

Cuprins

Prefață	vii
1 Introducere în Matlab	1
1.1 Gestionarea unei sesiuni Matlab	2
1.1.1 Lansarea și închiderea unei sesiuni Matlab	2
1.1.2 Comenzi <code>help</code>	2
1.1.3 Comenzi sistem	4
1.2 Constante. Variabile. Expresii aritmetice	4
1.2.1 Constante	5
1.2.2 Instrucțiunea <code>format</code>	5
1.2.3 Constante speciale	6
1.2.4 Variabile	7
1.2.5 Operatori aritmetici	8
1.2.6 Funcții predefinite	9
1.2.7 Comenzi pentru gestiunea variabilelor și a spațiului de lucru	10
1.3 Instrucțiuni de atribuire	12
1.4 Instrucțiuni de citire și scriere	13
1.4.1 Instrucțiunea <code>input</code>	13
1.4.2 Instrucțiunea <code>ginput</code>	13
1.4.3 Instrucțiunea <code>fprintf</code>	13
1.5 Operatori relaționali și operatori logici	14
1.5.1 Operatori relaționali	14
1.5.2 Operatori logici	14
1.6 Fișiere script (m-file)	15
1.6.1 Comenzi pentru gestionarea fișierelor	17
1.7 Grafică bidimensională	17
1.7.1 Instrucțiunea <code>plot</code>	17
1.7.2 Comenzi și instrucțiuni de gestionare a graficelor	22
1.7.3 Instrucțiunile <code>polar</code> și <code>ezpolar</code>	25

1.7.4	Instrucțiunea stairs	28
1.7.5	Instrucțiunile bar și barh	31
1.7.6	Instrucțiunea subplot	32
1.7.7	Instrucțiunea fplot	33
1.7.8	Instrucțiunea ezplot	35
1.7.9	Instrucțiunea fill	35
1.8	Instrucțiuni de ciclare și control	36
1.8.1	Instrucțiunea if	37
1.8.2	Instrucțiunea switch	38
1.8.3	Instrucțiunea while	39
1.8.4	Instrucțiunea for	40
1.9	Instrucțiuni de întrerupere	41
1.9.1	Instrucțiunea try...catch	41
1.9.2	Instrucțiunea pause	42
1.9.3	Instrucțiunea return	42
1.9.4	Instrucțiunea break	42
1.9.5	Instrucțiunea error	43
1.10	Funcții (proceduri) în Matlab	43
1.10.1	Definirea și structura unei funcții Matlab	44
1.10.2	Apelul unei funcții (proceduri) Matlab	45
1.10.3	Subfuncții	45
1.10.4	Funcția feval	47
1.10.5	Comanda echo	48
1.11	Instrucțiuni de evaluare a eficienței	48
1.11.1	Instrucțiunea flops	48
1.11.2	Instrucțiunile tic și toc	49
1.12	Grafică tridimensională	49
1.12.1	Instrucțiunea plot3	49
1.12.2	Instrucțiunea ezplot3	51
1.12.3	Instrucțiunile meshgrid, mesh și surf	51
1.12.4	Instrucțiunile contour și contourf	54
1.12.5	Instrucțiunile ezcontour și ezcontourf	56
1.12.6	Instrucțiunile ezmesh, ezsurf, ezmeshc și ezsurfc	57
1.12.7	Instrucțiunile bar3 și bar3h	58
2	Elemente de teoria probabilităților	61
2.1	Câmp de probabilitate	62
2.2	Variabile aleatoare	62
2.3	Funcție de repartiție	63
2.4	Legi de probabilitate de tip discret	64

2.4.1	Funcțiile Matlab pdf și cdf	65
2.4.2	Legi de probabilitate de tip discret clasice	67
2.5	Legi de probabilitate continue	73
2.5.1	Legi de probabilitate continue clasice	74
2.5.2	Legi de probabilitate continue statistice	81
2.5.3	Funcția Matlab normspec	87
2.5.4	Distribuție marginală	88
2.5.5	Funcție de repartiție condiționată	89
2.5.6	Funcție de supraviețuire. Funcție hazard	92
2.6	Caracteristici numerice	94
2.6.1	Valoare medie. Dispersie (varianță). Covarianță	94
2.6.2	Funcții Matlab pentru valoare medie și dispersie	99
2.6.3	Valoare medie condiționată. Dispersie (varianță) condiționată	99
2.6.4	Funcție caracteristică	103
2.6.5	Momente	103
2.6.6	Mediana. Cuartile. Cuantile	104
2.6.7	Funcția Matlab icdf	104
2.6.8	Funcția Matlab disttool	107
2.7	Șiruri de variabile aleatoare	107
2.7.1	Inegalitatea lui Cebîșev	107
2.7.2	Tipuri de convergență	108
2.7.3	Legea numerelor mari	109
2.7.4	Teoreme limită	109
3	Statistică descriptivă	113
3.1	Concepte de bază ale statisticii	114
3.2	Culegerea, prezentarea și prelucrarea datelor statistice	115
3.2.1	Generarea numerelor aleatoare în Matlab	115
3.2.2	Tabele statistice	119
3.2.3	Funcțiile tabulate și crosstab	125
3.2.4	Funcțiile caseread, casewrite, tblread și tblwrite	127
3.2.5	Reprezentări grafice	129
3.2.6	Funcții Matlab pentru reprezentarea grafică a datelor statistice	134
3.2.7	Funcția randtool	147
3.2.8	Funcțiile pie și pie3	147
3.3	Parametrii distribuțiilor statistice	148
3.3.1	Parametri statistici ce măsoară tendința	148
3.3.2	Funcțiile mean, geomean și harmean	149
3.3.3	Funcția trimmean	149
3.3.4	Funcția median	150

3.3.5	Parametri statistici ce măsoară dispersarea	150
3.3.6	Funcția prctile	151
3.3.7	Funcția moment	155
3.3.8	Funcțiile var și std	155
3.3.9	Funcțiile range, iqr, mad, skewness și kurtosis . .	155
3.3.10	Funcțiile max și min	156
3.3.11	Funcțiile sort și sortrows	156
3.3.12	Funcțiile sum, prod, cumsum și cumprod	157
3.3.13	Funcția diff	157
3.3.14	Funcțiile nanmean, nanmedian, nanstd, nanmin, nanmax și nansum	158
3.3.15	Corecțiile lui Sheppard	158
3.3.16	Funcția grpstats	163
3.3.17	Funcția boxplot	164
3.4	Corelație și regresie	165
3.4.1	Funcțiile cov și corrcoef	179
3.4.2	Funcțiile lsline, reflin și gline	181
3.4.3	Funcțiile polyfit și refcurve	182
3.4.4	Funcțiile polyval, polyvalm, poly și roots	183
3.4.5	Funcția polytool	184
3.4.6	Funcția nlinfit	185
3.4.7	Funcția nlintool	186
3.4.8	Coeficienții Spearman și Kendall	187
4	Teoria selecției	197
4.1	Tipuri de selecție	197
4.2	Funcții de selecție	198
4.2.1	Media de selecție	199
4.2.2	Momente de selecție	201
4.2.3	Coeficient de corelație de selecție	207
4.2.4	Funcție de repartiție de selecție	216
4.2.5	Funcția cdfplot	223
5	Teoria estimăției	229
5.1	Funcție de verosimilitate	229
5.2	Funcții de estimăție	236
5.2.1	Funcții de estimăție absolut corecte	237
5.2.2	Funcții de estimăție corecte	239
5.2.3	Funcții de estimăție eficiente	241
5.2.4	Estimatori optimali	247

5.3	Metode pentru estimarea parametrilor	251
5.3.1	Metoda momentelor	251
5.3.2	Metoda verosimilității maxime	252
5.3.3	Metoda minimului χ^2	256
5.3.4	Metoda intervalelor de încredere	257
5.3.5	Metoda intervalelor de încredere pentru selecții mari	275
5.3.6	Funcții Matlab privind estimația	277
6	Verificarea ipotezelor statistice	285
6.1	Concepte de bază	285
6.2	Testul Z privind media teoretică	286
6.2.1	Funcțiile $zscore$ și $ztest$	290
6.3	Puterea unui test	291
6.4	Testul T (Student) privind media teoretică	304
6.4.1	Funcția $ttest$	307
6.5	Testul raportului verosimilităților	308
6.6	Testul χ^2 privind dispersia teoretică	312
6.7	Testul F (Fisher–Snedecor) pentru compararea dispersiilor	318
6.8	Teste pentru compararea mediilor	322
6.8.1	Dispersii cunoscute	323
6.8.2	Dispersii egale necunoscute	325
6.8.3	Funcția $ttest2$	327
6.8.4	Dispersii diferite necunoscute	328
6.8.5	Observații perechi	332
6.9	Testul χ^2 pentru concordanță	334
6.10	Testul χ^2 pentru compararea mai multor caracteristici	346
6.11	Testul χ^2 pentru tabele de contingență	349
6.12	Testul de concordanță al lui Kolmogorov	353
6.12.1	Funcția $kstest$	360
6.12.2	Funcția $lillietest$	362
6.13	Testul Kolmogorov–Smirnov	363
6.13.1	Funcția $kstest2$	366
	Anexa I	369
	Anexa II	370
	Anexa III	371
	Anexa IV	372

Anexa V	376
Bibliografie	377
Index	383

Prefață

Aspecte primare ale statisticii se pierd în negura timpului. Astăzi, putem spune, fără a greși, că statistică facem fiecare dintre noi, cu știință sau fără. Gândul la ziua de mâine, cum ne organizăm și cum folosim timpul și mijloacele materiale de care dispunem, are la bază date (statistice) reținute (culese) și prelucrate în mod simplu sau mai complex. Ne dăm seama, prin urmare, că statistica are o largă răspândire, de la utilizarea ei în mod empiric, până la apelarea la metodele fundamentate matematic.

Statistică făceau chinezii antici, care dispuneau de date relative la populație, pământuri și recolte, dar și egiptenii, care efectuau cadastrări (operații de determinare a unor proprietăți agricole și imobiliare cu toate caracteristicile lor) și numărări de populație. Numărarea populației era cunoscută și la evrei, de exemplu, în texte biblice sunt menționate numărări de populație, ca cea efectuată de Moise, privind bărbații buni de oaste. La grecii antici erau, de asemenea, efectuate numărări de populație, ale bunurilor necesare pentru scopuri militare, pentru așezarea impozitelor și evaluarea bogățiilor. Toate acestea au culminat cu recensămintele și anchetele romane.

Culegerea de date privind resursele umane și materiale aveau la bază o gândire practică, folosirea acestora în scopuri fiscale, militare sau administrative. Putem spune că toate acestea erau utilizate în descrierea *statului*.

Tratarea științifică a datelor relative la descrierea statului se întâlnește în Germania (sec. XVII–XVIII), iar denumirea de *statistică* apare în cursul lui Martin Schmeitzel (1679–1757), intitulat "Collegium politico-statisticum" (Universitatea din Halle). Unii istorici atribuie întâietate privind această denumire lui Gottfried Achenwall (1719–1772), care a introdus învățământul "Staatskunde" la Universitatea din Göttingen.

În Anglia, în aceeași perioadă, în afara universităților, exista ca disciplină descriptivă a statului ceea ce se numea *aritmetică politică*. Aritmetica politică se ocupa în special de cercetarea fenomenelor demografice.

Un moment important în dezvoltarea statisticii îl reprezintă cristalizarea calculului probabilităților.

Calculul probabilităților și fundamentarea riguroasă din punct de vedere matema-

tic a teoriei probabilităților reprezintă o problemă fundamentală în stabilirea locului și rolului acestei discipline matematice în evantaiul larg al matematicii, precum și al fundamentării teoretice a statisticii matematice.

Desigur, faptul că cineva, cu diferite ocazii, folosește cuvântul *probabil*, nu înseamnă că este o persoană inițiată în calculul probabilităților. Am putea, eventual, accepta că sunt percepute anumite aspecte empirice ale calculului probabilităților, când, relativ la un anumit fenomen care prezintă un anumit grad de nedeterminare, se emit afirmații de forma: *este puțin probabil ca...*, *este foarte probabil ca...*, *este improbabil ca...*. De la astfel de afirmații și până la cunoașterea și înțelegerea în profunzime a teoriei probabilităților este un drum lung, care este de un larg interes atât din punct practic cât și teoretic, și care prezintă o atracție deosebită, nu numai pentru matematicieni, dar și pentru specialiști din alte discipline ale cunoașterii umane: fizică, chimie, biologie, economie, medicină, inginerie etc. Aceste aspecte, de-a lungul timpului, se regăsesc în apariția, dezvoltarea și fundamentarea axiomatică a teoriei probabilităților, care se îmbină și se continuă cu diferite domenii de aplicabilitate, dintre care statistica ocupă un rol privilegiat.

Este unanim acceptat că noțiunea de probabilitate își are originea în jocurile de noroc. În antichitate, la egipteni, greci și romani erau cunoscute jocuri de noroc cu zaruri, care s-au perpetuat până în zilele de astăzi. Se pare că cel mai vechi zar datează din jurul anilor 3500 î.e.n. și avea fețele numerotate la întâmplare de la 1 la 6, față de zarul actual, care are fețele numerotate astfel ca suma numerelor de pe fețele opuse să fie 7. Proliferarea jocurilor de noroc s-a produs nu numai prin aria de răspîndire, dar și prin diversitatea acestora. Ajunge să amintim doar jocurile de ruletă care sunt instalate în marile metropole ale lumii.

Problemele teoretice relative la jocurile de noroc, în special privind șansele parțenerilor, au atras atenția savanților timpului. Astfel, fundamentarea elementară a calculului probabilităților poate fi atribuită francezilor P. de Fermat (1601–1665) și B. Pascal (1623–1662). Corespondența din anul 1654, purtată de cei doi, conține primele aspecte privind calculul șanselor în jocurile de noroc. O parte din problemele considerate au fost formulate și prezentate lui Pascal de către Cavalerul de Méré, un renumit amator al jocurilor de noroc din acea vreme.

