ și $\varepsilon = 0,01$)?

Indicatie. Este vorba de un test de omogenitate cu $r = 3$ și $t = 2$. Avem că $\chi^2 = 5,77$ pentru o variabilă aleatoare χ^2 cu 2 grade de libertate. Cu pragul de 0,1 trebuie să respingem ipoteza, dar cu pragul de 0,01 trebuie să o acceptăm.

Lecția 13

Alte teste neparametrice

13.1 Testul de concordantă Kolmogorov-Smirnov

Acest test este din multe puncte de vedere "mai bun" decât testul χ^2 . El se aplică bine și pentru esantioanele mici ($n \leq 25$). Dacă pentru testul χ^2 trebuia să grupăm datele, pentru testul K-S nu este nevoie decât să calculăm funcția de repartiție empirică asociată selecției efectuate. El poate fi utilizat bine și pentru a compara două distribuții.

Presupunem că populația P are ca funcție de repartiție teoretică funcția $F_T(x)$. Efectuăm un sondaj de volum n : $\{x_1, \dots, x_n\}$ și notăm cu $F_S(x)$ funcția de repartiție empirică asociată acestui sondaj (vezi Lecția 8). Pentru a măsura deviația funcției $F_T(x)$ de la funcția $F_S(x)$ se introduce statistica lui Kolmogorov: $D = \max_x |F_S(x) - F_T(x)|$. Dacă populația P are într-adevăr repartiția $F_T(x)$ atunci se cunoaște distribuția v.a. D (vezi Tabelul XIX). Pe prima coloană în acest tabel avem volumul de selecție. Pe prima linie orizontală avem 5 valori ale pragului de semnificație: $\alpha = 0,80; 0,85; 0,90; 0,95$ și $0,99$.

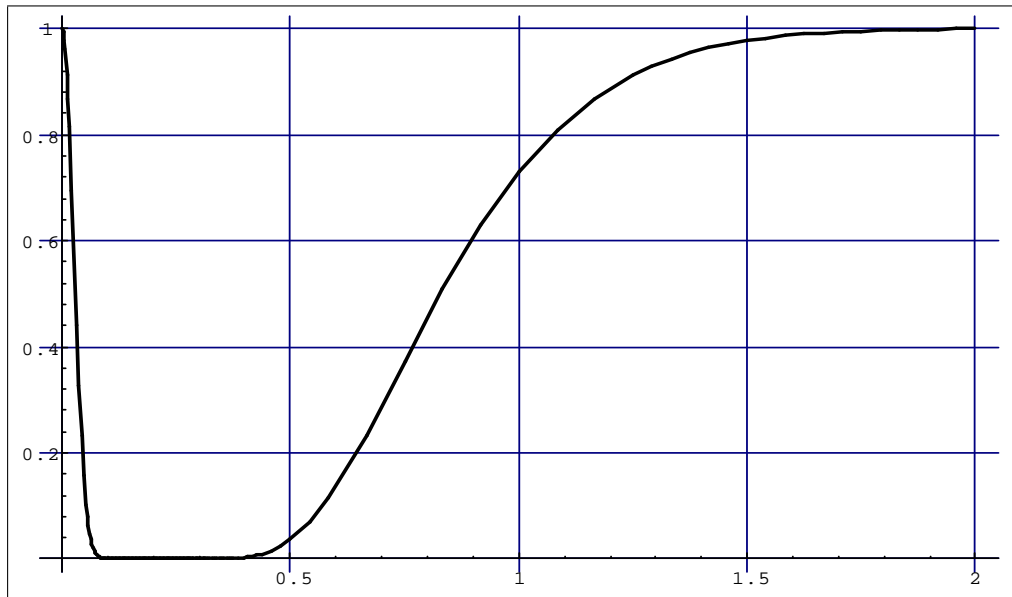
Tabelul XIX este construit pe baza teoremei lui A. N. Kolmogorov:

Teorema 13.1 *Fie X o v.a., $F_T(x)$ funcția ei de repartiție considerată continuă și $F_n(x)$ o funcție empirică de repartiție asociată unei selecții de volum $n: \{x_1, \dots, x_n\}$, dintr-o populație cu v.a. X . Atunci avem relația*

$$\lim_{n \rightarrow \infty} P \left(D \leq \frac{\lambda}{\sqrt{n}} \right) = K(\lambda) = \sum_{k=-\infty}^{\infty} (-1)^k \cdot \exp(-2k^2 \lambda^2) \quad (13.1)$$

pentru orice $\lambda > 0$.

Graficul funcției $K(\lambda)$ apare în figura următoare:

Graficul funcției $K(\lambda)$

Exemplul 13.2 Se face sondajul $\{-2, -1, -1, 0, 0, 1, 2, 2\}$ dintr-o populație P . Să se testeze cu testul $K-S$ dacă populația este normală cu media 0 și dispersia 2 (date și nu estimate!) cu pragul de semnificație $\alpha = 0,80$.

Soluție Trebuie să construim funcția de repartiție empirică, $F_S(x)$ — notată de noi cu $F_n^*(x)$ în Lectia 8.

$$F_S(x) = \begin{cases} 0, & x \leq -2 \\ \frac{1}{8}, & x \in (-2, -1] \\ \frac{3}{8}, & x \in (-1, 0] \\ \frac{5}{8}, & x \in (0, 1] \\ \frac{7}{8}, & x \in (1, 2] \\ 1, & x > 2 \end{cases}$$

Deoarece $P(X < x) = F_T(x)$ și cum v.a. $\frac{X}{\sigma} = \frac{X}{\sqrt{2}}$ este normală redusă (de tipul $N(0,1)$) avem că $P(X < x) = P\left(\frac{X}{\sigma} < \frac{x}{\sigma}\right) = \Phi\left(\frac{x}{\sigma}\right)$. Deci $F_T(x) = \Phi\left(\frac{x}{\sigma}\right)$, unde Φ este funcția lui Laplace tabelată în Tabelul I. Prin urmare va trebui să calculăm $\Phi\left(\frac{x}{\sqrt{2}}\right)$ în punctele $x = -2, -1, 0, 1, 2$. Cum $\Phi(-y) = 1 - \Phi(y)$, rămâne numai (deoarece $\Phi(0) = 0,5$) să calculăm din Tabelul I $\Phi\left(\frac{2}{\sqrt{2}}\right) = \Phi(\sqrt{2}) = \Phi(1,41) = 0,92$ și $\Phi\left(\frac{1}{\sqrt{2}}\right) = \Phi\left(\frac{\sqrt{2}}{2}\right) = \Phi(0,70) = 0,75$. Prin urmare $\Phi\left(-\frac{2}{\sqrt{2}}\right) = 1 - 0,92 = 0,08$ și $\Phi\left(-\frac{1}{\sqrt{2}}\right) = 1 - 0,75 = 0,25$. Calculăm acum diferențele

dintre $F_S(x)$ și $F_T(x)$ pe fiecare interval ce apare în definiția funcției $F_S(x)$. Găsim

x	$-\infty$	-2	-1	0	1	2	∞
$F_S(x)$	0	0	$\frac{1}{8}$	$\frac{3}{8}$	$\frac{5}{8}$	$\frac{6}{8}$	1
$F_T(x)$	0	$0,08$	$0,25$	$0,5$	$0,75$	$0,92$	1
$F_S(x) - F_T(x)$	0	$-0,08$	$-0,125$	$-0,125$	$-0,125$	$-0,170$	0

Prin urmare $D = \max |F_S(x) - F_T(x)| = 0,170$. Cuantila pentru $\alpha = 0,80$ și $n = 8$ (din Tabelul XIX) este $D_\alpha = 0,358$. Cum $P(D < 0,358) = 0,80$ rezultă că putem accepta ipoteza că populația P este normală cu pragul de semnificație 80% (deoarece $0,170 < 0,358$).

Observația 13.3 Acest test se poate aplica când avem de comparat două distribuții:

$$D = \max |F_{S_1}(x) - F_{S_2}(x)|$$

etc., unde cele două funcții empirice de repartiție care apar corespund celor două distribuții.

Exercițiul 13.4 Folositi testul $K-S$ pentru exemplul 12.1

13.2 Testul lungimilor (secventelor)

Fie X și Y două variabile aleatoare. Vrem să testăm următoarea ipoteză:

H : X și Y au aceeași lege de distribuție.

Pentru aceasta considerăm un esantion de volum m al v.a. X : (X_1, \dots, X_m) și un esantion de volum n al v.a. Y : (Y_1, \dots, Y_n) . Considerăm acum sirul de $m + n$ variabile aleatoare $(X_1, \dots, X_m, Y_1, \dots, Y_n)$. Presupunem că ele sunt independente. Dacă X și Y ar avea aceeași lege de distribuție am putea "amesteca" oricum aceste variabile și raționamentele noastre nu s-ar schimba la trecerea de la (X_1, \dots, X_n) la (Y_1, \dots, Y_m) .

Introducem următoarea v.a. R . Facem câte o selecție de volum m din X : x_1, \dots, x_m și de volum n din Y : y_1, \dots, y_n considerăm sirul de numere $x_1, \dots, x_m, y_1, \dots, y_n$. Asezăm acum aceste numere în ordinea crescătoare și nu ne interesează decât faptul că ele sunt din X sau din Y . Vom nota o astfel de situație sub forma: $\omega = \text{XXYXYYYXXXXYY}$. Aici $m =$ numărul X -lor, adică $m = 7$, $n =$ numărul Y -lor, adică $n = 6$. În sirul ω cel mai mic număr din sirul $x_1, \dots, x_m, y_1, \dots, y_n$, este din (x_1, \dots, x_m) , următorul în ordine crescătoare este tot din acest sir, al treilea în ordine crescătoare este din (y_1, \dots, y_n) , etc. Cel mai mare număr din ω este din (y_1, \dots, y_n) . Dacă cumva $x_{i_1} = x_{i_2}$ le asezăm unul după altul în orice ordine, etc.

Vom numi *secvență (lungime)* în sirul ω orice subsir format numai din X sau numai din Y . De exemplu în ω avem următoarele secvențe:

$$\begin{array}{cccccc} \text{XX} & \text{Y} & \text{X} & \text{YYY} & \text{XXXX} & \text{YY} \\ \hline \text{S1} & \text{S2} & \text{S3} & \text{S4} & \text{S5} & \text{S6} \end{array}$$

Prima secventa este S1: XX

A doua secventa este S2: Y

A treia secventa este S3: X

\vdots

A sasea secventa este S6: YY.

Pentru selectia noastră v.a. R ia valoarea 6, adică numărul secventelor de X-și sau de Y-ci care apar după rearanjarea în ordine crescătoare a esantionului $(x_1, \dots, x_m, y_1, \dots, y_n)$.