Prima carte de introducere în calculul probabilităților este *De Ratiociniis in Aleae Ludo* (*Cum se raționează în jocurile cu zarul*), apărută în anul 1657 și are ca autor pe olandezul Ch. Huygens (1629–1695).

Totuși, se remarcă faptul că italianului G. Cardano (1501–1576) i se datorează lucrarea intitulată *Liber de Ludo Aleae* (*Cartea despre jocurile cu zarul*), care a fost cunoscută și publicată la mai bine de o sută de ani de la dispariția sa.

Un destin asemănător a avut și lucrarea lui James Bernoulli (1654–1705), *Ars Conjectandi* (*Arta de a face presupuneri*), care a fost publicată de către fratele său

John Bernoulli în anul 1713. Se remarcă în conținutul cărții, ceea ce se numește azi legea numerelor mari, unul din rezultatele strălucitoare ale teoriei probabilităților.

Un alt rezultat de aceeași dimensiune, teorema limită centrală, se găsește în cartea *Doctrine of Chances*, apărută în anul 1718 și care este opera lui A. de Moivre (1667–1754), ceea ce, se cunoaște, este strâns legată de legea normală de probabilitate.

Descoperirea legii normale de probabilitate este atribuită, de cele mai multe ori, lui K. F. Gauss (1777–1855), și care este considerată ca fiind specifică erorilor de măsurare. Trebuie remarcat faptul că în anul 1808, cu un an înainte de descoperirea lui Gauss, topograful american Robert Adrain (1775–1843) publicase o lucrare în care a sugerat utilizarea legii normale pentru descrierea erorilor de măsurare.

Rezultatele lui Moivre, privind teorema limită centrală, sunt extinse de către P. - S. de Laplace (1770–1820), stabilind de asemenea legătura cu legea normală de probabilitate, în cartea *Théorie Analytique des Probabilités* (*Teoria analitică a probabilităților*), apărută în anul 1812.

Un statistician belgian de seamă al sec. XIX, unanim recunoscut pentru rezultatele de pionerat în domeniul statisticii matematice este A. Quetelet (1796–1874). De numele lui se leagă noțiuni ale statisticii ca: repartiție, medie, dispersie, observare de masă, regularitate. Pentru el, statistica reprezenta singura metodă ce se poate aplica fenomenelor de masă.

Un salt remarcabil în dezvoltarea teoriei probabilităților a fost făcut prin contribuția matematicienilor ruși L. P. Cebîșev (1821–1884), A. A. Markov (1856–1922) și A. M. Leapunov (1857–1918). Primul a introdus noțiunea de variabilă aleatoare, reformulând și generalizând legea numerelor mari în limbajul variabilelor aleatoare. Folosind noțiunea de variabilă aleatoare, A. A. Markov și A. M. Leapunov au obținut apoi alte rezultate privind legea numerelor mari și teorema limită centrală.

Pasul decisiv privind fundamentarea modernă a teoriei probabilităților este făcut prin lucrarea matematicianului rus A. N. Kolmogorov (1903–1989), apărută în anul 1933. A. N. Kolmogorov, folosind teoria măsurii, reușește să construiască modelul axiomatic al teoriei probabilităților.

Se remarcă de asemenea preocupări și rezultate importante obținute de alți matematicieni ca: S. D. Poisson (1781–1840), É. Borel (1871–1956), S. N. Bernstein (1880–1968), R. von Mises (1883–1953), A. I. Hincin (1894–1959).

Sfârșitul sec. XIX și începutul sec. XX se consideră a fi începutul statisticii moderne. Este momentul când se trece de la etapa descriptivă, la interpretarea analitică a fenomenelor de masă și obținerea de concluzii inductive pe baza observațiilor empirice. Statisticieni de renume, care au pus bazele statisticii matematice sunt considerați a fi englezii F. Galton (1822–1911) și K. Pearson (1857–1936). Opera acestora este continuată de celebrul statistician R. A. Fisher (1890–1962). Dintre reprezentanții școlii engleze de statistică amintim pe G. U. Yule (1871–1951), W. S. Gosset (1876–

1937), C. E. Spearman (1883–1945), M. Kendall (1907–1983).

Pentru detalii privind istoricul statisticii recomandăm lucrarea [38].

Astăzi, statistica matematică, având un suport teoretic al teoriei probabilităților, are o dezvoltare și aplicabilitate deosebită. Care este motivul aplicării statisticii matematice pe o scară largă și în atâtea domenii ale cunoașterii umane? Am putea da două răspunsuri principale. În primul rând, fiecare cercetător sau persoană care analizează fenomene ale lumii reale, modelează astfel de fenomene, caută să pună la baza metodelor sale de cercetare științifică, de evaluare a fenomenelor reale, un instrument de investigare obiectiv. Și este cvasiunanim recunoscut că matematica are astfel de instrumente de investigare, inclusiv prin metodele statisticii matematice. În al doilea rând, metodele statisticii matematice sunt ușor de aplicat. Greutatea constă numai în fundamentarea riguroasă din punct de vedere matematic a acestora. Acest aspect este ascuns celui care aplică statistica matematică. În materialul pe care îl prezentăm, considerăm noțiunile (conceptele) de bază în inițierea studiului și aplicării statisticii matematice. Se consideră partea “ascunsă” a statisticii matematice, adică fundamentarea teoretică, cât de cât riguroasă, a statisticii matematice, dar și partea care prezintă modul de aplicare a metodelor statisticii matematice.

Dacă în tot acest demers avem la dispoziție și mijlocul modern de lucru, reprezentat astăzi printr-un calculator performant sau mai puțin performant, atunci statistica devine foarte atractivă. De aceea au fost elaborate o serie de produse soft, care prezintă într-o concepție unitară reprezentarea, prelucrarea și studiul datelor statistice. Amintim aici câteva astfel de produse soft: *Statgraphics*, *Statistica*, *S-Plus*, *SAS*, *StatXact*, *SPSS*, *Stata*, *GraphPad*, *ViSta* etc. Remarcăm de asemenea că unele produse soft elaborate în alte scopuri decât cel al prelucrării datelor statistice, conțin proceduri privind prelucrarea mai elementară sau mai complexă a datelor statistice. În această categorie se încadrează și produsul *Matlab*, care va fi invocat în prezenta lucrare, dar tot în această categorie putem aminti și *Mathematica*, *Maple* etc.

Capitolul 1

Introducere în Matlab

MATLAB (MATrix LABoratory) este un sistem interactiv de nivel și performanță înalte pentru efectuarea de calcule științifice și tehnice, precum și de analiză și vizualizare a datelor și care dispune de un limbaj propriu de programare.

Prima versiune a fost realizată în anii 70 de către *Cleve Moler*, iar versiunea utilizată în cele ce urmează este *Matlab 6.0*.

Se pare că acest produs are o largă răspândire din mai multe motive:

- interactivitatea sistemului;
- elementele de bază sunt matricele, după cum spune și denumirea, ale căror tipuri și dimensiuni nu se declară;
- grafica puternică și foarte ușor de utilizat;
- posibilitatea integrării unor proceduri proprii în limbaje de programare cum ar fi *Fortran*, *C*;
- arhitectura modulară, care permite adăugarea unor proceduri și funcții proprii, ceea ce a permis deja elaborarea unor pachete de proceduri și funcții, numite *toolbox*-uri, orientate pe anumite domenii, cum ar fi
 - Statistica – *Statistics Toolbox*;
 - Funcții spline – *Spline Toolbox*;
 - Analiză wavelet – *Wavelet Toolbox*;
 - Procesarea semnalelor – *Signal Processing Toolbox*;
 - Procesarea imaginilor – *Image Processing Toolbox*;
 - Modelare, simulare și analiza sistemelor dinamice – *Simulink*;

- Calcul simbolic – *Symbolic/Extended Symbolic Math Toolbox*;
- Ecuații cu derivate parțiale – *PDE Toolbox*;
- Optimizare – *Optimization Toolbox*.

1.1 Gestionarea unei sesiuni Matlab

1.1.1 Lansarea și închiderea unei sesiuni Matlab

Trebuie să remarcăm de la început că sistemul Matlab dispune de o puternică bază de *help*-uri, dar care poate fi activată doar după ce sistemul Matlab a fost lansat.

Lansarea sistemului Matlab se efectuează prin comenzi ce sunt specifice sistemului de operare pe care este instalat. Vom considera în continuare sistemul de operare Windows, caz în care lansarea se face prin activarea

- *icon*-ului *Matlab* sau
- *aplicației Matlab* din directorul în care a fost instalat sistemul Matlab.

În urma efectuării acestei comenzi se deschide o fereastră în care apare prompterul specific sistemului Matlab

```
>>
```

Sistemul Matlab devine în acest fel interactiv, adică la fiecare comandă sau funcție tastată și acceptată de Matlab (editarea comenzii sau funcției încheindu-se prin *Enter*), sistemul o execută și afișează pe ecran rezultatul, dacă este cazul.

Există posibilitatea de trecere de la modul interactiv la cel programat, prin tastarea numelui unui fișier (cu extensia *.m*), care conține un program în sistemul Matlab (succesiune de instrucțiuni Matlab) și care mod se încheie odată cu terminarea executării programului respectiv, după care se revenie la modul interactiv.

Încheierea unei sesiuni Matlab se poate face, fie prin comenzi specifice sistemului de operare Windows, fie prin tastarea uneia din comenzile

```
>>quit  
>>exit
```

1.1.2 Comenzi *help*

Comanda *help* permite utilizatorilor să aibă acces la o mare cantitate de informație (în limba engleză) relativă la întreg sistemul Matlab.

Există diferite variante pentru lansarea comenzii *help*, anume prin:

- tastarea comenzii *help* urmată de domeniul (*topic*) în care suntem interesați, adică

```
>>help topic
```

Remarcăm faptul că parametrul `topic` poate lipsi. Astfel, lansarea comenzii

```
>>help
```

are ca efect afișarea pe ecran a unei liste ce cuprinde domeniile care operează în versiunea curentă Matlab

HELP topics:

matlab\general	-	General purpose commands.
matlab\ops	-	Operators and special characters.
matlab\lang	-	Programming language constructs.
matlab\elmat	-	Elementary matrices and matrix manipulation.
matlab\elfun	-	Elementary math functions.
matlab\specfun	-	Specialized math functions.
matlab\matfun	-	Matrix functions - numerical linear algebra.
matlab\datafun	-	Data analysis and Fourier transforms.
matlab\audio	-	Audio support.
matlab\polyfun	-	Interpolation and polynomials.
matlab\funfun	-	Function functions and ODE solvers.
matlab\sparfun	-	Sparse matrices.
matlab\graph2d	-	Two dimensional graphs.
matlab\graph3d	-	Three dimensional graphs.
matlab\specgraph	-	Specialized graphs.
matlab\graphics	-	Handle Graphics.
matlab\uitools	-	Graphical user interface tools.
matlab\strfun	-	Character strings.
matlab\iofun	-	File input/output.
matlab\timefun	-	Time and dates.
matlab\datatypes	-	Data types and structures.
matlab\verctrl	-	Version control.
matlab\winfun	-	Windows Operating System Interface Files (DDE/ActiveX)
matlab\demos	-	Examples and demonstrations.
toolbox\local	-	Preferences.
toolbox\stats	-	Statistics Toolbox.

For more help on directory/topic, type "help topic".

- tastarea comenzii `lookfor` urmată de un cuvânt cheie (`key`), în care suntem interesați, adică

```
>>lookfor key
```

va lista prima linie de comentariu din fiecare fișier cu extensia `.m`, care conține cuvântul cheie `key`, precedată de numele fișierului.

- tastarea comenzii `helpwin` urmată de un domeniu (`topic`) în care suntem interesați, adică

```
>>helpwin topic
```

produce deschiderea unei ferestre cu informații relative la domeniul precizat.

- tastarea comenzii `helpdesk`, adică

```
>>helpdesk
```

produce deschiderea unei ferestre prin care se permite căutarea în baza `help`-ului Matlab.

1.1.3 Comenzi sistem

Sistemul Matlab operează cu o mulțime de comenzi pe care le întâlnim și în alte sisteme. Enumerăm o parte dintre acestea:

- `mkdir` – deschiderea unui subdirector nou în directorul curent;
- `cd` – schimarea directorului curent;
- `dir` – lista conținutului directorului curent;
- `version` – afișarea versiunii sistemului Matlab instalat;
- `!com` – ieșire temporară din sistemului Matlab instalat și execuția comenzii Windows precizată prin `com`;
- `demo` – lansarea execuției de programe demonstrative.

Înainte de începerea unei sesiuni este de dorit să se creeze un director prin comanda `mkdir`, pentru salvarea fișierelor, iar apoi trecerea în acest director prin comanda `cd`.

1.2 Constante. Variabile. Expresii aritmetice

Având în vedere că sistemul Matlab operează numai cu un singur tip de obiect, matrice numerică, care ar putea avea elemente numere complexe, rezultă că orice scalar este considerat ca o matrice cu o linie și o coloană. În plus, reprezentarea internă a elementelor matricelor și operațiile cu acestea se fac în dublă precizie. Deoarece reprezentarea internă este prestabilită, avem doar posibilitatea de precizare a tipului (formatului) de afișare (ieșire) respectiv de intrare a datelor.

1.2.1 Constante

Introducerea constantelor se poate face:

- de la tastatură;
- prin generarea cu ajutorul unor instrucțiuni Matlab;
- dintr-un fișier creat cu ajutorul unui editor.

La intrare constantele numerice pot fi de tipul:

- *întreg*: 2002, -2002;
- *real*: 3.14159, 0.314e1, -.314e+01, 31.4e-1, 31.4e-01;
- *complex*: a+bi, unde a și b sunt două constante întregi sau reale, iar i reprezintă unitatea imaginară ($i = \sqrt{-1}$).

1.2.2 Instrucțiunea format

Sistemul Matlab dispune de câteva șabloane (formatări) standard pentru afișarea (extragera) datelor numerice. Instrucțiunea `format` prin care se precizează formatul standard are sintaxa

```
>>format tip
```

unde `tip` poate lua una din valorile `short`, `long`, `short e`, `long e`, `short g`, `long g`, `hex`, `+`, `bank`, `rat`, `loose`, `compact`.

Formatele `short` și `long` afișează pe ecran datele numerice în virgulă fixă, respectiv cu 5 cifre zecimale și 15 cifre zecimale (precedate de semnul minus, dacă sunt negative).

Formatele `short e` și `long e` afișează pe ecran datele numerice în virgulă flotantă, respectiv cu 5 cifre zecimale și 15 cifre zecimale (precedate de semnul minus, dacă sunt negative).

Formatele `short g` și `long g` afișează pe ecran datele numerice cu cea mai potrivită reprezentare, virgulă fixă sau virgulă flotantă, respectiv cu 5 cifre zecimale și 15 cifre zecimale (precedate de semnul minus, dacă sunt negative).

Formatele `rat`, `bank`, `hex` și `+` afișează pe ecran respectiv datele numerice sub formă de fracție ordinară, sub formă bancară (cu două zecimale), sub formă hexazecimală, iar în ultimul caz, prin `+` dacă numărul este pozitiv, respectiv `-`, dacă numărul este negativ și spațiu dacă este zero.

Formatele `loose` și `compact` permite afișare datelor la două rânduri (format implicit), respectiv la un rând în cazul al doilea.

Pentru a afla la un moment dat, care format este în funcție se poate apela la comanda

```
>>get(0,'Format')
```

1.2.3 Constante speciale

Sistemul Matlab are câteva constante speciale, care au denumiri specifice (identificatori specifici):

- `inf` – rezultatul împărțirii la 0;
- `NaN` (Not a Number) – rezultatul operației de tipul 0/0;
- `realmax` – cel mai mare număr real pozitiv reprezentat pe dublu cuvânt ($1.7977e+308$);
- `realmin` – cel mai mic număr real pozitiv reprezentat pe dublu cuvânt ($2.2251e-308$);
- `pi` – $\pi=3.14159\dots$ (valoarea reprezentată pe dublu cuvânt);
- `eps` – epsilonul mașinii ($eps=2.2204e-16=2^{-54}$);
- `date` – data curentă;
- `clock` – timpul curent;
- `calendar` – calendarul lunii curente;
- `i`, `j` – unitatea imaginară ($i, j=\sqrt{-1}$);
- `zeros(n)` și `zeros(m,n)` – generează matricea pătratică de ordin n , respectiv matricea cu m linii și n coloane, având toate elementele nule;
- `ones(n)` și `ones(m,n)` – generează matricea pătratică de ordin n , respectiv matricea cu m linii și n coloane, având toate elementele 1;
- `eye(n)` și `eye(m,n)` – generează matricea pătratică de ordin n , respectiv matricea cu m linii și n coloane, având pe diagonala principală 1 și 0 în rest (când matricea este pătratică avem matricea unitate);
- `magic(n)` – generează matrice magică pătratică de ordin n .

Pentru introducerea elementelor unei matrice, linie cu linie, elemente ce pot fi exprimate și prin expresii aritmetice, acestea sunt incluse între paranteze drepte. Elementele fiecărei linii sunt separate prin spațiu sau virgulă, iar liniile sunt separate prin punct-virgulă (;) sau prin trecerea la rândul următor.

Astfel matricea

$$\begin{pmatrix} 3 & 5 & 2 \\ 9 & 1 & 4 \end{pmatrix}$$

se poate introduce de la tastatură prin


```
>>[3 5 2;9 1 4]
```

sau

```
>>[3,5,2;9,1,4]
```

Dacă elementele matricei prezintă regularitatea, că elementele unei linii sunt date printr-o progresie aritmetică, atunci avem o introducere prescurtată. De exemplu, introducerea matricei

$$\begin{pmatrix} 3 & 5 & 7 & 9 \\ 9 & 6 & 3 & 0 \end{pmatrix}$$

se poate face de la tastatură prin

```
>>[3:2:9;9:-3:0]
```

1.2.4 Variabile

Variabilele sunt matrice ce sunt utilizate într-un program Matlab sau în mod interactiv și ale căror valori pot fi modificate. Numele lor fiind precizate prin *identificatori*, care sunt dați de orice combinație de litere, cifre zecimale și liniuța de subliniere, din care primul caracter este literă: Matrice1, mat, Cluj_Napoca.

Apelarea și vizualizarea unui element al unei variabile (matrice) se face prin specificarea numelui variabilei urmat de lista indicilor corespunzători, separați prin virgulă, închisă între paranteze. De exemplu, $A(i, j)$ pune în evidență elementul din linia i și coloana j al matricei A .

Facem observația că matricele în sistemul Matlab nu acceptă decât indici strict pozitivi.

Utilizând operatorul două puncte ($:$) se pot obține submatrice ale unei matrice. De exemplu, $A(1:3, 4:5)$ extrage submatricea formată din primele trei linii și coloanele 4 și 5 ale matricei A , iar $A(:, 4)$ și $A(1, :)$ extrag respectiv coloana a patra a matricei A și prima linie a matricei A . O formă mai generală de obținere a unei submatrice se realizează prin apelul de forma $A([l], [c])$, unde l și c reprezintă liste de indici de linie și indici de coloane, separați prin virgulă sau spațiu. Astfel va fi obținută submatricea formată din liniile și coloanele precizate prin l , respectiv c ale matricei A . Dacă se introduce comanda $A(:)$, atunci matricea A este transformată într-un vector coloană, ce conține coloanele matricei.

Alte combinații ale acestor tipuri de specificare ale elementelor unei matrice pot fi considerate.

Mai remarcăm aici comanda

```
>>reshape(A,m,n)
```

care transformă matricea A într-o matrice cu m linii și n coloane. Dacă A nu are $m \times n$ elemente, atunci apare eroare la executare.

1.2.5 Operatori aritmetici

Expresiile aritmetice se construiesc după cunoscutele reguli ale limbajelor de programare evolute, folosind operatorii aritmetici cunoscuți, desigur cu unele operații speciale specifice sistemului Matlab.

Ordinea execuțiilor operațiilor aritmetice este de asemenea cea cunoscută și pe care o reamintim aici, începând cu cel mai înalt nivel :

- prima dată sunt executate operațiile aritmetice dintre paranteze;
- urmează ridicarea la putere, iar dacă sunt mai multe operații consecutive de acest tip, atunci ordinea execuției este de la dreapta la stânga;
- următorul nivel de prioritate include operațiile de înmulțire și împărțire, iar dacă sunt mai multe operații consecutive de acest tip, atunci ordinea execuției este de la stânga la dreapta;
- ultimul nivel de prioritate conține operațiile de adunare și scădere, iar dacă sunt mai multe operații consecutive de acest tip, atunci ordinea execuției este de la stânga la dreapta.

Operatorii aritmetici matriceali de care dispune sistemul de bază Matlab sunt:

- $A+B$, $A-B$ sau `plus(A,B)`, `minus(A,B)` – adună, respectiv face diferența matricelor A și B de aceleași dimensiuni, iar dacă unul dintre operanzi e un scalar, acesta este extins la matricea constantă dată prin scalarul respectiv, având dimensiunile celeilalte matrice, după care se efectuează operația;
- $A*B$ sau `mtimes(A,B)` – înmulțește matricele A și B , pentru care numărul coloanelor matricei A este același cu numărul liniilor matricei B , iar dacă una din matrice e un scalar, acesta va înmulți toate elementele celeilalte matrice;
- A^B sau `mpower(A,B)` – dacă B este un scalar p întreg pozitiv, atunci se obține puterea p a matricei pătratică A , iar dacă A este un scalar α , atunci se ridică scalarul α la puterea B , matricea B fiind pătratică, adică, pentru exemplificare, dacă $\alpha = e$, atunci $e^B = \sum_{n=0}^{\infty} \frac{B^n}{n!}$;
- $A.*B$, $A.^B$, $A.\backslash B$, $A./B$, sau formele echivalente, respectiv `times(A,B)`, `power(A,B)`, `ldivide(A,B)`, `rdivide(A,B)` – efectuează operațiile de înmulțire, de ridicare la putere, împărțire inversă, împărțire directă, toate de tipul element cu element (ca la operațiile de adunare și scădere), pentru matricele de aceleași dimensiuni, iar dacă una dintre matrice e un scalar, aceasta este extinsă la matricea constantă dată prin scalarul respectiv, având dimensiunile celeilalte matrice, după care este efectuată operația;

- $A \setminus B$ sau `mldivide(A,B)` – are ca rezultat matricea ce reprezintă soluția ecuației matriceale $A * X = B$, când matricea pătratică A este inversabilă, iar dacă B este un vector coloană se obține soluția sistemului liniar corespunzător;
- A / B sau `mrdivide(A,B)` – are ca rezultat matricea ce reprezintă soluția ecuației matriceale $X * B = A$, când matricea pătratică B este inversabilă;
- $[A \ B]$, $[A ; B]$ (*concatenarea*) – produce concatenarea matricelor A și B pe orizontală, respectiv pe verticală, cu condiția ca la concatenarea pe orizontală trebuie ca numărul liniilor celor două matrice să fie același, iar la concatenarea pe verticală se impune ca numărul coloanelor celor două matrice să fie același.

1.2.6 Funcții predefinite

Sistemul Matlab dispune de numeroase funcții matematice. Remarcăm faptul că toate acestea, având în vedere că matricele sunt obiectele de bază în Matlab, acționează asupra matricelor.

Încercăm o clasificare a acestora:

- `abs` (modul), `sqrt` (radical), `sign` (semn), `conj` (conjugat), `angle` (argument al numărului complex), `imag` (parte imaginară), `real` (parte reală);
- `sin` (sinus), `cos` (cosinus), `tan` (tangentă), `cot` (cotangentă), `csc` (cosecantă), `sec` (secantă);
- `asin` (arcsinus), `acos` (arccosinus), `atan` (arctangentă), `acot` (arccotangentă), `acsc` (arccosecantă), `asec` (arcsecantă), `sinh` (sinus hiperbolic), `cosh` (cosinus hiperbolic), `tanh` (tangentă hiperbolică), `coth` (cotangentă hiperbolică), `csch` (cosecantă hiperbolică), `sech` (secantă hiperbolică), `asinh` (arcsinus hiperbolic), `acosh` (arccosinus hiperbolic), `atanh` (arctangentă hiperbolică), `acoth` (arccotangentă hiperbolică), `acsch` (arccosecantă hiperbolică), `asech` (arcsecantă hiperbolică);
- `exp` (exponențială), `log` (logaritm natural), `log10` (logaritm zecimal), `log2` (logaritm în baza doi);
- `round` (rotunjire la cel mai apropiat întreg), `floor` (rotunjire spre $-\infty$), `fix` (rotunjire spre zero), `ceil` (rotunjire spre $+\infty$).

Facem observația că pentru funcțiile trigonometrice se impune ca argumentele să fie date în radiani.

Dacă argumentul funcțiilor predefinite este matrice, atunci funcția se aplică fiecărui element al matricei.

- `atan2(A,B)` (are același efect cu `atan(A.\B)`), `rem(A,B)` (restul împărțirii lui A la B).

Dacă A și B sunt matrice, atunci trebuie să fie de aceeași dimensiune, iar funcția se aplică element cu element.

- `size(A)` (dimensiunea matricei A), `length(A)` (lungimea vectorului A), `ndims(A)` (numărul dimensiunilor matricei A), `disp(A)` (afișează elementele matricei A fără numele matricei), `det(A)` (determinantul matricei pătratice A), `inv(A)` (inversa matricei pătratice A), `rank(A)` (rangul matricei A), `A'` sau `transpose(A)` (transpusa matricei A, dacă matricea este cu elemente numere reale, dacă are elemente numere complexe, atunci se efectuează matricea transpusă a conjugatelor elementelor), `A.''` (transpusa matricei A), `diag(A)` (diagonala matricei pătratice A, iar dacă A este vector, generează matricea pătratică diagonală cu elementele diagonalei date de vector), `trace(A)` (urma matricei pătratice A, adică $\sum_i a_{ii}$).

Trebuie să remarcăm că dacă A este un vector linie, atunci `A'` este un vector coloană și invers, dacă A este un vector coloană, atunci `A'` este un vector linie.

- `factorial(n)` (factorialul lui n), `perms(1:n)` sau `perms(v)` (generează toate permutările primelor n numere naturale, respectiv generează toate permutările componentelor vectorului v), `nchoosek(n,k)` sau varianta `nchoosek(v,k)` (calculează numărul combinărilor de n elemente luate câte k, respectiv generează toate combinările elementelor vectorului v luate câte k), `combnk(v,k)` (generează toate combinările elementelor vectorului v luate câte k), `factor(n)` (factorii primi ai lui n repetați de atâtea ori cât reprezintă puterea acestora), `primes(n)` (numerele prime mai mici decât n), `gcd(cmmdc)`, `lcm(cmmmc)`.