Dacă ipoteza H este adevărată, functia de repartitie a v.a. R , $F(x) = P(R \leq x)$ se calculează tinând seama de următoarele observatii de natură combinatorică:

Dacă de exemplu $m \leq n$,

$$P(R = 2s) = \frac{2}{C_{m+n}^m} \cdot C_{m-1}^{s-1} \cdot C_{m-1}^{s-1}, \text{ pentru } 1 \leq s \leq m;$$

$$P(R = 2s + 1) = \frac{1}{C_{m+n}^m} \cdot [C_{m-1}^s \cdot C_{m-1}^{s-1} + C_{n-1}^s \cdot C_{n-1}^{s-1}], \text{ pentru } 1 \leq s \leq m;$$

$$P(R = 2m + 1) = \frac{1}{C_{m+n}^m} \cdot C_{n-1}^m, \text{ dacă } m < n;$$

$$P(R = r) = 0$$

în toate celelalte cazuri rămase.

Se demonstrează că statistica v.a. R este asimptotică (m, n mari) gaussiană cu

$$M(R) = 1 + \frac{2mn}{m+n}, D(R) = \frac{2mn(2mn - m - n)}{(m+n)^2(m+n-1)} \quad (13.2)$$

Să notăm cu

$$T = \frac{R - M(R)}{\sqrt{D(R)}} \quad (13.3)$$

Variabila T tinde către o v.a. gaussiană redusă.

Un test de prag ϵ se construiește astfel:

— Fie r cel mai mare întreg x astfel încât $F(x) \leq \epsilon$.

— Testul va fi: Respinge $H \iff R^* \leq r$, unde R^* este valoarea efectivă a v.a. R obținută din sondaje (= numărul secventelor). Pragul este ϵ și acest test este cel mai puternic pentru acest prag.

Desigur că am putea construi și alte teste plecând de la formula (13.3) în analogie cu testele de semnificatie sau ca testul χ^2 . Pentru m și n mici se calculează direct $F(x)$ după formulele (13.1), iar pentru m și n mari se folosește (13.3).

Exemplul 13.5 Se cântăresc 10 mere de două calități A și B și se găsesc următoarele rezultate (în grame):

$$\begin{array}{l} A: 192 \quad 197 \quad 207 \quad 182 \quad 191 \\ B: 212 \quad 201 \quad 209 \quad 214 \quad 203 \end{array}$$

Ne permit oare aceste rezultate să respingem ipoteza H: greutatea unui măr urmează aceeași lege de distribuție indiferent de calitatea lui A sau B (cu pragul $\epsilon=0,05$)? Calculați probabilitatea erorii de primul tip pentru acest test. Utilizați și aproximarea gaussiană și comparați rezultatele.

Soluție Aici avem $m = n = 5$; $C_{m+n}^m = 252$. Să notăm $N(x) = C_{m+n}^m \cdot P(R=x) = 252 \cdot P(R=x)$. Avem $N(2s) = 2 \cdot (C_4^{s-1})^2$, pentru $1 \leq s \leq 5$,
 $N(2s+1) = 2 \cdot C_4^{s-1} \cdot C_4^{s-1}$, pentru $1 \leq s \leq 4$ și $N(x) = 0$ în celelalte cazuri. Deducem de aici legea de distribuție pentru v.a. R:

$$\begin{array}{cccccccccc} x & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 \\ 252 \cdot P(R=x) & 2 & 8 & 32 & 48 & 72 & 48 & 32 & 8 & 2 \end{array}$$

Căutăm cel mai mare x cu $F(x) \leq 0,05$, sau $252 F(x) \leq 12,60$. Cum avem $252 F(x) = \sum_{i \leq x} N(i)$, găsim că $252 F(2) = 2$; $252 F(3) = 10$; $252 F(4) = 42$.

Cel mai mare x este deci 3 și testul devine: Respinge $H \iff R^* < 3$. Probabilitatea erorii de primul tip pentru acest test este $P(R \leq 3) = F(3) = \frac{10}{252} = 0,0397$. Să utilizăm acum aproximarea gaussiană estimând $M(R)$ și $D(R)$ direct din selecție.

$$M(R) = 1 + \frac{2 \cdot 5 \cdot 5}{5 + 5} = 6.$$

$$D(R) = \frac{2 \cdot 5 \cdot 5 \cdot 40}{10 \cdot 9} = 2,222.$$

Avem că

$$P(R \leq 3,5) = P\left(\frac{R - 6}{\sqrt{2,222}} \leq \frac{-2,5}{\sqrt{2,222}}\right) = 0,047$$

Să observăm că rezultatele sunt comparabile chiar dacă m și n și sunt mici.

În exemplul nostru $\omega = \text{XXXXYYXYYY}$, deoarece selecția ordonată este:

$$182 \quad 191 \quad 192 \quad 197 \quad 201 \quad 203 \quad 207 \quad 209 \quad 212 \quad 214$$

Avem deci 4 secvențe. Cum $R^* = 4 > 3$ va trebui să acceptăm ipoteza H cu pragul $\epsilon = 0,05$.

13.3 Testul lui Wilcoxon I (cazul observatiilor necuplate)

Definiția 13.6 Fie X și Y două v.a. Variabila aleatoare Y se zice stochastic superioară v.a. X dacă $(\forall) z \in \mathbf{R}$ avem că $P(Y \leq z) \leq P(X \leq z)$, cu inegalitate strictă cel puțin pentru un $z = z_1$. Presupunem că avem alternativele:

- sau X și Y au aceeași lege
- sau Y este stochastic superioară lui X .

Fie ipoteza H : X și Y au aceeași lege. Ca și la testul secvențelor construim sirul $\omega = XXYYX\dots$, de exemplu $\omega = XYXXYYXYXY$. Statistica lui Wilcoxon T este v.a. care are ca valoare $T =$ suma numerelor care arată locurile pe care le ocupă X în sirul ω . Aici avem $T = 1 + 3 + 4 + 7 = 15$. Testul este de forma: Respinge $H \iff T \leq t$, unde t este valoarea lui T din selecție. Dacă H este verificată, $P(T \leq t) = \frac{1}{C_{m+n}^m}$ înmulțit cu numărul probelor favorabile relației $T \leq t$ (pentru m și n mici se poate calcula direct $F(t)$). Aici lucrăm cu selecțiile $X_1, \dots, X_m, Y_1, \dots, Y_n$. Dacă m și n sunt mari, T este aproape gaussiană cu:

$$M(T) = \frac{m(m+n+1)}{2} \quad D(T) = \frac{mn(m+n+1)}{12} \quad (13.4)$$

Exemplul 13.7 Se încearcă un tratament medical nou asupra unui grup de persoane de aceeași vârstă, bolnave grav cu o boală de tip cardiac. S-a notat timpul (în ani) după care aceste persoane tratate mai trăiau încă, cât și timpul după care alte persoane de aceeași vârstă, bolnave de aceeași boală, dar netratate, mai trăiau. S-au obținut următoarele rezultate:

Tratate: 1, 2 6, 3 6, 5 7, 8 11, 2 15, 6
 Netratate: 0, 4 3, 5 4, 8 6, 7

Testati ipoteza H potrivit căreia tratamentul nu prelungeste viața unui bolnav (prag 0,05).

Soluție Aplicăm testul lui Wilcoxon. X este v.a. care exprimă durata de viață a unui bolnav netratat și Y a unui bolnav tratat. Este clar că v.a. Y este stochastic superioară v.a. X . Ordonăm crescător sirul $\{x_1, \dots, x_m, y_1, \dots, y_n\}$ și găsim $\omega = XYXXYYXYXY$. Aici $T = 1 + 3 + 4 + 7 = 15$. Avem $m = 4$ și $n = 6$. Să numărăm cazurile favorabile condiției $T \leq 15$. Un caz îl avem mai sus: (1, 3, 4, 7), adică am așezat într-un sir pozițiile lui X în ω . Cazurile favorabile vor fi: (1, 2, 3, 4), (1, 2, 3, 5), (1, 2, 3, 6), (1, 2, 3, 7), (1, 2, 3, 8), (1, 2, 3, 9), (1, 2, 4, 5), (1, 2, 4, 6), (1, 2, 4, 7), (1, 2, 4, 8), (1, 2, 5, 6), (1, 2, 5, 7), (1, 3, 4, 5), (1, 3, 4, 6), (1, 3, 4, 7) și (1, 3, 5, 6). Sunt în total 16 posibilități favorabile. Numărul tuturor posibilităților este $C_{10}^4 = 210$ (= numărul modurilor în care putem așeza patru litere X într-un sir de 10 litere X și Y). Dacă H ar fi adevărată ar trebui să avem $P(T \leq 15) = \frac{16}{210} = 0,076$. Prin urmare, cum $0,076 > 0,05$ trebuie să acceptăm ipoteza H cu pragul 0,05 și trebuie să o respingem cu pragul 0,1.

13.4 Testul semnelor

Presupunem că avem două v.a. X și Y definite pe aceeași categorie de probe, astfel încât $P(X=Y)=0$. Vrem să testăm ipoteza: $H: P(Y>X) \geq P(X>Y)$ contra ipotezei $K: P(Y>X) < P(X>Y)$.

Pentru aceasta se fac n observatii independente $(X_1, Y_1), \dots, (X_n, Y_n)$ ale v.a. (X, Y) . Se notează cu V numărul acelor cupluri (X_i, Y_i) pentru care $Y_i < X_i$. Testul este de forma:

Respinge $H \iff V \leq v$, unde v este valoarea lui V din sondajul respectiv.

Dacă H este adevărată cea mai mare valoare posibilă pentru $P(V \leq v)$ se obține dacă presupunem că V satisface legea lui Bernoulli $B(n, \frac{1}{2})$.

Exemplul 13.8 *O firmă vrea să testeze un nou ingredient adăugat unei creme antisolare. Se fac testări pe 7 voluntari și pe spatele fiecăruia se aplică cremă antisolară astfel: pe jumătatea superioară se aplică crema veche, iar pe jumătatea inferioară se aplică crema cu ingredientul respectiv. Se expun la soare cei 7 voluntari și se observă măsura în care pielea lor se înnegreste. Se obține tabelul următor:*

Voluntarul nr.	1	2	3	4	5	6	7
crema veche	42	51	31	61	44	55	48
crema nouă	38	53	36	52	33	49	36

Solutie Notăm cu Y v.a. corespunzătoare înnegririi pielii cu vechea cremă și cu X v.a. pentru noua cremă, deoarece se presupune că prin adăogarea ingredientului valorile lui X vor fi în general mai mici decât cele ale lui Y . Avem $n=7$. Aplicăm testul semnelor. Numărul cuplurilor (X, Y) pentru care $Y < X$ este $V=2$. Dacă H este adevărată V are o distribuție $B(7, \frac{1}{2})$ și deci

$$P(V \leq 2) = \left(\frac{1}{2}\right)^7 [C_7^0 + C_7^1 + C_7^2] = \frac{29}{128} = 0,23 > 0,1$$

Este deci foarte probabil să găsim pentru V valori mai mari decât 2. Prin urmare, cu acest test al semnelor (V reprezintă câte semne "+" avem în diferența $X-Y$) ipoteza H trebuie acceptată (cu pragul 0,1).