1.2.7 Comenzi pentru gestiunea variabilelor și a spațiului de lucru

Sistemul Matlab se poate lansa din orice director, acesta devenind directorul curent. Desigur este de dorit ca utilizatorul să lucreze într-un director propriu. Așa cum am văzut, un astfel de director poate fi creat până nu se intră în sistemul Matlab, ca apoi acesta să fie lansat din acest director sau se poate lansa la început sistemul Matlab, iar apoi prin comenzile `cd` și `mkdir` să se ajungă într-un director al utilizatorului sau să se creeze un director propriu nou.

Toate operațiile sunt executate în directorul curent.

Dacă se doresc informații relative la variabilele existente la un moment dat în directorul curent se pot lansa una din comenzile:

- `who` – produce lista variabilelor curente;
- `whos` – produce lista variabilelor curente, împreună cu informații privind aceste variabile (dimensiune, număr de elemente, memorie ocupată, tipul);
- `clear` – șterge variabilele și funcțiile temporare din directorul curent (ștergerea unei singure variabile se realizează prin comanda `clear` urmată de numele variabilei);
- `diary file` – salvează toate informațiile apărute pe ecran în timpul derulării unei sesiuni în fișierul `file`, cu excepția graficelor;
- `save` – salvează variabilele de lucru pe disc;
- `load` – încarcă variabilele de lucru de pe disc.

Comenzile

```
>>save file  
>>save file x y z
```

vor salva în fișierul `file.mat` toate variabilele din sesiune curentă, respectiv numai variabilele `x`, `y`, `z`.

Diferența dintre comanda `diary` și comanda `save` este că fișierul obținut prin `diary` se poate modifica (edita) cu ajutorul unui editor, pe când celalalt nu se poate edita.

De asemenea, remarcăm faptul că în fișierul creat prin comanda `diary` conține toate informațiile până la întâlnirea comenzii

```
>>diary off
```

Dacă înainte de închiderea unei sesiuni Matlab nu se execută o comandă `save`, atunci toate variabilele împreună cu conținutul lor se pierd.

Într-o altă sesiune Matlab, se pot recupera variabilele din sesiunea precedentă, dacă au fost salvate prin comanda `save`, prin lansarea comenzii `load`. Anume, dacă fișierul salvat este `file.mat` atunci comanda

```
>>load file
```

va încărca variabilele cu conținutul lor pentru sesiunea curentă.

Comanda `load` poate fi utilizată și pentru încărcarea unui fișier ASCII creat cu alt editor. Astfel comanda

```
>>load A.dat -ascii
```

va avea ca efect pentru sesiunea curentă a generării unei variabile cu numele `A` și care conține ca valori datele din fișierul `A.dat`.

1.3 Instrucțiuni de atribuire

Sistemul Matlab dispune de trei tipuri de instrucțiuni de atribuire:

```
expresie
variabila = expresie
[lista de variabile] = func
```

Expresia se compune, după reguli cunoscute în programare, cu ajutorul operațiilor aritmetice și al funcțiilor, prin operarea asupra constantelor și asupra variabilelor.

Dacă se consideră prima formă a instrucțiunii de atribuire, atunci rezultatul este păstrat într-o variabilă de lucru numită `ans` și care rămâne nemodificată până la execuția unei alte instrucțiuni de atribuire de această formă. Forma a doua face această atribuire, variabilei din partea stângă a semnului de egalitate.

Forma a treia este specifică pentru cazul în care se apelează o *procedură* (funcție Matlab), care returnează mai mult de o valoare.

Există expresii specifice sistemului Matlab, care nu au corespondent în alte limbaje evolute. Dacă vrem să atribuim unei variabile valori care sunt obținute ca și elemente ale unei progresii aritmetice, atunci putem folosi una din instrucțiunile de atribuire

```
>>vi:pas:vf
>>x = vi:pas:vf
```

unde `vi`, `pas`, `vf` sunt expresii aritmetice, care pot fi evaluate la momentul executării instrucțiunii. În urma executării acestor instrucțiuni variabila `ans` respectiv `x` va conține valorile numerice începând cu `vi` până la `vf` cu pasul `pas`.

Dacă valoarea lui `vf` este mai mică decât cea a lui `vi` și `pas` are valoare negativă, atunci rezultatul instrucțiunii de atribuire este *matricea vidă*, având sintaxa `[]`.

Dacă `pas=1`, atunci se pot considera formele echivalente

```
>>vi:vf
>>x = vi:vf
```

Un alt tip specific de instrucțiune de atribuire a sistemului Matlab este acela când variabila primește valoarea unei submatrice. De exemplu, instrucțiunile de atribuire

```
>>x = A(:,2)
>>y = A(end,1:10)
```

au ca efect obținerea în `x` a vectorului coloană format din coloana a doua a matricei `A`, iar `y` va fi un vector linie format din primele 10 elemente ale ultimei linii a matricei `A`.

Orice instrucțiune poate fi continuată pe linia următoare prin trei puncte (`...`).

O linie poate să conțină mai multe instrucțiuni separate prin virgulă sau prin punct-virgulă (`;`).

Dacă o instrucțiune este urmată de punct-virgulă (`;`), atunci valoarea variabilei nu este afișată după executarea instrucțiunii, în caz contrar valoarea curentă a variabilei va fi afișată.

1.4 Instrucțiuni de citire și scriere

1.4.1 Instrucțiunea `input`

Sintaxa comenzii `input` este:

```
>>v = input('s')
```

Prin executarea acestei instrucțiuni se afișează pe ecran șirul de caractere dat prin `s`, așteptându-se tastarea datelor ce se doresc a fi atribuite pentru variabila `v`. Când se încheie operațiunea de tastarea datelor, se tastează `Enter`, drept consecință se efectuează operația de atribuire.

1.4.2 Instrucțiunea `ginput`

Instrucțiunea `ginput` permite introducerea de date cu ajutorul *mouse*-ului dintr-o fereastră ce are trasate axele de coordonate.

Există următoarele forme ale instrucțiunii `ginput`:

```
>>[x,y] = ginput(n)
>>[x,y] = ginput
>>[x,y,button] = ginput(n)
```

Executarea unei astfel de instrucțiuni produce introducerea coordonatelor unor puncte din figura curentă, `n` când se folosesc prima și ultima variantă, respectiv un număr nedefinit, când se folosește varianta a doua. La început apare pe ecran în fereastra figurii curente cursorul realizat prin două drepte verticale. Prin poziționări succesive, cu ajutorul *mouse*-ului, ale cursorului pe anumite puncte, urmate de *click*-uri, coordonatele acestora sunt introduse respectiv în vectorii `x` și `y`. Pentru prima și ultima formă trebuie efectuate astfel de operații succesive în număr de `n`, iar pentru a doua succesiunea de operații se încheie prin tastarea comenzii `Enter`. Forma a treia va returna și vectorul `button`, care specifică ce buton al *mouse*-lui a fost folosit la fiecare operație (1, 2 sau 3, începând de la stânga).

1.4.3 Instrucțiunea `fprintf`

Afișarea pe ecran și nu numai, folosind un format propriu, se face cu ajutorul instrucțiunii `fprintf`. Vom prezenta în continuare numai modul de afișare pe ecran folosind această instrucțiune.

Cea mai simplă sintaxă a acestei instrucțiuni este

```
>>fprintf('f',x)
```

prin care valoarea variabilei `x` este afișată pe ecran conform formatului `f`.

Formatul `f` poate fi, de exemplu: `%md`, `%md\n`, `%m.zf`, `%m.zf\n`, `%m.ze`, `%m.ze\n`, unde `md`, `m.zf` și `m.ze` exprimă respectiv faptul că valoarea lui `x` se afișează pe `m` poziții ca un întreg, ca un real cu `z` cifre zecimale, ca un real normalizat

cu z cifre zecimale. Parametrul `\n` specifică faptul că după afișare se trece la rândul următor (ține locul comenzii `Enter`).

1.5 Operatori relaționali și operatori logici

1.5.1 Operatori relaționali

- `A==B` sau `eq(A,B)`, `A~=B` sau `ne(A,B)`, `A<B` sau `lt(A,B)`, `A>B` sau `gt(A,B)`, `A<=B` sau `lte(A,B)`, `A>=B` sau `gte(A,B)` – compară element cu element matricele de aceleași dimensiuni `A` și `B` și produce o matrice ce are valorile 1 și 0, când condiția este adevărată respectiv falsă, iar dacă unul din operanzi este un scalar, acesta este extins la matricea constantă dată prin scalarul respectiv, având dimensiunile celeilalte matrice, după care se efectuează operația relațională.

1.5.2 Operatori logici

- `A&B` sau `and(A,B)` (*conjuncția*) – compară element cu element matricele de aceleași dimensiuni `A` și `B` și produce o matrice ce are valorile 1 și 0, când ambele elemente sunt diferite de zero, respectiv cel puțin unul din cele două elemente este zero;
- `A|B` sau `or(A,B)` (*disjuncția*) – compară element cu element matricele de aceleași dimensiuni `A` și `B` și produce o matrice ce are valorile 1 și 0, când cel puțin unul din cele două elemente este zero, respectiv ambele elemente sunt zero;
- `xor(A,B)` (*disjuncția exclusivă*) – compară element cu element matricele de aceleași dimensiuni `A` și `B` și produce o matrice ce are valorile 1 și 0, când unul și numai unul din cele două elemente este zero, respectiv ambele elemente sunt zero sau ambele sunt diferite de zero;

Dacă unul din operanzi este un scalar, acesta este extins la matricea constantă dată prin scalarul respectiv, având dimensiunile celeilalte matrice, după care se efectuează operația logică.

- `~A` sau `not(A)` (*negația*) – produce o matrice de aceleași dimensiuni cu cele ale matricei `A` având valorile 1 și 0, când elementul corespunzător al matricei `A` este zero, respectiv când este diferit de zero.
- `any(A)` – dacă `A` este vector, returnează valoarea 1, dacă există componente ale vectorului diferite de zero și 0 în caz contrar, iar dacă `A` este o matrice, operează

asupra fiecărei coloane a matricei și returnează corespunzător un vector de 1 și 0 cu lungimea egală cu cea a numărului coloanelor matricei A;

- `all(A)` – dacă A este vector, returnează valoarea 1, dacă toate componentele vectorului sunt diferite de zero și 0 în caz contrar, iar dacă A este o matrice, operează asupra fiecărei coloane a matricei și returnează corespunzător un vector de 1 și 0 cu lungimea egală cu cea a numărului coloanelor matricei A.

1.6 Fișiere script (m-file)

Pentru a executa o secvență de instrucțiuni Matlab (program), fără intervenția utilizatorului, se editează un fișier ASCII cu extensia `.m`, numit (*fișier script*) sau *m-file*, care să conțină instrucțiunile programului. Lansarea execuției programului se face prin tastarea numelui (fără extensia `.m`), după care se tastează `Enter`. Astfel, dacă utilizatorul și-a creat un fișier script cu numele `exemplu.m`, atunci linia de comandă prin care se lansează execuția este

```
>>exemplu
```

Fișierul `exemplu.m` ar putea avea următorul conținut

```
m = input('m (nr. natural > 0):');
A = magic(m);
d = ndims(A)
[l c] = size(A)
B = A*ones(m,1);
C = B'
```

În urma lansării execuției acestui program, fără a putea interveni până la încheierea execuției, se produc următoarele operații:

- execută instrucțiunea din prima linie, prin afișarea pe ecran a mesajului

```
m (nr. natural > 0):
```

și se așteaptă tastarea unui număr întreg pozitiv, după care prin comanda `Enter` se trece la instrucțiunea din linia următoare;

- instrucțiunea din linia a doua va genera matricea magică pătratică A de ordin m, care nu va fi afișată pe ecran deoarece instrucțiunea se termină cu simbolul punct-virgulă;
- se execută funcția `ndim`, iar rezultatul este atribuit variabilei d, care va avea valoarea 2 și va fi afișat pe ecran, deoarece instrucțiunea nu se termină cu punct-virgulă, pe două rânduri

```
d =
2
```

- următoarea instrucțiune atribuie variabilelor l și c numărul liniilor și coloanelor matricei A , care în cazul de față coincid cu valoarea lui m , iar valorile lui l și c vor fi afișate pe ecran;
- B va fi un vector coloană cu m componente și conține produsul dintre matricea A și vectorul coloană de m componente, toate având valoarea 1, iar dacă avem în vedere că A este o matrice magică, toate componentele lui B vor fi egale;
- ultima instrucțiune ajută la afișarea pe ecran a componentelor vectorului linie C , care este vectorul transpus al vectorului coloană B .

Editarea unui fișier script se poate face cu orice editor de text, dar sistemul Matlab dispune de un editor propriu.

Editorul Matlab se poate lansa din fereastra de lucru Matlab prin comenzi specifice sistemului Windows, alegând din meniul `File` opțiunea `New` sau `Open`, dacă se dorește crearea unui nou fișier, respectiv se dorește editare unui fișier deja existent. Același rezultat se obține prin comanda

```
edit file.m
```

în care `file.m` este opțional. Dacă lipsește acest argument atunci urmează crearea (editarea) unui fișier nou, iar dacă acest argument este prezent, fișierul cu acest nume urmează să fie reactualizat (editat).

Editorul sistemului Matlab poate fi lansat din orice altă parte prin activarea unui fișier cu extensia `.m`.

Editorul Matlab are capacitatea de a efectua unele verificări sintactice, de exemplu corespondența paranteză deschisă – paranteză închisă, dar este prevăzut și cu un set de comenzi (*debugger*) pentru depanarea programului prin evaluarea și modificarea unor variabile și expresii, crearea și ștergerea unor puncte de întrerupere etc.

Comenzile specifice sistemului de depanare Matlab sunt:

- `dbstop` – crearea unui punct de întrerupere (*breakpoint*);
- `dbclear` – ștergerea unui punct de întrerupere;
- `dbstatus` – listează punctele de întrerupere;
- `dbstack` – listează stiva apelurilor;
- `dbcont` – reluarea execuției programului;
- `dbstep` – execuția uneia sau mai multe instrucțiuni;
- `dbtype` – generează lista fișierelor cu extensia `.m`, împreună cu numărul liniilor acestora;

- `dbup` – schimbă spațiul de lucru curent, cu cel al unui fișier cu extensia `.m`, care a fost apelat;
- `dbdown` – efectuează operația inversă comenzii `dbup`;
- `dbmex` – execută depanarea unui fișier cu extensia `.mex` (creat cu unul din limbajele Fortran și C);
- `dbquit` – încheierea depanării.

1.6.1 Comenzi pentru gestionarea fișierelor

- `type file` – listează conținutul fișierului `file` din directorul curent;
- `delete file` – șterge fișierului `file` din directorul curent;
- `what` – listează fișierele cu extensiile `.m` din directorul curent;
- `which numef` – listează calea directorului în care se află funcția `numef`;

Comanda

```
>>what director
```

extrage lista fișierelor directorului specificat `...\toolbox\matlab\director`

Lista acestor directoare poate fi vizualizată prin comanda `help` fără nici un parametru.

1.7 Grafică bidimensională

Sistemul Matlab dispune de o largă gamă de instrucțiuni sau proceduri (funcții) pentru reprezentarea grafică a datelor, care permit o apelare simplă și eficientă.

1.7.1 Instrucțiunea `plot`

Formele acceptate de sistemul Matlab pentru instrucțiunea `plot`, privind reprezentarea grafică în plan sunt

```
plot(y)
plot(x,y)
plot(x1,y1,x2,y2,...)
plot(y,s)
plot(x,y,s)
plot(y1,s1,y2,s2,...)
plot(x1,y1,s1,x2,y2,s2,...)
```

În cazul primei forme, dacă `y` este un vector de lungime `m`, se reprezintă grafic linia poligonală ce trece prin punctele de coordonate (i, y_i) , $i = \overline{1, m}$.

Programul 1.7.1. Să considerăm programul, care reprezintă grafic linia poligonală ce unește punctele de coordonate (i, y_i) , $i = \overline{1, m}$, unde y reprezintă elementele diagonalei unei matrice magice de ordin m :

```
m = input('m:');
y = diag(magic(m));
plot(y)
```

În urma execuției programului, pentru $m=10$, se obține graficul din Figura 1.1.

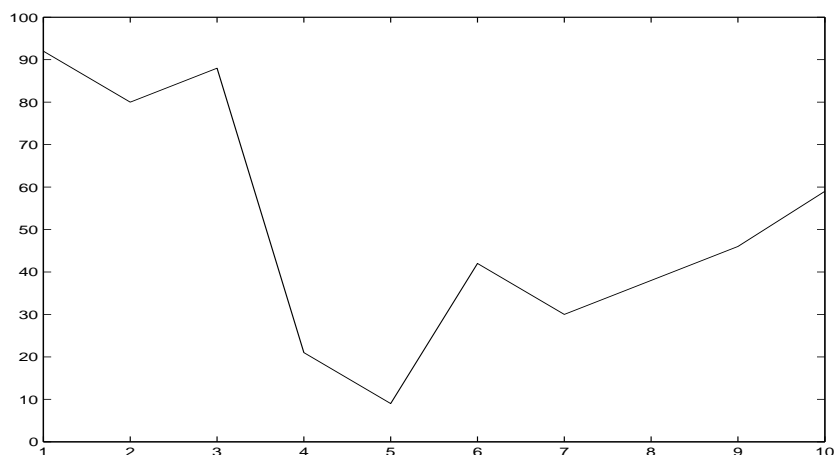


Figura 1.1: Linie poligonală

Dacă y este o matrice de tipul (m, n) , atunci se reprezintă grafic n linii poligonale, corespunzând celor n coloane ale matricei y , anume $(i, y_{ij})_{i=1}^m$, $j = \overline{1, n}$.

Programul 1.7.2. Să considerăm programul, care reprezintă grafic liniile poligonale ce unesc respectiv punctele de coordonate $(i, y_{ij})_{i=1}^m$, pentru fiecare $j = \overline{1, n}$, unde y este o matrice magică de ordin m , iar n este un număr întreg pozitiv cel mult m :

```
m = input('m:');
n = input('n (n<=m):');
y = magic(m);
plot(y(:,1:n))
```

În urma execuției programului, pentru $m=10$ și $n=2$, se obține graficul din Figura 1.2.

Dacă y este de tip complex, atunci prima formă este echivalentă cu a doua formă a instrucțiunii `plot`, adică cu

```
plot(real(y), imag(y))
```

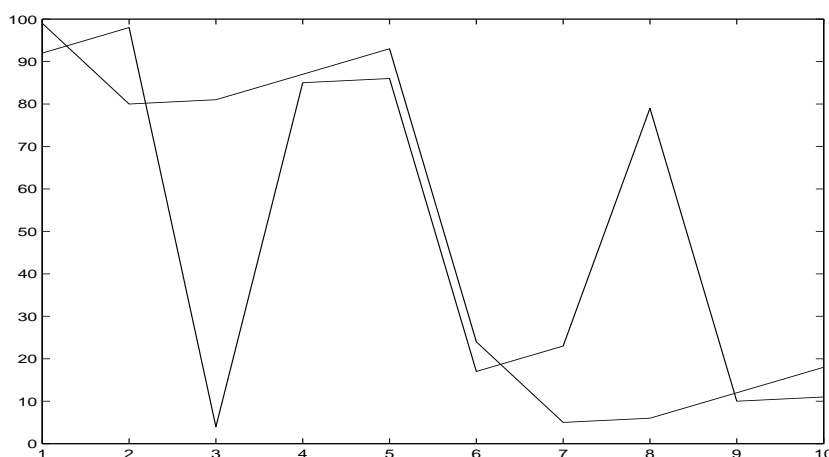


Figura 1.2: Linii poligonale

Programul 1.7.3. Să considerăm programul, care reprezintă grafic funcția cu valori complexe $z = x \cos(2\pi x) + ix^2 \sin(2\pi x)$, $x \in [0, 1]$.

```
m = input('m:');
h = 1/m; x = 0:h:1;
z = x.*cos(2*pi*x)+i*x.^2.*sin(2*pi*x);
plot(z)
```

În urma execuției programului, pentru $m=100$, se obține graficul din Figura 1.3.

Remarcăm faptul că pentru toate celelalte forme, dacă avem parametri complecși, partea imaginară este ignorată.

Dacă se utilizează a doua formă, iar x și y sunt vectori de aceeași lungime m , atunci se reprezintă grafic linia poligonală ce trece prin punctele de coordonate (x_i, y_i) , $i = \overline{1, m}$.

Programul 1.7.4. Fie programul care reprezintă grafic funcția $\sin(2\pi x)$ pe intervalul $[0, 1]$:

```
m = input('m:');
h = 1/m; x = 0:h:1;
plot(x, sin(2*pi*x))
```

Execuția programului, pentru $m=200$, generează graficul din Figura 1.4.

În cazul în care x și y sunt matrice având același tip (m, n) , se reprezintă grafic n linii poligonale, corespunzând celor n coloane ale matricelor, adică $(x_{ij}, y_{ij})_{i=1}^m$, pentru fiecare $j = \overline{1, n}$. Dacă x este un vector de lungime m , iar y este un vector de lungime m sau matrice de tipul (m, n) , se reprezintă grafic liniile poligonale ce trec

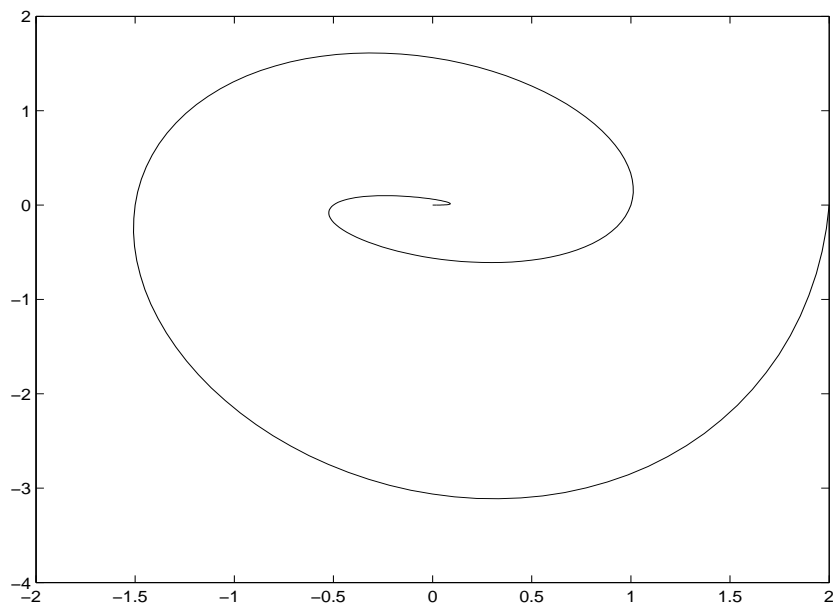


Figura 1.3: $z = x \cos(2\pi x) + ix^2 \sin(2\pi x)$

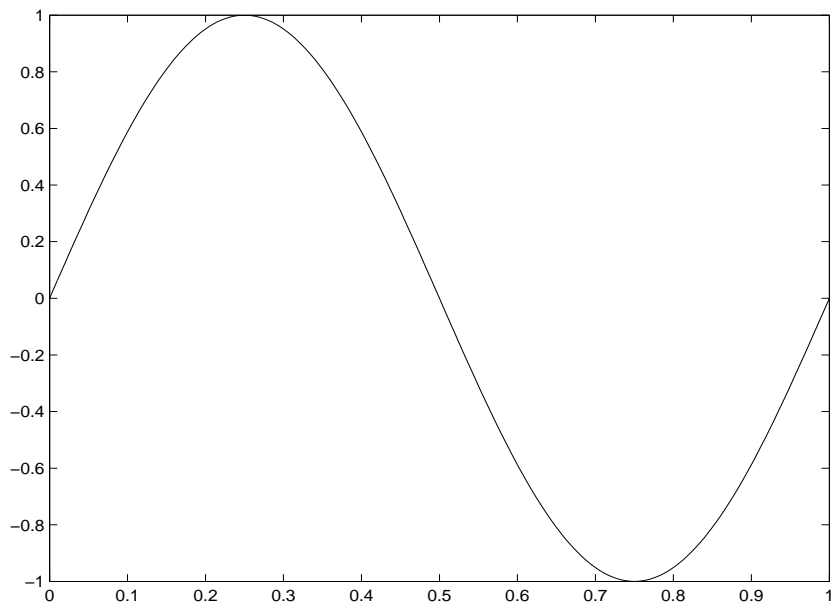


Figura 1.4: $y = \sin(2\pi x)$

prin punctele de coordonate $(x_i, y_i)_{i=1}^m$, respectiv n linii poligonale ce trec respectiv prin punctele date prin $(x_i, y_{ij})_{i=1}^m$, pentru fiecare $j = \overline{1, n}$.

Forma a treia a instrucțiunii `plot` implică reprezentarea grafică pe aceeași figură, conform formei precedente, a grupelor de date $(x_1, y_1), (x_2, y_2)$ ș.a.m.d.

Programul 1.7.5. Vom prezenta două variante de programe, care reprezintă pe aceeași figură graficele funcțiilor $\sin(2\pi x)$ și $\cos(2\pi x)$ pe intervalul $[0, 1]$:

```
m = input('m:');
h = 1/m; x = 0:h:1;
y = 2*pi*x;
plot(x, sin(y), x, cos(y))
```

respectiv

```
m = input('m:');
h = 1/m; x = 0:h:1;
y = 2*pi*x;
y = [sin(y)' cos(y)']; x=[x' x'];
plot(x, y)
```

Executând una din cele două variante, cu $m=200$, se generează graficul din Figura 1.5.

Remarcăm faptul că prima variantă ar fi potrivită dacă cele două funcții se reprezintă grafic pe același domeniu, iar a doua variantă când cele două funcții sunt reprezentate grafic pe domenii diferite.

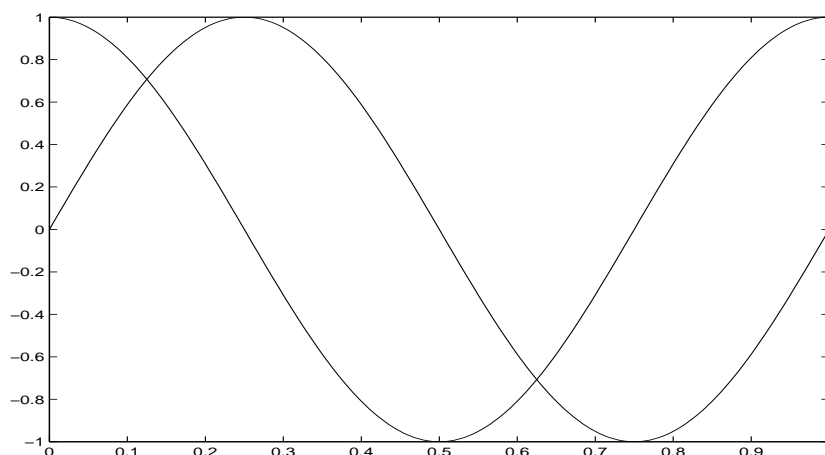


Figura 1.5: $y_1 = \sin(2\pi x)$, $y_2 = \cos(2\pi x)$

Pentru o reprezentare netedă este de dorit ca numărul punctelor ce generează liniile poligonale să fie suficient de mare.

Pe de altă parte, observăm că prin primele trei tipuri de instrucțiune `plot` nu se specifică nici o caracteristică a liniilor reprezentate grafic. Sistemul Matlab face o alegere implicită a acestor caracteristici.