13.5 Testul lui Wilcoxon II (cazul observatiilor cuplate)

Fie X, Y două v.a. definite pe aceeași categorie de probe. Fie $Z=Y-X$. Vrem să testăm ipoteza $H: Z$ este o variabilă aleatoare simetrică, adică are aceeași lege ca și v.a. $-Z$, contra ipotezei $K: Z$ este stochastic superioară lui $-Z$. Să observăm că cele 2 ipoteze *nu* sunt contrare și deci *nu* acoperă gama de posibilități.

Pentru aceasta facem n observatii independente $(X_1, Y_1), \dots, (X_n, Y_n)$ ale v.a. (X, Y) . Aranjăm apoi în ordinea crescătoare a modulelor lor diferențele $Z=Y-X$ și nu reținem decât semnele lor. Găsim de exemplu $\omega=(- - - + - + + +)$. Notăm cu W suma numerelor ce exprimă poziția semnelor minus. În cazul nostru $W=1+2+3+5=11$. Testul are forma următoare: Respinge $H \iff W \leq w$, unde w este valoarea v.a. W obținută din sondaj.

Dacă H este adevărată, mulțimea Ω a celor 2^n posibilități pentru semnul "–" este înzestrată cu legea uniformă (este unica lege probabilistică pe o mulțime finită de evenimente, care face ca aceste evenimente să fie egal probabile). Dacă n nu este prea mare legea W se obține prin numărarea directă a cazurilor favorabile. Dacă n este mare W este aproximativ gaussiană cu:

$$M(W) = \frac{n(n+1)}{4} \quad iD(W) = \frac{n(n+1)(2n+1)}{24} \quad (13.5)$$

Exemplul 13.9 *Reluăm exemplul 13.6 și vrem să-i aplicăm testul Wilcoxon II ($\epsilon=0,1$).*

• Ordonăm valorile v.a. $Y-X$ în ordinea crescătoare a valorilor lor absolute: $-2, 4, -5, 6, 9, 11, 12$. Găsim sirul de semne $\omega=(- + - + + +)$. Suma indicilor cu semnul minus este $W=1+3=4$. Numărul cazurilor posibile este $2^7=128$. Aici $w=4$. Cazurile favorabile ($W \leq 4$) sunt: $(+ + + + + +)$, $(- + + + + +)$, $(+ - + + + +)$, $(+ + - + + +)$, $(+ + + - + +)$, $(- - + + + +)$, $(- + - + + +)$, adică 7. De aici avem că $P(W \leq 4) = \frac{7}{128} = 0,054$.

Cu pragul 0,1 va trebui să respingem ipoteza H .

Observația 13.10 *Dacă comparăm rezultatul obisnuit cu acela din Exemplul 13.6 aparent găsim o contradicție. Aceasta se explică deoarece în testul semnelor nu ținem seama decât de numărul semnelor minus și nu de poziția lor în sirul ω . În testul lui Wilcoxon se ține seama că aceste semne sunt plasate la început, și nu oriunde în sirul ω . Acest lucru face ca valoarea V să fie în general mult mai mare decât valoarea W . De aici apare clar că rezultatul testului Wilcoxon II sr trebui să fie "mai demn" de crezut de firmă decât rezultatul testului semnelor. Comparatia dintre cele două teste se face de obicei de la caz la caz și se interpretează rezultatul potrivit situației particulare studiate. Evident că aici, pentru firmă, este convenabil testul Wilcoxon II și nu testul semnelor, care nu pare concludent. În plus, în testul Wilcoxon II se presupune "ceva" în plus de la început (K nu este alternativa ipotezei H).*

13.6 Exerciții

1. Se dau două esantioane independente de volum 20 din v.a. X și Y :

X: 147 193 238 225 252 143 178 209
 259 263 226 179 253 262 181 169
 210 233 248 194

Y: 240 254 192 157 168 170 207 222
 201 215 217 243 172 183 197 241
 182 163 173 167

Cu fiecare din pragurile $\epsilon=0,05$; $\epsilon=0,1$, testati ipoteza H : X și Y au aceeași lege,

1) cu testul lungimilor;

2) cu testul Wilcoxon I.

Explicati eventualele contradicții.

Indicatie 1) $R=15$. Prin aproximarea gaussiană pentru $\epsilon=0,05$ găsim

$$P(R \leq 15) \cong 0,0388$$

Se respinge deci H .

2. Suma indicilor lui X este $T=462$. Folosim aproximarea gaussiană pentru $\epsilon=0,05$ și găsim

$P(T \leq 462) > 0,5$, deci ipoteza H se acceptă. Aici trebuie să respingem H deoarece se vede clar că cele două legi "sunt" departe una de alta. Testul Wilcoxon I "nu a mers" deoarece nerealizarea lui $H \implies Y$ este stochastic superioară lui X , lucru evident neadevărat, din sondaj. Deci trebuie să acceptăm pe H , dar cu rezerve. Probabil că cele două legi nu sunt aceleasi dar se "întrepătrund".

3. Se testează un medicament nou pe un lot de 13 soareci și se obțin următoarele rezultate relative la o anumită analiză ("mare" înseamnă înrăutățirea stării individului):

soareci netratati: 45 88 16 6 28 122 62 13
 soareci tratati: 23 104 2 9 30

Folositi testul lui Wilcoxon I pentru a testa ipoteza H : medicamentul nu dă rezultate. (prag $\epsilon=0,05$ și $\epsilon=0,1$).

Indicatie. Se foloseste aproximarea gaussiană și se găseste că $P(T \leq 30) \approx 0,255$. Ipoteza se acceptă.

4. Douăzeci de stupi cu albine se lasă pe aceeași perioadă a anului în două zone diferite A (zece stupi) și B (zece stupi) timp de 20 de ani. Se observă câte kilograme de miere se obțin de la ei în fiecare an în cele 2 zone:

Anul	A	B	Anul	A	B
1	68,3	72,5	11	32,2	31,9
2	60,1	56,0	12	63,3	58,1
3	52,2	55,8	13	54,2	52,7
4	41,7	39,2	14	47,0	46,2
5	32,0	31,4	15	91,9	90,2
6	30,9	35,5	16	56,1	55,4
7	39,3	39,2	17	79,6	75,1
8	42,0	41,1	18	81,2	86,6
9	37,7	43,3	19	78,4	75,3
10	33,5	31,7	20	46,6	43,8

Să se testeze cu testul semnelor, apoi cu testul Wilcoxon II ($\epsilon=0,05$) ipoteza: H : cele două zone melifere A și B sunt tot atât de productive.

Indicatie. Fie X v.a. ce măsoară greutatea mierii provenită din zona A și Y v.a. core-spunzătoare zonei B.

Testul semnelor $Z=Y-X$ conduce la legea binomială $B(20, \frac{1}{2})$ și deci

$$P(V \leq 4) = 0,00591 < 0,05$$

, lucru ce conduce la respingerea ipotezei H .

Testul Wilcoxon II conduce la $W=71$. Aproximăm W cu legea gaussiană și găsim

$$P(W \leq 71) = P\left(\frac{W - 105}{\sqrt{717,5}} \leq -1,27\right) = 0,102 > 0,05$$

, deci ipoteza H se acceptă în acest caz. Deoarece avem puține "minusuri" (doar 4) vom prefera testul Wilcoxon II. Considerăm numai o pură întâmplare că avem puține minusuri. În general, când numărul minusurilor în testul semnelor este mic, nu putem să ne bazăm pe acest test. El este "slab" în acest caz.

Lecția 14

Analiza dispersiei și analiza regresiei

14.1 Analiza dispersiei

Vom analiza aici cea mai simplă problemă dispersională.

Problemă Se consideră s variabile aleatoare gaussiene X_1, \dots, X_s de aceeași dispersie necunoscută σ^2 . Se notează cu $m_i = M(X_i)$. Vrem să testăm următoarea ipoteză: $H: m_1 = m_2 = \dots = m_s$, adică toate v.a. au aceeași medie.

Presupunem că avem pentru fiecare $i = 1, \dots, s$ câte un esantion de volum n_i : $X_{i1}, X_{i2}, \dots, X_{in_i}$, al v.a. X_i . Presupunem că toate cele $n = n_1 + n_2 + \dots + n_s$ v.a. X_{11}, \dots, X_{sn_s} sunt independente. Notăm cu

$$\bar{X}_i = \left(\frac{1}{n_i} X_{i1} + \dots + X_{in_i} \right), \text{ pentru } i=1, \dots, s \quad (14.1)$$

$$\bar{X} = \frac{1}{n} \sum_{\substack{1 \leq i \leq s \\ 1 \leq j \leq n_i}} X_{ij} = \sum_{1 \leq i \leq s} \frac{n_i}{n} \bar{X}_i \quad (14.2)$$

$$Q_A = \sum_{1 \leq i \leq s} n_i (\bar{X}_i - \bar{X})^2 = \left(\sum_{1 \leq i \leq s} n_i \bar{X}_i^2 \right) - n \bar{X}^2 \quad (14.3)$$

$$Q_R = \sum_{\substack{1 \leq i \leq s \\ 1 \leq j \leq n_i}} (X_{ij} - \bar{X}_i)^2 = \left(\sum_{\substack{1 \leq i \leq s \\ 1 \leq j \leq n_i}} X_{ij}^2 \right) - \sum_{1 \leq i \leq s} n_i \bar{X}_i^2. \quad (14.4)$$

• Se știe că statistica Q_R/σ^2 se spune unei legi Pearson cu $n - s$ grade de libertate și că variabila aleatoare.

• $U_{ij} = \sqrt{\frac{(n_i+n_j)(n-s)}{n_i n_j}} \frac{[\bar{X}_i - \bar{X}_j - (m_i - m_j)]}{\sqrt{Q_R}}$ are o distribuție Student cu $n - s$ grade de libertate pentru orice i, j ca mai sus.

• De asemenea, statistica $W = \frac{(n-s)Q_A}{(s-1)Q_R}$ se supune unei legi Fisher-Snedecor (distribuția F) cu $\left\{ \begin{matrix} s-1 \\ n-s \end{matrix} \right\}$ grade de libertate, dacă ipoteza H este adevărată.

Această ultimă observație va constitui esența testului următor: Respinge $H \iff W \geq w$, unde w este cel mai mare număr astfel încât $P(FS \geq w) \leq \epsilon$, unde FS este v.a. Fisher-Snedecor cu $\left\{ \begin{matrix} s-1 \\ n-s \end{matrix} \right\}$ grade de libertate, iar ϵ este un prag de semnificație, considerat mic, de exemplu $\epsilon = 0,05; 0,1; \text{etc.}$ Este posibil ca în tabele să găsim cuantilele de ordin 0,95; 0,9; etc. Se trece atunci la probabilitatea evenimentului contrar, etc.