Ultimele patru forme ale instrucțiunii `plot` permit utilizatorului specificarea unor caracteristici ale liniilor reprezentate grafic cu ajutorul parametrului `s`.

Parametrul de tipul `s` este un șir de cel mult trei caractere, prin care se precizează culoarea, tipul de marcaj al punctelor și tipul de linie.

Culoarea este specificată cu unul din simbolurile următoare: `y` – galben (yellow), `m` – violet (magenta), `c` – ciclamen (cyan), `r` – roșu (red), `g` – verde (green), `b` – albastru (blue), `w` – alb (white), `k` – negru (black).

Tipul de marcaj este dat prin unul din simbolurile: `.` – punct (point), `o` – cerculeț (circle), `x` – cruciuliță diagonală (x-mark), `+` – cruciuliță (plus), `*` – steluță (star), `s` – pătrățel (square), `d` – romb (diamond), `v` – triunghi cu vârf în jos (triangle down), `^` – triunghi cu vârf în sus (triangle up), `<` – triunghi cu vârf spre stânga (triangle left), `>` – triunghi cu vârf spre dreapta (triangle right), `p` – steluță în cinci colțuri (pentagram), `h` – steluță în șase colțuri (hexagram).

Următoarele patru tipuri de linie sunt disponibile: continuă (solid) (`-`), punctată (dotted) (`:`), întreruptă (dashed) (`--`), linie-punct (dashdot) (`-.`).

De exemplu, dacă `s='c:'`, graficul este trasat cu linie punctată ciclamen, iar dacă `s='bd'` punctele sunt marcate prin semnul romb de culoare albastră, dar nu sunt unite între ele.

Programul 1.7.6. Vom reprezenta grafic funcțiile $\sin(2\pi x)$ și $\cos(2\pi x)$ pe intervalul $[0, 1]$, prima funcție prin linie continuă, iar a doua prin linie întreruptă, iar punctele de pe cele două curbe, care au abscisele multipli de $\frac{1}{6}$, să fie marcate prin cerculețe, respectiv prin steluțe.

```
m = input('m:');
h = 1/m; x = 0:h:1; y = 2*pi*x;
Y = [sin(y)' cos(y)'];
h = 1/m; p = 0:1/6:1; pp = 2*pi*p;
plot(x,y,p,sin(pp),'o',p,cos(pp),'*')
```

Pentru `m=200`, graficul este cel din Figura 1.6.

1.7.2 Comenzi și instrucțiuni de gestionare a graficelor

Sistemul Matlab dispune de comenzi și instrucțiuni prin care sunt gestionate și specificate anumite caracteristici ale graficelor produse privind titlul, etichetarea axelor, inserarea unor texte, marcarea unor rețele, delimitarea graficului etc:

- `hold on` – păstrează graficul curent, un nou grafic va fi suprainprimat;
- `hold off` – graficul curent este abandonat, fără a fi șters;

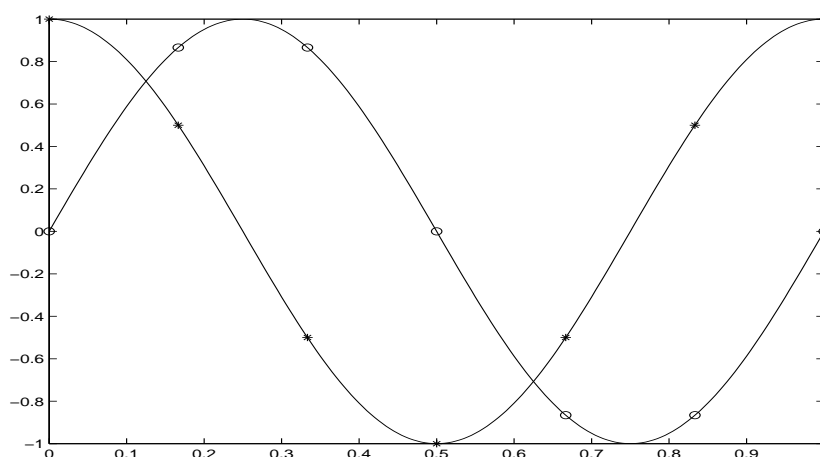


Figura 1.6: $y_1 = \sin(2\pi x)$ și $y_2 = \cos(2\pi x)$

- `hold` – trece la modul `off`, dacă sistemul se afla în modul `on`, respectiv la modul `on`, când sistemul se află în modul `off`;
- `clf` – șterge graficul curent;
- `cla` – șterge graficul curent, păstrând axele;
- `figure(n)` – activează fereastra cu graficul `n` anterior obținut;
- `axis([xmin xmax ymin ymax])` – fixează limitele pentru axele absciselor și ordonatelor pentru figura curentă;
- `v = axis` – întoarce un vector linie ce conține cele patru limite pentru axele de coordonate, care pot fi și `-Inf`, `Inf`;
- `axis auto` – consideră forma implicită a limitelor axelor de coordonate;
- `axis manual` – păstrează limitele curente, astfel dacă este utilizată instrucțiunea `hold`, figurile următoare vor avea aceleași limite pentru axele de coordonate;
- `axis tight` – atribuie pentru limitele axelor cele mai apropiate valori ce rezultă din datele ce se reprezintă grafic;
- `axis ij` – configurează axele în modul `matrix`, adică originea axelor de coordonate este colțul din stânga sus, axa `i` fiind verticală și marcată de sus în jos, iar axa `j` este orizontală și este marcată de la stânga la dreapta;

- `axis xy` – configurează axele în modul cartezian, mod implicit;
- `axis equal` – unitățile de măsură pe cele două axe au aceeași mărime;
- `axis image` – la fel cu `axis equal`, cu excepția că figura este restrânsă la domeniul datelor;
- `axis square` – pune unitățile de măsură pe cele două axe de coordoante astfel încât figura să fie prezentată într-un pătrat;
- `axis normal` – rearanjează axele în forma implicită;
- `axis off` – elimină axele de coordoante;
- `axis on` – inserează axele de coordoante;
- `title('text')` – afișează textul specificat în partea de sus a graficului;
- `legend('str1','str2',...)` – afișează legenda prin care se efectuează corespondența dintre fiecare tip de curbă a graficului și textele specificate prin `str1`, `str2` etc.;
- `legend('str1','str2',...,p)` – afișează legenda de tipul precizat în pozițiile respectiv înafara axelor ($p=-1$), în interiorul axelor, dar pe cât e posibil să nu acopere figurile trasate ($p=0$), în partea dreaptă sus ($p=1$), care este poziția implicită, în partea stângă sus ($p=2$), în partea stângă jos ($p=3$), în partea dreaptă jos ($p=4$);
- `legend off` – elimină legenda;
- `xlabel('text')` – etichetează axa absciselor cu textul precizat;
- `ylabel('text')` – etichetează axa ordonatelor cu textul precizat;
- `text(x,y,'text')` – inserează textul precizat prin `text`, în punctul de coordonate (x,y) , iar dacă x și y sunt vectori (de aceeași lungime), atunci inserează textul respectiv în toate punctele de coordoante (x_i, y_i) , $i = 1, 2, \dots$, pe când dacă textul este dat printr-o matrice cu același număr de linii cu cel al lungimilor vectorilor x și y , atunci textul din linia i este inserat în punctul de coordonate (x_i, y_i) , $i = 1, 2, \dots$.
- `gtext('text')` – afișează graficul, împreună cu cursorul plus, unde urmează să fie inserat textul, cursorul putând fi poziționat cu ajutorul *mouse*-ului sau tastele cu săgeți, după care prin apăsarea oricărei taste (*mouse* sau *keyboard*), se obține inserarea textului;

- `grid` – adaugă pe grafic o rețea rectangulară;
- `grid off` – grafic fără rețea rectangulară;
- `zoom` – permite mărirea unei părți a figurii pentru urmărirea anumitor detalii.

Se fixează cursorul *mouse*-ului pe poziția în care suntem interesați a fi mărită. Prin *dublu-click* pe butonul din stânga se produce o mărire a zonei de două ori, iar prin *dublu-click* pe butonul din dreapta se produce o micșorare a zonei de două ori. Dacă se păstrează apăsat butonul din dreapta și se deplasează cursorul, se obține un dreptunghi, care prin eliberarea butonului *mouse*-ului, va umple întreaga figură.

Prin instrucțiunile

```
zoom off
zoom out
```

se dezactivează instrucțiunea `zoom`, respectiv se revine la dimensiunile inițiale ale figurii.

Comenzile de tip `axis` se impun a fi inserate după instrucțiunea `plot`.

Menționăm că textele introduse pe grafic admit sintaxa sistemului \LaTeX simplificat.

Programul 1.7.7. Vom considera un exemplu de reprezentare grafică prin *mascarea* unor părți a graficului. Să reprezentăm grafic cu linie continuă pe intervalul $[0, 6]$ funcția care coincide cu $\sin(\pi x)$, când valorile acestei funcții sunt pozitive, iar în rest să fie zero. Pe aceeași figură, să reprezentăm prin linie întrerupă și $\sin(\pi x)$.

```
m = input('m:');
h = 1/m; x = 0:h:6;
y = sin(pi*x); Y=ge(y,0).*y;
plot(x,y,'k--',x,Y,'k-')
```

Pentru $m=200$, graficul este cel din Figura 1.7.

1.7.3 Instrucțiunile `polar` și `ezpolar`

Instrucțiunea `polar` este analogă cu instrucțiunea `plot`, numai că este considerată reprezentarea curbei în coordoante polare, iar instrucțiunea `ezpolar` are sintaxa

```
ezpolar('f',[a,b])
```

și are ca efect reprezentarea grafică a funcției în coordonate polare, dată prin expresia algebrică f , pe intervalul $[a, b]$. Al doilea parametru este opțional. Valoarea implicită este $[0, 2\pi]$.

De exemplu, comanda

```
ezpolar('1+cos(theta)')
```

va produce graficul din Figura 1.8. Mai remarcăm faptul că parametrul f poate fi numele unei funcții Matlab predefinite sau a utilizatorului.

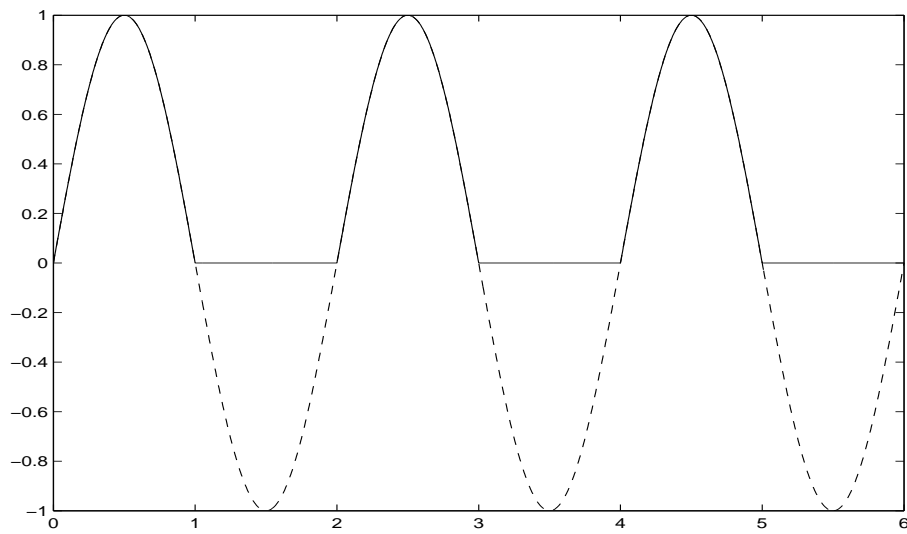


Figura 1.7: $y = \sin(\pi x)$

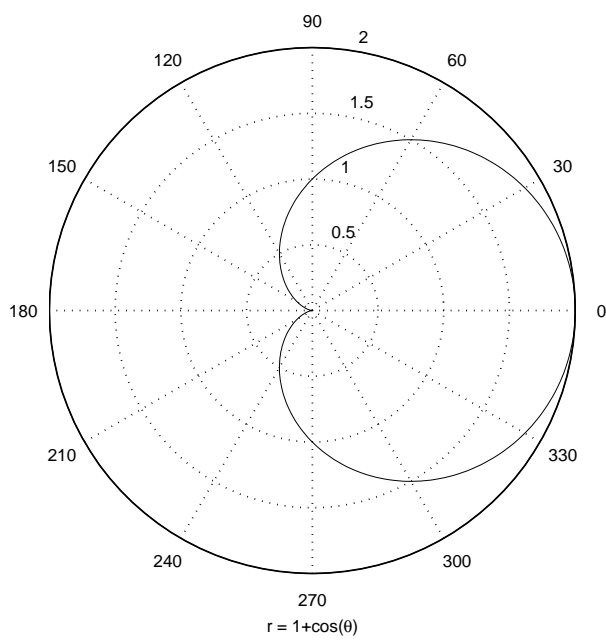


Figura 1.8: $r = 1 + \cos(\theta)$, $\theta \in [0, 2\pi]$

Programul 1.7.8. Următorul program

```

clf; t=0:0.01:pi;
r=sin(3*t).*exp(-0.3*t);
polar(t,r,'k-'), grid

```

are ca efect graficul din Figura 1.9.

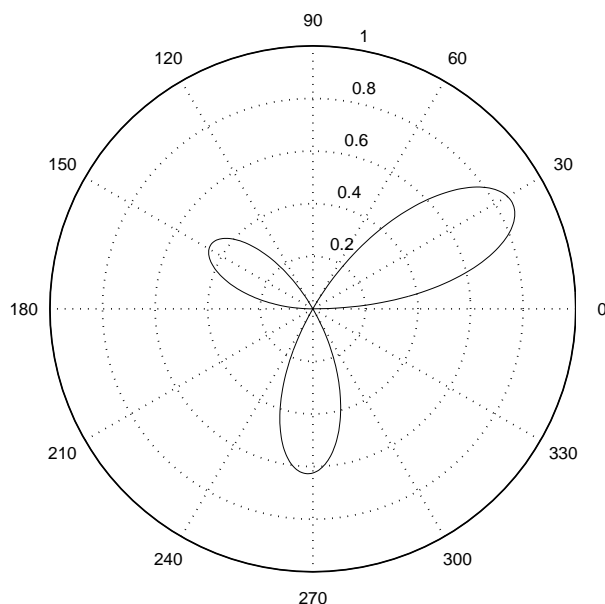


Figura 1.9: $r = \sin(3t) e^{-0.3t}$

Instrucțiunea colormap

Sistemul Matlab dispune de un set de gestiune a culorilor în reprezentările grafice. Instrucțiunea `colormap` permite definirea unui set propriu pentru gestiunea culorilor. Lansarea unei instrucțiuni se face prin comanda

```

colormap(map)
colormap('default')

```

A doua comandă setează setul de culori la forma standard (implicită).

Parametrul `map` este o matrice cu trei coloane, fiecare linie conținând elemente numere din intervalul $[0, 1]$, prin care se definește câte o culoare. Modul de definire a culorii, după această regulă, se numește RGB (Red-Green-Blue), și specifică în cadrul culorii, respectiv intensitățile culorilor roșu, verde și albastru.

Sistemul Matlab dispune de asemenea de unele seturi de culori, care pot fi selectate și specificate de utilizator. O parte din acestea sunt obținute prin:

```
colormap spring
```

generează nuanțe de culori între violet (magenta) și galben;

```
colormap summer
```

generează nuanțe de culori între verde și galben;

```
colormap autumn
```

generează nuanțe de culori între roșu, trecând prin portocaliu spre galben;

```
colormap winter
```

generează nuanțe de culori între albastru și verde;

```
colormap gray
```

generează nuanțe de culori gri.

1.7.4 Instrucțiunea `stairs`

Sistemul Matlab dispune de procedura `stairs` pentru reprezentarea grafică a funcțiilor în scară și care acceptă formele instrucțiunii `plot`.

Instrucțiunea de forma

```
stairs(y)
```

reprezintă grafic funcția în scară cu valorile precizate prin vectorul y de lungime m , iar abscisele sunt considerate numerele de la 1 la m . Dacă y este o matrice de tipul (m, n) vor fi reprezentate grafic n funcții în scară, câte una pentru fiecare coloană a matricei y .

Dacă se consideră instrucțiunea de forma

```
stairs(x,y)
```

se va reprezenta grafic funcția în scară cu valorile precizate prin vectorul y și cu abscisele corespunzătoare date prin vectorul x (componentele lui x trebuie să fie crescătoare).

Dacă y și x sunt matrice de același tip (m, n) , atunci se vor reprezenta grafic n funcții în scară, câte una pentru fiecare din coloanele celor două matrice (elementele coloanelor matricei x trebuie să fie crescătoare). Dacă matricea x este un vector coloană cu m componente crescătoare, atunci funcțiile în scară corespunzătoare celor n coloane ale matricei y se consideră că au aceleași abscise precizate prin vectorul x .

Instrucțiunea de tipul

```
stairs(x,y,s)
```

permite prezentarea tipului liniei cu care se trasează graficul cu ajutorul parametrului s și care are sintaxa cea precizată la instrucțiunea `plot`.

Prin urmare, putem spune că toate variantele folosite în cadrul instrucțiunii `plot` au corespondente pentru instrucțiunea `stairs`.

Programul 1.7.9. Fie funcția lui Haar, numită în analiza wavelet *funcție wavelet mamă*

$$\psi(x) = \begin{cases} 1, & \text{dacă } 0 \leq x < \frac{1}{2}, \\ -1, & \text{dacă } \frac{1}{2} \leq x < 1, \\ 0, & \text{în rest.} \end{cases}$$

Folosind funcția ψ se definește familia de funcții wavelet

$$\psi_{j,k}(x) = 2^{\frac{j}{2}} \psi(2^j x - k)$$

unde j și k se numesc *indice de dilatare*, respectiv *indice de translatare*. Dacă se are în vedere formula de definiție pentru funcția ψ , putem scrie

$$\psi_{j,k}(x) = \begin{cases} 2^{\frac{j}{2}}, & \text{dacă } \frac{1}{2^j}k \leq x < \frac{1}{2^j}k + \frac{1}{2^{j+1}}, \\ -2^{\frac{j}{2}}, & \text{dacă } \frac{1}{2^j}k + \frac{1}{2^{j+1}} \leq x < \frac{1}{2^j}k + \frac{1}{2^j}, \\ 0, & \text{în rest.} \end{cases}$$

Vom scrie un program care să reprezinte grafic funcția f definită segmentar cu ajutorul funcțiilor $\psi_{j,k}$, $k = \overline{r, s}$, unde j , r și s sunt fixați ($r < s$).

Funcția f este o funcție în scară, având punctele de discontinuitate $\frac{k}{2^{j+1}}$, pentru $k = \overline{2r, 2s+2}$, cu valorile corespunzătoare acestor puncte, $2^{j/2}$, pentru k par, respectiv $-2^{j/2}$, pentru k impar.

Cu aceste precizări, avem următorul program Matlab pentru reprezentarea grafică a funcției f :

```
j = input('j=');
r = input('r=');
s = input('s (>r)=');
rs = 2*s-2*r+3;
fprintf('rs=2s-2r+3=%2d\n',rs);
y = input('rs valori alternative [1, -1 ... 1]:');
y = 2^(j/2)*y;
x = 1/2^j*[r:1/2:s+1];
stairs(x,y)
```

Execuția programului, pentru $j=1$, $s=1$, $r=-2$, produce graficul din Figura 1.10.

Singurul inconvenient la execuția acestui program este legat de instrucțiunea `input` prin care se inițializează vectorul y . Când numărul funcțiilor $\psi_{j,k}$ prin care se definește funcția f este mare, atunci rs este mare și introducerea listei lui y este anevoioasă.

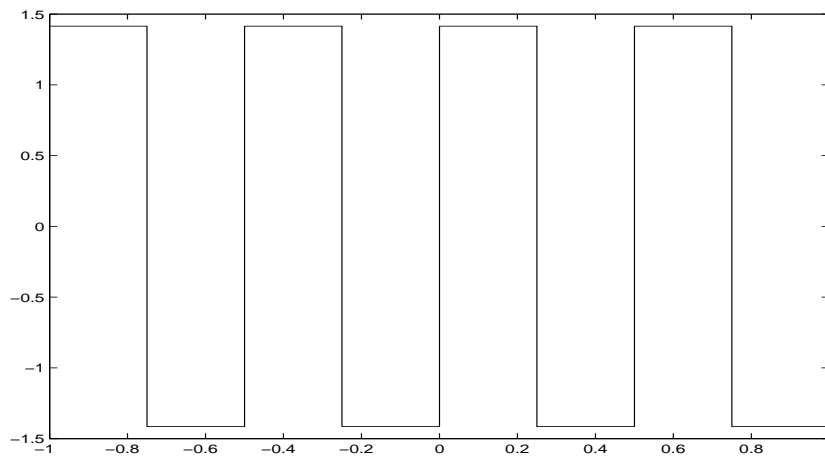


Figura 1.10: $f = (\psi_{1,-2}, \psi_{1,-1}, \psi_{1,0}, \psi_{1,1}, \psi_{1,2})$

O variantă a programului precedent, pentru care nu este necesară inițializarea lui y de la tastatură, ar putea fi următorul program:

```
j = input('j=');
r = input('r=');
s = input('s (>r)=');
rs = 2*s-2*r+3;
y = (-ones(1,rs)).^[2*r:2*s+2];
y = 2^(j/2)*y;
x = 1/2^j*[r:1/2:s+1];
stairs(x,y)
```

Programul 1.7.10. Să reluăm reprezentarea grafică a funcțiilor wavelet Haar și să scriem un program care să reprezinte pe aceeași figură două secvențe de funcții wavelet, pentru două valori consecutive ale indicelui de dilatare j , iar secvențele alese să acopere intervalul $[-1, 1]$. Pentru aceasta trebuie ca $r = -2^j$, iar $s = 2^j - 1$. Totodată să trasăm graficul primei secvențe cu linie continuă și punctele să le marcăm prin cerculețe, iar pentru a doua secvență să alegem linie punctată pentru reprezentarea grafică, iar punctele să fie marcate cu stelute în cinci colțuri:

```
clf, hold on, axis off
plot([-1.2 1.2],[0 0], 'k:')
j = input('j=');
r = -2^\{ }j; s=-r-1;
rs = 2*s-2*r+3;
y = (-ones(1,rs)).^[2*r:2*s+2];
y = 2^(j/2)*y;
x = 1/2^j*[r:1/2:s+1];
stairs(x,y, 'ko-')
```



```

j = j+1; r = -2^j; s=-r-1;
rs = 2*s-2*r+3;
y = (-ones(1,rs)).^[2*r:2*s+2];
y = 2^(j/2)*y;
x = 1/2^j*[r:1/2:s+1];
stairs(x,y,'kp--')

```

Execuția programului, pentru $j=0$, produce graficele din Figura 1.11. Se observă

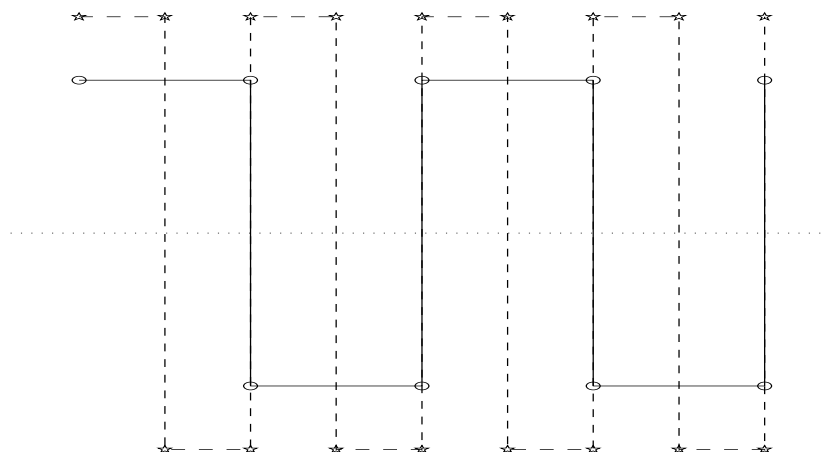


Figura 1.11: $(\psi_{0,-1}, \psi_{0,0}), (\psi_{1,-2}, \psi_{1,-1}, \psi_{1,0}, \psi_{1,1})$

că programul mai elimină axele sistemului Matlab, dar adaugă axa absciselor trasată punctat.

1.7.5 Instrucțiunile `bar` și `barh`

Aceste două instrucțiuni, după cum spun și denumirile, reprezintă anumite date numerice cu ajutorul barelor (batoanelor) verticale și respectiv orizontale.

Forme ale acestor instrucțiuni sunt:

```

bar(y,w,'tip','color')
bar(x,y,w,'tip','color')
barh(y,w,'tip','color')
barh(x,y,w,'tip','color')

```

argumentele `w`, `'tip'` și `'color'` sunt opționale, dar când sunt specificate, trebuie să păstreze această ordine.

Parametrul `x` este un vector de lungime `m`.

Dacă `y` este un vector de aceeași lungime ca și `x`, atunci vor fi reprezentate grafic în fiecare din punctele de abscise x_i , $i = \overline{1, m}$, câte o bară de înălțimile date prin y_i ,

$i = \overline{1, m}$, care pot să fie și negative. Dacă parametrul x lipsește, atunci se consideră implicit că x are componentele ca fiind primele m numere naturale.

Dacă y este o matrice de tipul (m, n) , reprezentările sunt efectuate pentru fiecare din cele n coloane ale matricei, adică în fiecare punct x_i , $i = \overline{1, m}$, de pe axa absciselor vor fi reprezentate câte n bare (batoane).

Parametrul w specifică lățimea barelor, implicit fiind $w = 0.8$, iar pentru $w > 1$ barele se suprapun.

Dacă parametrul `tip=grouped`, care este valoarea implicită, atunci avem reprezentările precizate mai sus, iar dacă `tip=stacked` barele (batoanele) vor fi stivuite pe verticală.

Parametrul `color` specifică culoarea barelor (batoanelor) și are una din următoarele valori: `r` (roșu), `g` (verde), `b` (albastru), `y` (galben), `m` (violet), `c` (ciclamen), `k` (negru), `w` (alb).

Tot ce s-a spus despre `bar` se poate spune și pentru `barh`, numai că orientarea barelor este pe orizontală.

1.7.6 Instrucțiunea `subplot`

Sistemul Matlab are posibilitatea de a împărți o fereastră, pentru reprezentare grafică, în $m \times n$ zone (ferestre), cu ajutorul instrucțiunii `subplot`.

Forma generală a instrucțiunii este

```
subplot(m,n,p)
```

unde m , n și p sunt numere naturale, care exprimă faptul că fereastra inițială se împarte în $m \times n$ ferestre aranjate pe m linii și n coloane, iar reprezentarea grafică, la momentul execuției instrucțiunii, se va face în fereastra p , numerotarea ferestrelor făcându-se pe linii.

Programul 1.7.11. Reprezentarea grafică poate fi făcută și cu ajutorul reprezentărilor parametrice ale curbelor.