• Dacă datele X_{ij} sunt mari se înlocuiesc acestea cu datele $aX_{ij} + b$, unde $a \neq 0, b \in \mathbf{R}$ sunt alese astfel încât numerele $aX_{ij} + b$ să devină mici. Prin această schimbare v.a. U_{ij} și W nu se modifică, deci testul decurge exact ca mai sus pentru noile date.

Exemplul 14.1 Pe patru soluri diferite A_1, A_2, A_3, A_4 se plantează orz. Se fac selecții de volume diferite din tulpini de orz ajunse la maturitate din cele patru soluri și se notează lungimea acestora (în cm):

A_1	A_2	A_3	A_4
380	350	354	376
376	358	360	344
360	356	362	342
368	376	352	372
372	338	366	374
366	342	372	360
374	366	362	
382	350	344	
	344	342	
	364	358	
		351	
		348	
		348	

Se notează cu X_i lungimea aleatoare a unei tulpini de orz de pe terenul A_i . Se presupune că X_i sunt gaussiene cu aceeași dispersie σ^2 . Fie pragul $\epsilon = 0,05$.

1) Testati ipoteza H : X_1, X_2, X_3, X_4 au aceeași medie.

2) Testati ipoteza H : X_2, X_3 și X_4 au aceeași medie.

Soluție Deoarece datele sunt mari le centrăm cu ajutorul transformării $Z_i = X_i - 330$. Obținem

un nou tabel:

A ₁	A ₂	A ₃	A ₄
50	20	24	46
46	28	30	14
30	26	32	12
38	46	22	42
42	8	36	44
36	12	42	30
44	36	32	
52	20	14	
	14	12	
	34	28	
		21	
		18	
		18	

1) Folosim o analiză dispersională pentru a testa ipoteza H. Aici avem $s=4$, $n_1=8$, $n_2=10$, $n_3=13$, $n_4=6$, $n=8+10+13+6=37$; $n_1\bar{Z}_1=338$; $n_2\bar{Z}_2=244$; $n_3\bar{Z}_3=188$; $n\bar{Z}=n_1\bar{Z}_1+\dots+n_4\bar{Z}_4=1099$.

Prin urmare $\sum_{1 \leq i \leq 4} n_i \bar{Z}_i^2 = 34449$, $n\bar{Z}^2 = 32640$, deci $Q_A = 1809$. Cum $\sum_{i,j} Z_{ij}^2 = 38229$, avem $Q_R = 3780$. De aici $W = 5,26$.

Probabilitatea ca o v.a. Fisher-Snedecor cu $\left\{ \begin{matrix} 3 \\ 33 \end{matrix} \right\}$ grade de libertate să ia o valoare asemănătoare lui W este inferioară lui 0,01. Prin urmare ipoteza H se respinge.

2) În acest caz $s'=3$, $n'=n_2+n_3+n_4=29$, $Q'_A=199$; $Q'_R=3400$ și $W'=0,76$. Această valoare a lui W' este prea mică, deci ipoteza se acceptă în acest caz

14.2 Analiza regresiei

Fie X și Y două variabile aleatoare și $X_1, \dots, X_n, Y_1, \dots, Y_n$ v.a. de selecție de același volum n . Se pune problema de a studia *legătura* (dacă există) dintre cele două v.a. X și Y numai din analiza unor cupluri de selecție de tipul $\{x_1, \dots, x_n\}, \{y_1, \dots, y_n\}$. Este posibil ca cele două v.a. să nu fie "corelate", adică *coeficientul de corelație* $\rho_{XY} = \frac{\text{cov}(X,Y)}{\sigma_X \sigma_Y} = \frac{M((X-\bar{X})(Y-\bar{Y}))}{\sigma_X \sigma_Y}$ să fie zero. Acest lucru nu înseamnă că nu poate exista o relație funcțională de forma $F(X,Y)=0$ între v.a. X și Y . Această egalitate poate să nu fie "deterministă". Sau poate să fie astfel, dar noi să nu putem descrie matematic această funcție de legătură. În Lectia 6 s-a arătat că X și Y sunt legate între ele printr-o relație liniară: $Y=aX+b$, sau $X=cY+d$, dacă și numai dacă $\rho_{XY}=\pm 1$. De regulă noi facem sondaje în urma cărora estimăm coeficientul de corelație printr-o formulă empirică de forma:

$$\rho_{XY}^* = \frac{\sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y})}{n \cdot \sigma_x^* \cdot \sigma_y^*} \quad (14.5)$$

unde $\bar{x} = m_x^* = \frac{\sum_{i=1}^n x_i}{n}$ și $\bar{y} = m_y^* = \frac{\sum_{i=1}^n y_i}{n}$, $\sigma_x^* = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n}}$, $\sigma_y^* = \sqrt{\frac{\sum (y_i - \bar{y})^2}{n}}$.

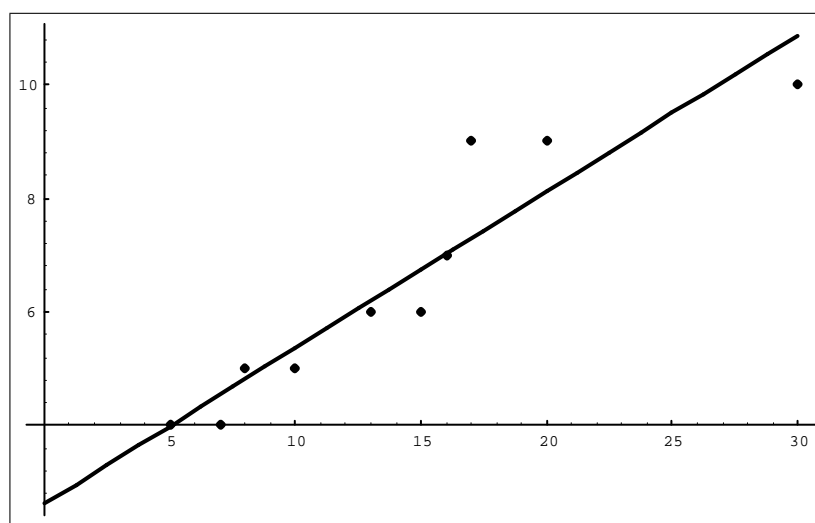
Dacă ρ_{XY}^* se apropie de +1 sau de -1 putem spera ca X și Y să fie corelate liniar. Dacă nu, se continuă investigația prin analize mai fine.

Să începem prin a examina următorul exemplu:

Exemplul 14.2 *Ne interesează care este legătura între numărul de ore afectate de un student de inteligență medie studiului Analizei matematice (lunar) și rezultatele obținute de acesta la examen. În urma unui sondaj efectuat pe 10 studenți s-au obținut următoarele rezultate:*

Student	Nr. ore: x	Nota: y
1	5	4
2	7	4
3	8	5
4	10	5
5	13	6
6	15	6
7	16	7
8	17	9
9	20	9
10	30	10

Punem într-un grafic aceste date:



Dreapta de regresie

Se pune problema dacă aceste puncte sunt "foarte" apropiate de o dreaptă. Mai exact, să notăm cu (x_i, y_i) , $i = \overline{1, n}$ valorile obținute dintr-un sondaj pentru v.a. X și Y . Există $b_0, b_1 \in \mathbb{R}$ astfel încât diferențele $e_i = y_i - b_0 - b_1 x_i$ să fie "mici"? Care sunt "cei mai buni" b_0 și b_1 care să facă acest lucru? Sau poate există b_0, b_1, \dots, b_k astfel încât diferențele $e_i = y_i - b_0 - b_1 x_i - b_2 x_i^2 - \dots - b_k x_i^k$ să fie mici pentru orice $i = \overline{1, n}$?

Definiția 14.3 *Ecuatia*

$$y_i = \beta_0 + \beta_1 x_i + e_i, i = \overline{1, n} \quad (14.6)$$

se numeste *model de regresie simplă (sau liniară)*, iar ecuatia

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \dots + \beta_k x_i^k + e_i \quad (14.7)$$

unde $i = \overline{1, n}$, se numeste *model de regresie multiplă*.

De exemplu, pentru $k = 2$ se numeste *regresie parabolică*, etc.

Noi ne vom ocupa aici în exclusivitate cu regresia liniară (simplă). Vom interpreta y_i ca fiind valorile unei v.a.. La fel vom interpreta valorile *erorilor* e_i . De asemenea vom interpreta β_0 și β_1 ca fiind valorile unor v.a. pe care le vom determina prin metoda celor mai mici pătrate.

14.2.1 Metoda celor mai mici pătrate (C. F. Gauss)

Dacă vrem să aproximăm $y_i \approx \beta_0 + \beta_1 x_i$ eroarea comisă este e_i , $i = \overline{1, n}$. Vom pune condiția ca suma pătratelor erorilor e_i să fie minimă:

$$S = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 = \text{minimă}. \quad (14.8)$$

(vezi și Lecția 8). Este ușor de arătat că $S = S(\beta_0, \beta_1)$ are un singur minim pentru $\beta_0 = b_0$, $\beta_1 = b_1$, unde b_0 și b_1 reprezintă soluția sistemului liniar

$$\frac{\partial S}{\partial \beta_0} = -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) = 0 \quad (14.9)$$

$$\frac{\partial S}{\partial \beta_1} = -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) x_i = 0 \quad (14.10)$$

Notăm cu $\bar{y} = n^{-1} \sum_{i=1}^n y_i$ și cu $\bar{x} = n^{-1} \sum_{i=1}^n x_i$.

Atunci, din (14.9) rezultă că :

$$b_0 = \bar{y} - b_1 \bar{x} \quad (14.11)$$

Dar b_0 și b_1 verifică și (14.10):

$$\sum_{i=1}^n y_i x_i - n b_0 \bar{x} - b_1 \sum_{i=1}^n x_i^2 = 0 \quad (14.12)$$

(14.11) și (14.12) ne conduc la expresia lui b_1 :

$$b_1 = \left(\sum_{i=1}^n y_i x_i - n \bar{x} \bar{y} \right) / \left(\sum_{i=1}^n x_i^2 - n \bar{x}^2 \right) \quad (14.13)$$

Nu este greu de arătat că b_1 se mai poate scrie sub formă "centrată":

$$b_1 = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{cov(x, y)}{\sigma_x^{*2}} \quad (14.14)$$

Cu această expresie a lui b_1 venim în (14.11) și găsim

$$b_0 = \bar{y} - \bar{x} \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) / \sum_{i=1}^n (x_i - \bar{x})^2 = \bar{y} - \bar{x} \frac{cov(x, y)}{\sigma_x^{*2}} \quad (14.15)$$

(Vezi și formulele corespunzătoare din Lecția 8). În exemplul 14.2 metoda celor mai mici pătrate ne dă $b_0 = 2,621$, $b_1 = 0,274$, $y = b_0 + b_1 x$ fiind dreapta cea mai apropiată de norul de puncte respectiv (vezi figura "Dreapta de regresie").

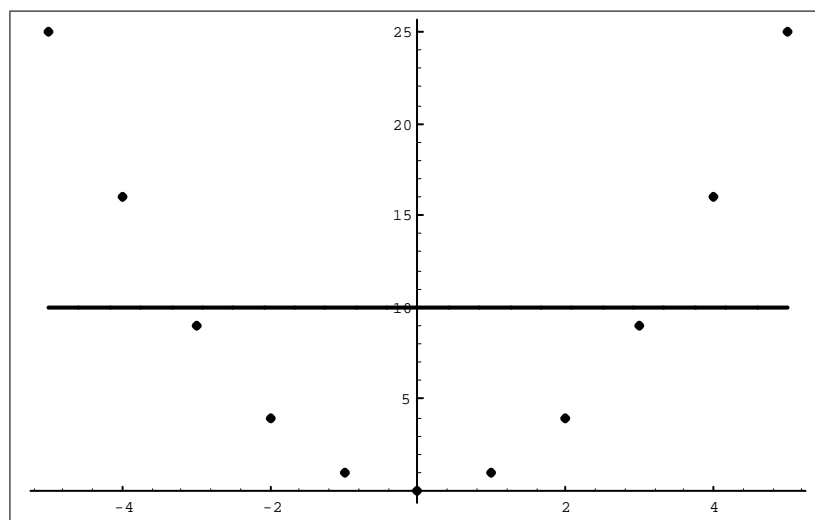
14.2.2 Condițiile Gauss–Markov pentru metoda celor mai mici pătrate

Condițiile Gauss–Markov sunt condiții naturale care se impun v.a. e_i , $i = \overline{1, n}$. Prima condiție cere ca media v.a. e_i să fie zero:

$$M(e_i) = 0, \quad \text{pentru orice } i = \overline{1, n} \quad (14.16)$$

Să observăm că oricum $\sum_{i=1}^n e_i = 0$, deci $\sum_{i=1}^n M(e_i) = 0$, în general.

Figura următoare ne arată un caz în care nu are loc (14.16).

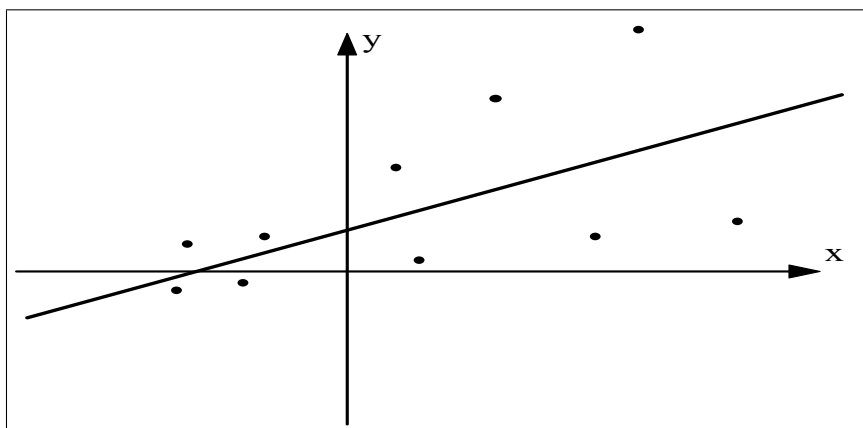


Un caz în care nu este îndeplinită prima condiție Gauss-Markov

A doua condiție Gauss-Markov cere ca dispersiile v.a. e_i să fie constante, adică

$$D(e_i) = \sigma^2 = \text{constantă, dar necunoscută}$$

Figura următoare ne prezintă o situație în care *nu* are loc (14.17) deoarece $D(e_i)$ cresc odată cu creșterea i -urilor.



Un caz în care nu este îndeplinită a doua condiție Gauss Markov

Uneori chiar un singur punct poate să facă condițiile (14.16) sau (14.17) neadevărate.

Dacă observațiile noastre sunt corelate unele cu altele *nu* putem face mai târziu aprecieri pertinente asupra lor. De aceea vom micsora numărul acestor observații până când acestea vor deveni necorelate.

Ultima condiție Gauss–Markov se referă tocmai la acest lucru. Se cere ca $\text{cov}(e_i, e_j) = 0$. Cum $M(e_i) = M(e_j) = 0$, rămâne doar condiția:

$$M(e_i e_j) = 0, \text{ pentru orice } i \neq j. \quad (14.17)$$

Definiția 14.4 *Metoda celor mai mici pătrate se spune că este o metodă bună dacă variabilele aleatoare erori, e_i , $i=1, 2, \dots, n$ îndeplinesc cele trei condiții Gauss–Markov (14.16), (14.17) și (14.18).*

14.2.3 Măsura deviației la metoda celor mai mici pătrate

Am văzut mai sus că $S = \sum_{i=1}^n e_i^2$ măsoară deviația adevăratelor y_i de la valorile estimate prin metodă, $\hat{y}_i = b_0 + b_1 x_i$, deoarece $e_i = y_i - \hat{y}_i$. Cum nu se dorește ca măsura deviației să depindă de unitatea de măsură, se lucrează cu altă mărime, oarecum relativă.

Definiția 14.5 *Fie modelul de regresie liniară $y_i = \beta_0 + \beta_1 x_i + e_i$, $i = 1, 2, \dots, n$.*

- a) *Dacă $\beta_0 \neq 0$ se ia ca măsură a deviației expresia: $R^2 = 1 - \frac{\sum_{i=1}^n e_i^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$.*
- b) *Dacă $\beta_0 = 0$ se ia ca măsură a deviației expresia: $R^2 = 1 - \frac{\sum_{i=1}^n e_i^2}{\sum_{i=1}^n y_i^2}$.*

• Media și varianta v.a. b_0 și b_1

Să interpretăm acum pe y_i și pe e_i ca variabile aleatoare, pe β_0 și β_1 ca niste constante (parametri), iar pe b_0 și b_1 ca variabile aleatoare care iau diferite valori la fiecare selecție în parte.

Teorema 14.6 *Fie modelul de regresie liniară $y_i = \beta_0 + \beta_1 x_i + e_i$, $i = \overline{1, n}$ în care v.a. e_i , $i = \overline{1, n}$ verifică cele 3 condiții Gauss–Markov. Atunci avem relațiile*

$$M(b_0) = \beta_0 \quad D(b_0) = \sigma^2 \left[n^{-1} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]$$

$$M(b_1) = \beta_1 \quad D(b_1) = \sigma^2 \Big/ \sum_{i=1}^n (x_i - \bar{x})^2 \quad (14.18)$$

Demonstratie Deoarece $\sum_{i=1}^n (x_i - \bar{x}) = 0$ rezultă că

$$\sum_{i=1}^n (y_i - \bar{y}) (x_i - \bar{x}) = \sum_{i=1}^n y_i (x_i - \bar{x}) \quad (14.19)$$

Din (14.20) și (14.14) rezultă că

$$\begin{aligned} b_1 &= \sum_{i=1}^n c_i y_i, \text{ unde} \\ c_i &= (x_i - \bar{x}) \Big/ \sum_{i=1}^n (x_i - \bar{x}) \end{aligned} \quad (14.20)$$

Cu aceste notatii avem:

$$\begin{aligned} \sum_{i=1}^n c_i &= 0 \\ \sum_{i=1}^n c_i x_i &= \sum_{i=1}^n c_i (x_i - \bar{x}) = 1, \text{ de unde} \\ \sum_{i=1}^n c_i^2 &= \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \end{aligned}$$

De aici, aplicând operatorul de medie relației (14.21) în care $y_i = \beta_0 + \beta_1 x_i + e_i$ găsim:

$$M(b_1) = \sum_{i=1}^n c_i M(y_i) = \beta_0 \underbrace{\sum_{i=1}^n c_i}_{=0} + \beta_1 \underbrace{\sum_{i=1}^n c_i x_i}_{=1} + \sum_{i=1}^n c_i \underbrace{M(e_i)}_{=0} = \beta_1 \quad (14.21)$$

Calculăm acum $D(b_1)$:

$$D(b_1) = \sum_{i=1}^n c_i^2 D(y_i) = \left(\sum_{i=1}^n c_i^2 \right) \sigma^2 = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad (14.22)$$

Acum, deoarece $M(\bar{y}) = \frac{\sum_{i=1}^n M(y_i)}{n} = \beta_0 + \beta_1 \bar{x}$, rezultă că

$$M(b_0) = M(\bar{y} - b_1 \bar{x}) = \beta_0 + \beta_1 \bar{x} - \bar{x} M(b_1) = \beta_0 \quad (14.23)$$

Calculăm acum $D(b_0)$:

$$b_0 = n^{-1} \sum_{i=1}^n y_i - \bar{x} \sum_{i=1}^n c_i y_i = \sum_{i=1}^n (n^{-1} - \bar{x} c_i) y_i, \text{ deci}$$

$$\begin{aligned} D(b_0) &= \sum_{i=1}^n [n^{-1} - \bar{x} c_i]^2 D(y_i) \\ &= \sigma^2 \sum_{i=1}^n [n^{-2} - 2n^{-1} \bar{x} c_i + \bar{x}^2 c_i^2] \\ &= \sigma^2 \left[n^{-1} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right] \end{aligned}$$

deoarece $\sum_{i=1}^n c_i = 0$.

În cazul în care $\beta_0 = 0$, avem direct din $\frac{\partial S}{\partial \beta_1} = 0$ că

$$b_1 = \frac{\sum_{i=1}^n y_i x_i}{\sum_{i=1}^n x_i^2} \quad e_i = y_i - b_1 x_i \quad (14.24)$$

(în general $\sum_{i=1}^n e_i \neq 0$).

Înlocuim în (14.25) pe y_i cu $\beta_1 x_i + e_i$ și găsim că

$$b_1 = \frac{\beta_1 \sum_{i=1}^n x_i^2}{\sum_{i=1}^n x_i^2} + \frac{\sum_{i=1}^n e_i x_i}{\sum_{i=1}^n x_i^2} = \beta_1 + \frac{\sum_{i=1}^n e_i x_i}{\sum_{i=1}^n x_i^2} \quad (14.25)$$

De aici rezultă că

$$D(b_1) = \sigma^2 \frac{\sum_{i=1}^n x_i^2}{\left(\sum_{i=1}^n x_i^2 \right)^2} = \frac{\sigma^2}{\sum_{i=1}^n x_i^2}$$

și că $M(b_1) = \beta_1$, deoarece $M(e_i) = 0$, pentru orice $i = \overline{1, n}$. \square

Corolarul 14.7 *Cu formele de mai sus, rezultă din Teorema 14.6 că b_0 este un estimator nedeplasat al parametrului β_0 și că b_1 este un estimator nedeplasat al parametrului β_1 . \square*

Formulele (14.19) arată că $D(b_0)$ și $D(b_1)$ contin pe σ^2 care este necunoscut. De obicei σ^2 se estimează cu estimatorul nedeplasat (verificarea nu este simplă!)

$$s^2 = (n-1)^{-1} \sum_{i=1}^n e_i^2 \quad (14.26)$$

14.2.4 Intervale de încredere și teste pentru β_0 și β_1

Presupunem că avem modelul de regresie liniară: $y_i = \beta_0 + \beta_1 x_i + e_i$, $i = \overline{1, n}$, unde (y_i) și (e_i) sunt v.a., iar β_0 și β_1 sunt considerați parametri statistici care au fost estimați mai sus prin b_0 și b_1 .

Presupunem că acest model îndeplinește condițiile Gauss–Markov. De asemenea facem presupunerea că v.a. (e_i) au o distribuție normală $N(0, \sigma^2)$. Atunci rezultă (vezi Lectia 4) că (y_i) au o distribuție normală $N(\beta_0 + \beta_1 x_i, \sigma^2)$. Cum b_0 și b_1 sunt combinații liniare de (y_i) -uri rezultă că și ele sunt v.a. cu mediiile și dispersiile date în formula (14.19).

Se poate arăta că v.a. $(b_j - \beta_j) / \sqrt{D(b_j)}$ este o v.a. Student cu $n-2$ grade de libertate, pentru $j = 0, 1$ (dacă $\beta_0 \neq 0$).

Folosim acum teoria testelor de semnificație și găsim că intervalul

$$\left(b_j - \sqrt{D(b_j)} T_{n-2, \alpha/2} \quad b_j + \sqrt{D(b_j)} T_{n-2, \alpha/2} \right)$$

este un interval de încredere pentru β_j , $j=0,1$, de $(1-\alpha) \times 100$ procente. Aici $T_{n-2, \alpha/2}$ este cuantila de ordin $\alpha/2$ a distribuției Student cu $n-2$ grade de libertate.

În lumina rezultatelor de mai sus vom studia pe scurt următoarea situație.

Fie X_1, X_2, \dots, X_n variabile aleatoare gaussiene independente cu aceeași dispersie σ^2 astfel încât să existe două constante a, b cu proprietatea: $M(X_i) = a + bt_i$, $1 \leq i \leq n$.

Ca și mai sus (înlocuim pe Y_i cu $M(X_i)$!) introducem următoarele notații:

$$\begin{aligned} \bar{t} &= \frac{1}{n} \sum_{i=1}^n t_i; \quad \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \\ S_t^2 &= \frac{1}{n} \sum_{i=1}^n (t_i - \bar{t})^2 = \frac{1}{n} \sum_i t_i^2 - \bar{t}^2; \end{aligned}$$

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

$$S_{tX}^2 = \frac{1}{n} \sum_i (t_i - \bar{t}) (X_i - \bar{X}) = \frac{1}{n} \sum_i t_i X_i - \bar{t} \bar{X}$$

Statisticile (care dau estimatori pentru a și b): $\hat{b} = S_{tX}^2 / S_t^2$ și $\hat{a} = \bar{X} - \bar{t} \hat{b}$ sunt v.a. gaussiene cu $M(\hat{b}) = b$,

$$M(\hat{a}) = a, \quad D(\hat{a}) = \frac{\sigma^2}{n} \left(1 + \frac{\bar{t}^2}{S_t^2} \right), \quad D(\hat{b}) = \frac{\sigma^2}{n S_t^2}$$

(vezi Teorema 1). Se poate arăta că statistica

$$\rho_{tX}^2 = \frac{1}{n} \sum_i (X_i - \hat{a} - \hat{b} t_i)^2 = \frac{S_t^2 \cdot S_X^2 - (S_{tX}^2)^2}{S_t^2} = S_X^2 - \hat{b} S_{tX}^2$$

este independentă de statisticile \hat{a} și \hat{b} . Se poate arăta de asemenea că $\frac{n}{\sigma^2} \rho_{tX}^2$ se supune unei legi Pearson (χ^2) cu $n-2$ grade de libertate. Această observație ne permite să construim un test de semnificație și un nou tip de intervale de încredere pentru parametrii a și b .

Mai mult, se poate arăta că statisticile

$$T = \frac{\sqrt{S_t^2} (\hat{b} - b)}{\sqrt{\frac{1}{n-2} \rho_{tX}^2}} = \frac{\sqrt{n-2} S_t^2 (\hat{b} - b)}{\sqrt{S_t^2 \cdot S_X^2 - (S_{tX}^2)^2}}$$

și

$$U = \frac{\sqrt{n(n-2)} S_t^2 (\hat{a} - a)}{\sqrt{\sum_i t_i^2} \sqrt{S_t^2 \cdot S_X^2 - (S_{tX}^2)^2}} = \frac{\sqrt{S_t^2} (\hat{a} - a)}{\sqrt{\frac{1}{n-2} \rho_{tX}^2 (S_t^2 + \bar{t}^2)}}$$

satisfac legea Student cu $n-2$ grade de libertate. Si ele pot fi folosite pentru cei doi parametri a și b .

Exemplul 14.8 *Un biolog studiază creșterea unei specii de plante, pe mai multe exemplare, într-un interval de timp dat. La începutul perioadei planta avea (în mm) înălțimea inițială t . La sfârșitul perioadei ea a avut înălțimea X . S-au făcut 10 probe:*

t	57	60	52	49	56	46	51	63	49	57
X	86	93	77	67	81	70	71	91	67	82

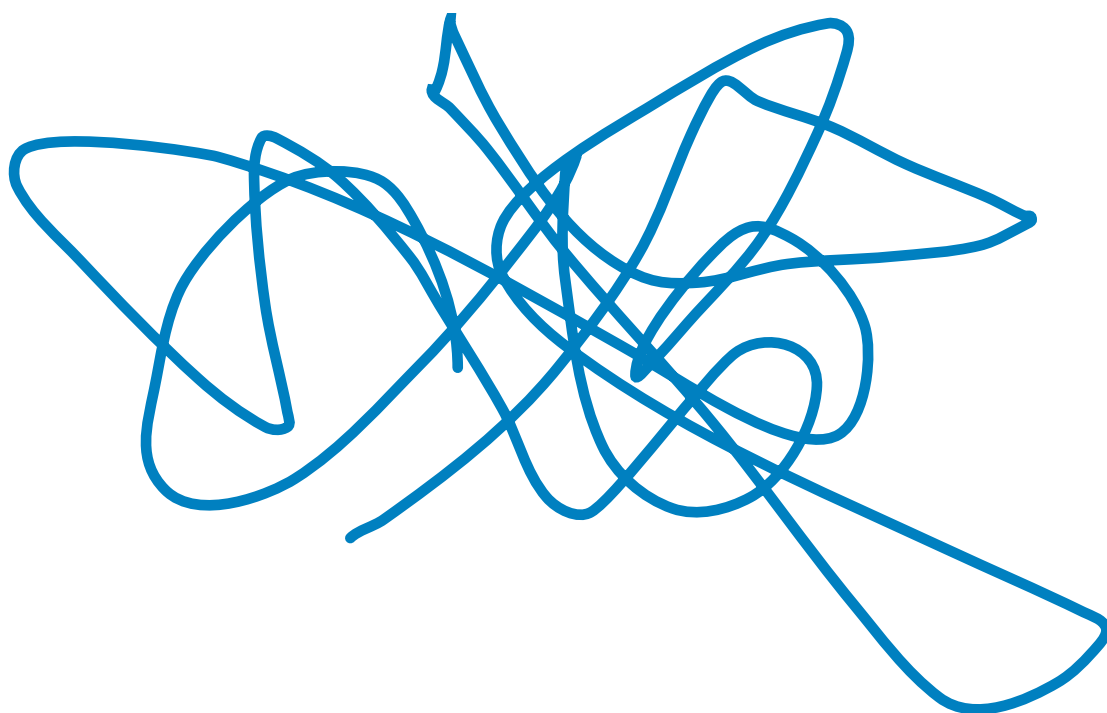
- 1) Găsiți estimatori punctuali pentru a , b și σ^2 .
- 2) Estimați înălțimea unei plante la sfârșitul perioadei dacă inițial ea a avut 52 mm.
- 3) Dati pentru b un interval de încredere de 95%.

Soluție 1) Folosim formulele de mai sus și găsim: $n=10$; $n\bar{t}=540$; $n\bar{X}=785$; $\sum_{i=1}^{10} t_i^2=29426$; $\sum_i X_i^2=62459$; $\sum_i t_i X_i=42836$; $\bar{t}=54$; $\bar{X}=78,5$; $nS_t^2=266$; $nS_{tX}^2=446$; $nS_X^2=836,5$. De aici se găsește pentru b estimarea $\hat{b}=1,677$ și pentru a , $\hat{a}=-12,06$.

Dar $n\rho_{tX}^2=88,5$. Folosim acum faptul că $\frac{n\rho_{tX}^2}{\sigma^2}$ se supune unei legi Pearson cu $n-2=8$ grade de libertate. Pentru σ^2 avem estimarea $\hat{\sigma}^2=\frac{n\rho_{tX}^2}{8}=11,06$.

2) O plantă cu $t=52$ la începutul perioadei, va avea la sfârșitul duratei lungimea $\hat{a}+52\hat{b}=75,14$.

3) Dacă $Y=\hat{b}$ și $Z=\sqrt{\frac{\rho_{tX}^2}{S_t^2(n-2)}}$, v.a. $T=\frac{Y-\hat{b}}{Z}$ se supune unei legi Student cu $n-2=8$ grade de libertate. Avem deci că $P(|T| > t_\varepsilon)=0,05$ ($=95\%$) pentru $t_\varepsilon=2,306$. Se obține de aici un interval de încredere de prag 95% pentru b : $\hat{b} - t_\varepsilon |Z| \leq b \leq \hat{b} + t_\varepsilon |Z|$. Dar $Z=\sqrt{\frac{n\rho_{tX}^2}{n(n-2)S_t^2}}$ și deci $t_\varepsilon |Z|=0,47$. Găsim în final intervalul $1,207 \leq b \leq 2,147$.



Bibliografie

1. Alain Cambrouze - Probabilités et Statistique, Press Universitaires de France, 1993
2. G. Ciucu, V. Craiu - Introducere teoria probabilităților și Statistica matematică, Editura Didactică și Pedagogică, București, 1971
3. Harald Cramer - Mathematical Methods of Statistics, Princeton University Press, 1946
4. W. Feller - An introduction to Probability Theory and Its Applications, Vol. I, John Wiley & Sons, Inc. 1960
5. B.V. Gnedenko - The theory of Probability, Mir, Moscow, 1969
6. M. Iosifescu, Gh. Mihoc, R. Theodorescu - Teoria probabilităților și Statistică matematică, Editura Tehnică, București, 1966
7. P. Jaquard - Probabilités, Statistique (Culegere de probleme), Masson, Paris, 1972
8. A. Krief, S. Levy - Calcul des Probabilités (Exercices), Hermann, Paris, 1982
9. D. Lungu, D. Chiocel - Metode probabilistice în calculul construcțiilor, Editura Tehnică, 1982
10. Ashis Sen, Muni Srivastava - Regression Analysis, Theory, Methods and Applications, Springer Texts in Statistics, Springer-Verlag, New-York, Inc. 1990
11. H. Ventsel - Théorie des probabilités, Edition Mir, Moscou, 1973
12. R.L. Winkler, W. L. Hays - Statistics, Holt, Reinhart and Winston, N.Y., 1975

x	0	1	2	3	4	5	6	7	8	9
0.0	0.0000	00399	00798	01197	01595	01994	02392	02790	03188	03586
0.1	00583	04380	04776	05172	05567	05962	06356	06749	07142	07535
0.2	07928	08322	08715	09108	09500	09892	10283	10674	11064	11453
0.3	11791	12172	12552	12930	13307	13683	14058	14431	14803	15173
0.4	15542	15910	16276	16640	17003	17364	17724	18082	18438	18793
0.5	19145	19497	19847	20194	20540	20884	21226	21566	21904	22240
0.6	22575	22907	23237	23565	23891	24215	24537	24857	25175	25490
0.7	25804	26115	26424	26730	27035	27337	27637	27935	28230	28524
0.8	28814	29103	29390	29673	29953	30230	30504	30775	31043	31307
0.9	31564	31829	32091	32351	32609	32864	33117	33368	33616	33861
1.0	34114	34375	34634	34890	35143	35394	35643	35890	36134	36374
1.1	36633	36880	37124	37366	37605	37841	38075	38306	38534	38759
1.2	38993	39224	39452	39677	39899	40118	40334	40547	40757	40964
1.3	41168	41375	41579	41780	41978	42173	42365	42554	42740	42923
1.4	43103	43284	43462	43637	43809	43978	44144	44307	44466	44622
1.5	44775	44933	45088	45240	45389	45535	45678	45818	45955	46089
1.6	46220	46353	46483	46610	46734	46855	46973	47088	47199	47307
1.7	47412	47519	47623	47724	47821	47916	48008	48097	48183	48266
1.8	48345	48429	48509	48586	48660	48731	48799	48864	48926	48985
1.9	49041	49097	49150	49200	49248	49293	49335	49374	49410	49444
2.0	49475	49511	49543	49572	49598	49621	49641	49658	49672	49683
2.1	49691	49703	49712	49719	49724	49727	49728	49727	49724	49719
2.2	49712	49706	49698	49688	49676	49661	49644	49624	49601	49576
2.3	49549	49523	49495	49464	49431	49395	49356	49314	49269	49221
2.4	49171	49121	49068	49012	48953	48891	48826	48758	48687	48612
2.5	48536	48461	48383	48302	48218	48131	48041	47948	47852	47753
2.6	47651	47553	47452	47348	47241	47131	47018	46902	46783	46660
2.7	46534	46408	46279	46147	46012	45874	45732	45587	45439	45287
2.8	45132	44977	44819	44658	44494	44326	44155	43980	43802	43621
2.9	43437	43253	43066	42876	42682	42485	42284	42079	41871	41660
3.0	41445	41229	40999	40765	40527	40285	40039	39789	39535	39277
3.1	39015	38749	38479	38205	37927	37645	37359	37068	36773	36474
3.2	36171	35874	35573	35268	34959	34646	34329	34007	33681	33351
3.3	33018	32681	32340	32005	31666	31323	30975	30623	30267	29907
3.4	29543	29179	28811	28439	28063	27683	27299	26911	26519	26123
3.5	25724	25318	24909	24496	24079	23658	23233	22804	22371	21934
3.6	21492	21056	20616	20172	19724	19272	18816	18356	17891	17421
3.7	16947	16477	16003	15525	15042	14554	14061	13563	13060	12552
3.8	12039	11526	11008	10485	9957	9424	8885	8340	7790	7234
3.9	6673	6112	5545	4972	4393	3808	3217	2620	2017	1408
4.0	859	519	39	519	39	519	39	519	39	519
4.1	519	39	519	39	519	39	519	39	519	39
4.2	39	519	39	519	39	519	39	519	39	519
4.3	519	39	519	39	519	39	519	39	519	39
4.4	39	519	39	519	39	519	39	519	39	519
4.5	519	39	519	39	519	39	519	39	519	39
4.6	39	519	39	519	39	519	39	519	39	519
4.7	519	39	519	39	519	39	519	39	519	39
4.8	39	519	39	519	39	519	39	519	39	519
4.9	519	39	519	39	519	39	519	39	519	39

I

TABLE

TABELUL I DISTRIBUȚIA NORMALĂ STANDARD

$$F(z) = \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} dt$$

• Exemple: $F(1,02) = 0,8461358$, $F(-3,40) = 1 - F(3,40) = 1 - 0,9995631 = 0,0003369$.

z	$F(z)$	z	$F(z)$	z	$F(z)$	z	$F(z)$
.00	.5000000	.38	.6495764	.72	.7642375	1.08	.8599289
.01	.5039894	.37	.6443098	.73	.7673049	1.09	.8621434
.02	.5079783	.36	.6450273	.74	.7703500	1.10	.8643339
.03	.5119665	.39	.6517347	.75	.7733726	1.11	.8665005
.04	.5159534	.40	.6554217	.76	.7763727	1.12	.8686431
.05	.5199388	.41	.6590970	.77	.7793501	1.13	.8707619
.06	.5239222	.42	.6627571	.78	.7823046	1.14	.8728568
.07	.5279032	.43	.6664022	.79	.7852361	1.15	.8749281
.08	.5318814	.44	.6700314	.80	.7881446	1.16	.8769756
.09	.5358564	.45	.6736448	.81	.7910299	1.17	.8789995
.10	.5398278	.46	.6772417	.82	.7938919	1.18	.8809999
.11	.5437953	.47	.6808224	.83	.7967306	1.19	.8829768
.12	.5477584	.48	.6843866	.84	.7995459	1.20	.8849303
.13	.5517168	.49	.6879331	.85	.8023375	1.21	.8868606
.14	.5556700	.50	.6914627	.86	.8051056	1.22	.8887676
.15	.5596177	.51	.6949747	.87	.8078498	1.23	.8906514
.16	.5635595	.52	.6984687	.88	.8105703	1.24	.8925123
.17	.5674949	.53	.7019440	.89	.8132671	1.25	.8943502
.18	.5714237	.54	.7054013	.90	.8159399	1.26	.8961653
.19	.5753454	.55	.7088403	.91	.8185887	1.27	.8979577
.20	.5792597	.56	.7122603	.92	.8212136	1.28	.8997274
.21	.5831662	.57	.7156612	.93	.8238146	1.29	.9014747
.22	.5870654	.58	.7190427	.94	.8263912	1.30	.9031996
.23	.5909581	.59	.7224047	.95	.8289439	1.31	.9049021
.24	.5948439	.60	.7257469	.96	.8314724	1.32	.9065825
.25	.5987203	.61	.7290691	.97	.8339768	1.33	.9082409
.26	.6025881	.62	.7323711	.98	.8364569	1.34	.9098773
.27	.6064479	.63	.7356527	.99	.8389129	1.35	.9114920
.28	.6102992	.64	.7389137	1.00	.8413447	1.36	.9130850
.29	.6140919	.65	.7421530	1.01	.8437524	1.37	.9146565
.30	.6179114	.66	.7453773	1.02	.8461368	1.38	.9162067
.31	.6217596	.67	.7485871	1.03	.8484950	1.39	.9177356
.32	.6255458	.68	.7517747	1.04	.8508300	1.40	.9192433
.33	.6293000	.69	.7549329	1.05	.8531409	1.41	.9207302
.34	.6330317	.70	.7580663	1.06	.8554277	1.42	.9221962
.35	.6368307	.71	.7611747	1.07	.8576903	1.43	.9236415

z	$F(z)$	z	$F(z)$	z	$F(z)$	z	$F(z)$
1.44	.9250063	1.77	.9616304	2.10	.9821356	2.43	.9924906
1.45	.9264707	1.78	.9624623	2.11	.9826708	2.44	.9926564
1.46	.9278550	1.79	.9632730	2.12	.9830070	2.45	.9928572
1.47	.9292191	1.80	.9640897	2.13	.9834142	2.46	.9930531
1.48	.9305674	1.81	.9648521	2.14	.9838226	2.47	.9932443
1.49	.9318873	1.82	.9656205	2.15	.9842224	2.48	.9934300
1.50	.9331923	1.83	.9663736	2.16	.9846137	2.49	.9936128
1.51	.9344753	1.84	.9671169	2.17	.9849966	2.50	.9937903
1.52	.9357446	1.85	.9678432	2.18	.9853713	2.51	.9939634
1.53	.9369916	1.86	.9685572	2.19	.9857379	2.52	.9941323
1.54	.9382198	1.87	.9692581	2.20	.9860966	2.53	.9942969
1.55	.9394292	1.88	.9699460	2.21	.9864474	2.54	.9944574
1.56	.9406201	1.89	.9706210	2.22	.9867906	2.55	.9946139
1.57	.9417924	1.90	.9712834	2.23	.9871263	2.56	.9947664
1.58	.9429466	1.91	.9719334	2.24	.9874545	2.57	.9949151
1.59	.9440826	1.92	.9725711	2.25	.9877753	2.58	.9950600
1.60	.9452007	1.93	.9731960	2.26	.9880884	2.59	.9952012
1.61	.9463011	1.94	.9738133	2.27	.9883932	2.60	.9953388
1.62	.9473839	1.95	.9744117	2.28	.9886902	2.70	.9957630
1.63	.9484493	1.96	.9750021	2.29	.9889803	2.80	.9961449
1.64	.9494974	1.97	.9755863	2.30	.9892750	2.90	.9964842
1.65	.9505285	1.98	.9761482	2.31	.9895559	3.00	.9967801
1.66	.9515428	1.99	.9767043	2.32	.9898206	3.20	.9973129
1.67	.9525403	2.00	.9772495	2.33	.9900909	3.40	.9978031
1.68	.9535213	2.01	.9777784	2.34	.9903581	3.60	.9982409
1.69	.9544860	2.02	.9783083	2.35	.9906133	3.80	.9986277
1.70	.9554345	2.03	.9788217	2.36	.9908665	4.00	.9989683
1.71	.9563671	2.04	.9793248	2.37	.9911060	4.50	.9993066
1.72	.9572838	2.05	.9798173	2.38	.9913437	5.00	.9995987
1.73	.9581849	2.06	.9803005	2.39	.9915758	5.50	.9998600
1.74	.9590706	2.07	.9807738	2.40	.9918025		
1.75	.9599408	2.08	.9812372	2.41	.9920237		
1.76	.9607961	2.09	.9816911	2.42	.9922397		

TABLELA II FRACȚIILE (CUANTILELE) DISTRIBUȚIEI T

- $v = 1, 2, \dots, 50, 40, \dots$ reprezintă numărul gradelor de libertate.
- $\alpha = 0,60, 0,75, 0,90, \dots$ reprezintă ordinul fracției
- dacă ordinul fracției este mai mic decât 0,50 ($\alpha = P(t) < 0,50$), atunci folosim formula $t = -u$, unde $F(u) = 1 - \alpha$, deoarece $F(-t) = 1 - F(t)$.
- Exemple: pentru $v = 11$, $\alpha = 0,95$, fracția este 1,796; pentru $v = 5$, $\alpha = 0,025$, fracția este -2,571.

$F(t)$

α	.60	.75	.90	.95	.975	.99	.995	.999
1	0.325	1.000	3.078	6.314	12.706	31.821	63.657	318.31
2	.289	0.816	1.886	2.920	4.303	6.965	9.923	22.326
3	.277	.765	1.638	2.353	3.182	4.541	5.841	10.213
4	.271	.741	1.533	2.133	2.776	3.747	4.604	7.173
5	0.267	0.727	1.476	2.015	2.571	3.365	4.032	5.893
6	.265	.718	1.440	1.943	2.447	3.143	3.707	5.208
7	.263	.711	1.415	1.895	2.365	2.998	3.499	4.785
8	.262	.706	1.397	1.860	2.306	2.896	3.355	4.501
9	.261	.703	1.383	1.833	2.282	2.821	3.250	4.297
10	0.260	0.700	1.372	1.812	2.225	2.764	3.169	4.144
11	.260	.697	1.363	1.796	2.201	2.718	3.106	4.025
12	.259	.695	1.356	1.782	2.179	2.681	3.055	3.930
13	.259	.694	1.350	1.771	2.160	2.650	3.012	3.852
14	.258	.692	1.345	1.761	2.145	2.624	2.977	3.787
15	0.258	0.691	1.341	1.753	2.131	2.602	2.947	3.733
16	.258	.690	1.337	1.746	2.120	2.583	2.921	3.686
17	.257	.689	1.333	1.740	2.110	2.567	2.898	3.646
18	.257	.688	1.330	1.734	2.101	2.552	2.878	3.610
19	.257	.688	1.328	1.729	2.093	2.539	2.861	3.579
20	0.257	0.687	1.325	1.725	2.086	2.528	2.845	3.552
21	.257	.686	1.323	1.721	2.080	2.518	2.831	3.527
22	.256	.686	1.321	1.717	2.074	2.508	2.819	3.505
23	.256	.685	1.319	1.714	2.069	2.500	2.807	3.485
24	.256	.685	1.318	1.711	2.064	2.492	2.797	3.467
25	0.256	0.684	1.316	1.708	2.060	2.485	2.787	3.450
26	.256	.684	1.315	1.706	2.056	2.479	2.779	3.435
27	.256	.684	1.314	1.703	2.052	2.473	2.771	3.421
28	.256	.683	1.313	1.701	2.048	2.467	2.763	3.408
29	.256	.683	1.311	1.699	2.045	2.462	2.756	3.396
30	0.256	0.683	1.310	1.697	2.042	2.457	2.750	3.385
40	.255	.681	1.303	1.684	2.021	2.423	2.704	3.307
60	.254	.679	1.296	1.671	2.000	2.390	2.660	3.232
120	.254	.677	1.289	1.658	1.980	2.358	2.617	3.160
∞	.253	.674	1.282	1.645	1.960	2.326	2.576	3.090

TABELA III FRACȚIILE (CUANTILELE) DISTRIBUȚIEI χ^2

* Exemplu: pentru 14 grade de libertate fracția de ordin 0,025 este 2,62872.

v	$P(\chi^2)$						
	.005	.01	.025	.05	.10	.25	.50
1	382704 · 10 ⁻¹⁰	157085 · 10 ⁻⁹	982069 · 10 ⁻⁸	393214 · 10 ⁻⁷	0.0157908	0.101308	0.451507
2	0.0100231	0.0201007	0.0506356	0.102587	0.210720	0.575354	1.385829
3	0.0717212	0.114832	0.216795	0.351846	0.584375	1.212534	2.36597
4	0.208989	0.297110	0.484419	0.710721	1.063623	1.92255	3.35670
5	0.411740	0.554300	0.831211	1.145476	1.61081	2.67460	4.35146
6	0.675727	0.872085	1.237347	1.63539	2.20413	3.45460	5.34812
7	0.989266	1.239043	1.68187	2.16736	2.83311	4.25485	5.34812
8	1.344419	1.646483	2.17073	2.73264	3.48954	5.07064	7.34412
9	1.734926	2.087912	2.73039	3.32611	4.16816	5.89833	8.34351
10	2.15585	2.55821	3.24697	3.94030	4.86518	6.73720	9.34182
11	2.60321	3.05347	3.81675	4.57481	5.57779	7.58412	10.340
12	3.07382	3.57058	4.40879	5.22603	6.30390	8.43832	11.3403
13	3.56503	4.10991	5.00874	5.89189	7.04150	9.29906	12.3348
14	4.07458	4.66043	5.62872	6.57093	7.78953	10.1663	13.3303
15	4.60094	5.22936	6.26214	7.26004	8.54875	11.0355	14.3369
16	5.14224	5.81221	6.90765	7.95164	9.31223	11.9122	15.3385
17	5.69724	6.40776	7.56418	8.67176	10.0852	12.7919	16.3391
18	6.26481	7.01491	8.23075	9.39046	10.8649	13.6753	17.3399
19	6.84398	7.63273	8.90655	10.1170	11.6509	14.5620	18.3375
20	7.43386	8.26040	9.59083	10.8608	12.4420	15.4518	19.3374
21	8.03386	8.89720	10.28203	11.5913	13.2395	16.3444	20.3372
22	8.64272	9.54249	10.9823	12.3380	14.0416	17.2396	21.3370
23	9.26042	10.19567	11.6886	13.0905	14.8479	18.1373	22.3369
24	9.88623	10.8564	12.4011	13.8484	15.6587	19.0372	23.3367
25	10.5197	11.5240	13.1197	14.6114	16.4734	19.9393	24.3356
26	11.1603	12.1981	13.8439	15.3791	17.2919	20.8434	25.3354
27	11.8076	12.8786	14.5733	16.1513	18.1138	21.7494	26.3350
28	12.4613	13.5648	15.3079	16.9279	18.9392	22.6572	27.3343
29	13.1211	14.2565	16.0471	17.7083	19.7677	23.5668	28.3332
30	13.7867	14.9536	16.7908	18.4926	20.5992	24.4770	29.3320
40	20.7065	22.1643	24.4331	26.5093	28.0505	33.6657	39.3354
50	27.9907	29.7067	32.3674	34.7642	37.6895	42.9121	49.3349
60	35.5346	37.4848	40.4817	43.1879	46.4589	52.2038	59.3347
70	43.2752	45.4418	48.7576	51.7393	55.3290	61.6063	69.3344
80	51.1720	53.5400	57.1532	60.3915	64.2778	71.1445	79.3343
90	59.1963	61.7541	65.8466	69.1250	73.2912	80.8247	89.3342
100	67.3276	70.0648	74.2219	77.9288	82.3581	90.1832	99.3341

V

v	F(α)						
	.75	.90	.95	.975	.99	.995	.999
1	1.32330	2.70554	3.85146	5.02380	6.63490	7.87944	10.829
2	2.77239	4.60517	5.99147	7.37776	9.21034	10.5966	13.816
3	4.10835	6.25139	7.81473	9.34940	11.3449	12.8381	16.266
4	5.38527	7.77944	9.48773	11.1433	13.2767	14.8602	18.467
5	6.58568	9.20035	11.0705	12.8325	15.0863	16.7496	20.515
6	7.74080	10.6446	12.5916	14.4494	16.8119	18.5470	22.458
7	8.83715	12.0170	14.0671	16.0128	18.4753	20.2777	24.322
8	10.2138	13.3616	15.5073	17.5346	20.0902	21.9530	26.125
9	11.3887	14.5837	16.9190	19.0228	21.6650	23.5803	27.877
10	12.5489	15.6871	18.3070	20.4831	23.2003	25.1682	29.588
11	13.7007	16.7730	19.6761	21.9200	24.7250	26.7209	31.261
12	14.8464	17.8494	21.0261	23.3367	26.2170	28.2495	32.904
13	16.0839	18.9119	22.3621	24.7356	27.6883	29.7504	34.528
14	17.1176	20.0042	23.6848	26.1190	29.1413	31.2193	36.123
15	18.2451	22.3072	24.9959	27.4864	30.5779	32.6613	37.697
16	19.3688	23.5418	26.2952	28.8454	31.9999	34.0872	39.252
17	20.4887	24.7690	27.5871	30.1910	33.4087	35.4985	40.789
18	21.6049	25.9894	28.8693	31.5204	34.8053	36.8964	42.312
19	22.7178	27.2030	30.1435	32.8423	36.1908	38.2822	43.820
20	23.8277	28.4120	31.4104	34.1596	37.5652	39.6568	45.312
21	24.9344	29.6161	32.6705	35.4739	38.9321	41.0210	46.787
22	26.0393	30.8132	33.9244	36.7807	40.2894	42.3756	48.248
23	27.1413	32.0069	35.1726	38.0767	41.6384	43.7213	49.698
24	28.2412	33.1963	36.4161	39.3641	42.9708	45.0585	51.130
25	29.3389	34.3816	37.6525	40.6463	44.2941	46.3878	52.550
26	30.4343	35.5631	38.8852	41.9232	45.6047	47.7099	53.952
27	31.5284	36.7412	40.1133	43.1944	46.9030	49.0249	55.346
28	32.6205	37.9159	41.3372	44.4607	48.2782	50.3323	56.732
29	33.7109	39.0875	42.5569	45.7222	49.5879	51.6356	58.102
30	34.7998	40.2560	43.7729	46.9792	50.8922	52.9320	59.463
40	45.6160	51.8030	55.7586	59.3417	63.6907	66.7630	73.402
50	56.3336	63.1671	67.3048	71.4202	76.1539	79.4900	86.061
60	66.9814	74.3970	79.0819	83.2976	88.3794	91.9517	99.607
70	77.5708	85.5271	90.5312	95.0231	100.423	104.215	112.317
80	88.1303	96.5782	101.879	106.620	112.329	116.321	124.399
90	98.6499	107.565	113.145	118.136	124.116	128.299	137.208
100	109.141	118.498	124.342	129.561	135.807	140.169	149.419