Pentru reprezentarea grafică a unui cerc cu centrul în originea axelor de coordoane am putea folosi secvența de instrucțiuni:

```
r = input('raza=');
t = 0:0.01:2*pi;
x = r*sin(t); y=r*cos(t);
subplot(2,1,1), plot(x,y)
axis normal
subplot(2,1,2), plot(x,y)
axis square
```

Pentru $r=2$, graficul din partea stângă a Figurii 1.12, se obține când nu s-a impus ca unitățile de măsură pe cele două axe de coordoanate să fie aceleași, iar graficul din partea dreaptă s-a obținut după ce s-a înlăturat acest neajuns cu instrucțiunea `axis square`.

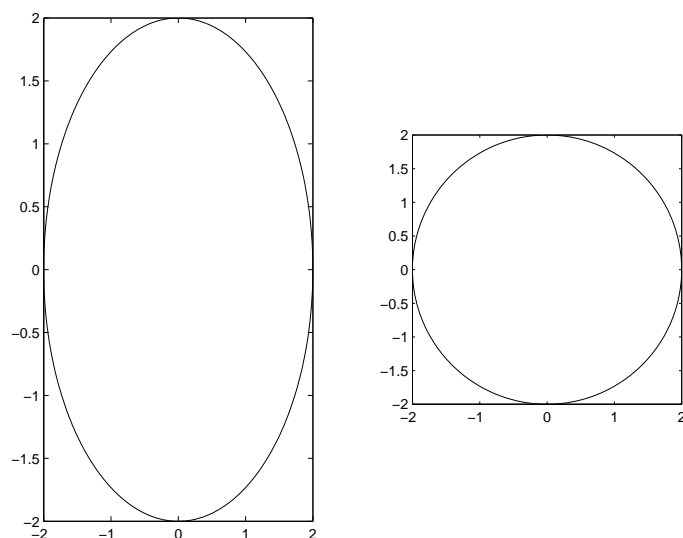


Figura 1.12: $x = r \cos t, y = r \sin t, r = 2, t \in [0, 2\pi]$

Programul 1.7.12. Pentru exemplificarea instrucțiunilor `bar`, `barh` și `subplot`, facem din nou apel la o matrice magică de ordin m și să reprezentăm în patru ferestre de tipul 2×2 , prin bare verticale grupate și prin bare verticale stivuite, valorile primelor două coloane ale matricei, respectiv prin bare orizontale grupate și prin bare orizontale stivuite valorile ultimelor două coloane:

```
m=input('m:'); A=magic(m);
subplot(221), bar(A(:,[1 2]),'grouped')
title('Primele doua coloane - grupate')
subplot(222), bar(A(:,[1 2]),'stacked')
title('Primele doua coloane - stivuite')
subplot(223), barh(A(:,[m-1 m]),'grouped')
title('Ultimele doua coloane - grupate')
subplot(224), barh(A(:,[m-1 m]),'stacked')
title('Ultimele doua coloane - stivuite')
colormap summer
```

Reprezentarea grafică, pentru $m=4$, este dată în Figura 1.13.

1.7.7 Instrucțiunea `fplot`

Pentru reprezentarea grafică a unei funcții definite cu ajutorul unei proceduri Matlab, se pot folosi una din următoarele variante ale instrucțiunii `fplot`:

```
fplot('numef',lim)
fplot('numef',lim,s)
fplot('numef',lim,tol)
```

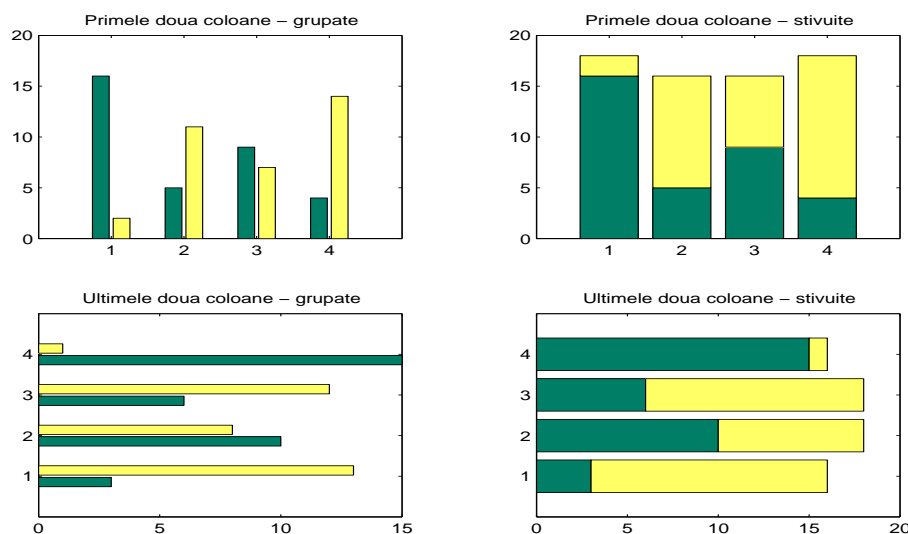


Figura 1.13: Bare verticale și bare orizontale

```
fplot('numef',lim,tol,s)
fplot('numef',lim,n)
[X,Y] = fplot('numef',lim,...)
```

În toate cazurile se reprezintă grafic funcția `numef`, cu limitele pentru axa absciselor specificate prin `lim=[xmin xmax]` sau `lim=[xmin xmax ymin ymax]`. Tipul liniei folosit în reprezentarea grafică este specificat prin parametrul `s`, ca și în cazul instrucțiunii `plot`, iar eroarea relativă folosită în reprezentarea grafică este precizată prin parametrul `tol`, care are valoarea implicită `tol=2e-3`. Parametrul `n` specifică faptul că graficul se trasează cu ajutorul a `n+1` puncte, valoarea implicită a lui `n` este `n=1`, iar valoarea maximă a pasului nu depășește $\frac{x_{\max}-x_{\min}}{n}$.

Remarcăm faptul că parametrii opționali `s`, `tol` și `n` pot fi luați în orice ordine.

Funcția `numef` poate să întoarcă, pentru un vector `x`, un vector `y` de aceeași lungime ca și `x` sau o matrice cu mai multe coloane și număr de linii ca și lungimea lui `x`. În primul caz, se va reprezenta grafic o singură curbă, iar în al doilea caz, pentru fiecare coloană a matricei `y` câte o curbă. Parametrul `numef` poate fi înlocuit și cu un șir de caractere de tipul `'[f(x),g(x),...]'`, unde `f,g,...` sunt nume de funcții Matlab, caz în care, pe aceeași figură, se vor reprezenta grafic funcțiile respective.

De exemplu, instrucțiunea

```
fplot('sin(2*pi*x), cos(2*pi*x)',[0 1])
```

va reprezenta grafic, pe aceeași figură, funcțiile $\sin(2\pi x)$ și $\cos(2\pi x)$, pe intervalul $[0, 1]$.

Ultima formă a instrucțiunii `fplot` nu produce nici un grafic, ci păstrează punctele graficului în X și Y , care pot fi apoi folosite efectiv, prin instrucțiunea `plot(X,Y)`, la reprezentarea grafică.

1.7.8 Instrucțiunea `ezplot`

Sintaxa instrucțiunii `ezplot` poate fi:

```
ezplot('f')
ezplot('f',[a,b])
ezplot('f',[a,b,c,d])
ezplot('x','y')
ezplot('x','y',[a,b])
```

Parametrul f reprezintă o expresie algebrică de una sau două variabile.

Dacă f este o expresie algebrică de o singură variabilă, atunci este reprezentată grafic funcția definită de această expresie, pe intervalul $[a, b]$. Valoarea implicită a acestui parametru este $[-2\pi, 2\pi]$.

Dacă f este o expresie algebrică de două variabile, atunci este reprezentată grafic funcția implicită definită prin $f=0$, pe domeniul $[-2\pi, 2\pi] \times [-2\pi, 2\pi]$, dacă se folosește prima formă de apel, pe domeniul $[a, b] \times [a, b]$, pentru a doua formă, respectiv pe $[a, b] \times [c, d]$, pentru a treia formă. Cele două variabile sunt considerate ca fiind ordonate lexicografic.

Ultimele două forme sunt folosite pentru reprezentările parametrice, x reprezentând abscisa, iar y ordonata. Argumentul acestor parametri variază în intervalul $[a, b]$, iar dacă acesta lipsește, atunci se consideră intervalul implicit $[0, 2\pi]$.

De exemplu, instrucțiunea

```
ezplot('sin(3*t)*cos(t)','sin(3*t)*sin(t)',[0,pi])
```

produce graficul din Figura 1.14. Mai remarcăm faptul că parametrii f , x și y , pot fi numele unor funcții Matlab predefinite sau ale utilizatorului.

1.7.9 Instrucțiunea `fill`

Instrucțiunea `fill` este folosită pentru umbrirea interiorului unui poligon specificat prin vârfurile sale:

```
fill(x,y,c)
```

unde x și y specifică vârfurile poligonului, iar c reprezintă culoarea dorită pentru interiorul poligonului.

Instrucțiunea poate fi utilizată pentru umbrirea unei părți dintr-o figură.

Programul 1.7.13. Să scriem un program care să umbrească aria cuprinsă între axa absciselor, curba lui Gauss și dreptele perpendiculare pe axa absciselor în punctele a și b ($-3 < a < b < 3$).

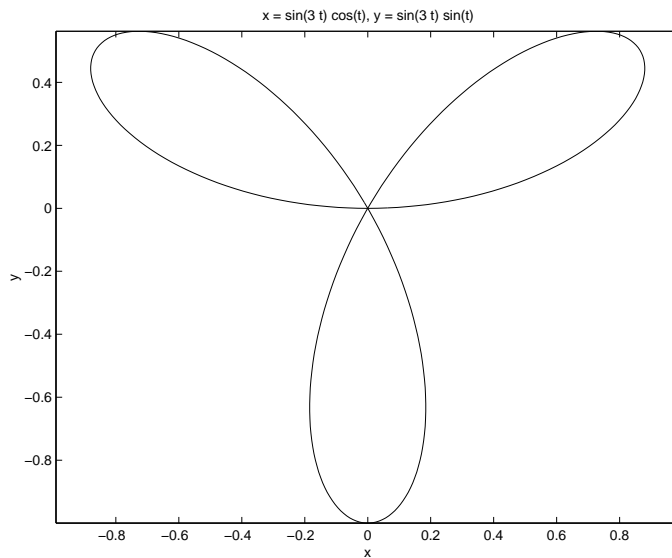


Figura 1.14: $x = \sin(3t) \cos(t)$, $y = \sin(3t) \sin(t)$, $t \in [0, \pi]$

Curba lui Gauss este dată prin funcția

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}, \quad x \in \mathbb{R}.$$

```
a = input('a='); b = input('b=');
xmin = -3; xmax = 3;
fplot('1/sqrt(2*pi)*exp(-x^2/2)', [xmin xmax])
x = a; y = 0;
t = a:0.01:b;
x = [x,t]; y = [y,1/sqrt(2*pi)*exp(-t.^2/2)];
x = [x,b]; y = [y,0];
fill(x,y,'c')
```

Pentru $a=-2$ și $b=1.5$, se obține graficul din Figura 1.15.

1.8 Instrucțiuni de ciclare și control

Ca și în celelalte limbaje de programare evolute, sistemul Matlab dispune de instrucțiuni, care permit repetarea unei secvențe de instrucțiuni de un număr definit sau indefinit de ori, precum și de instrucțiuni, care permit ieșirea din executarea secvențială a instrucțiunilor unui program.

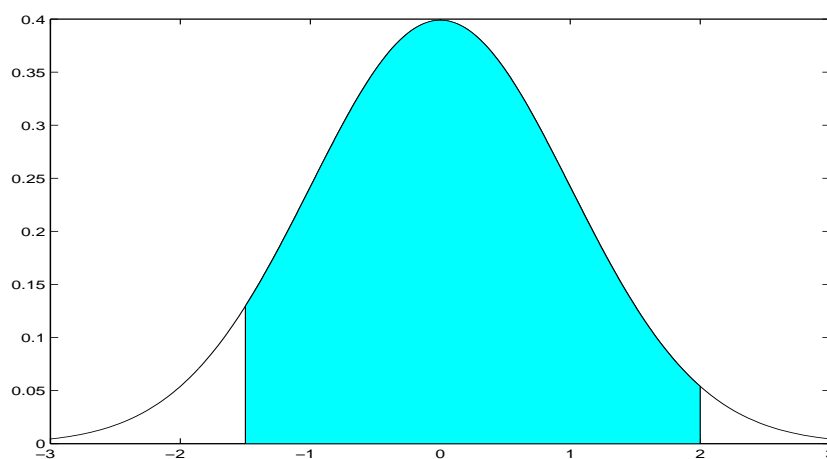


Figura 1.15: Curba lui Gauss

1.8.1 Instrucțiunea `if`

Trei forme de bază ale instrucțiunii `if` sunt:

```
if exp
    secv
end
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%

if exp
    secv
else
    secv1
end
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%

if exp
    secv
elseif exp1
    secv1
else
    secv2
end
```

unde `exp` și `exp1` sunt expresii logice ce pot fi evaluate în acest punct, iar `secv`, `secv1` și `secv2` sunt secvențe de instrucțiuni Matlab.

Instrucțiunea `if` sub prima formă face să se execute secvența de instrucțiuni `secv` numai dacă valoarea expresiei logice `exp` este adevărată, după care se trece la prima instrucțiune ce urmează liniei `end`.

A doua formă de instrucțiune `if` are ca efect executarea uneia dintre secvențele de instrucțiuni `secv` și `secv1`, în funcție de faptul că expresia logică `exp` este adevărată, respectiv falsă, după care se trece la instrucțiunea ce urmează după linia `end`.

Ultima formă are ca rezultat executarea uneia dintre secvențele de instrucțiuni `secv`, `secv1` și `secv2`, după cum valoarea lui `exp` este adevărată, a lui `exp` este falsă și a lui `exp1` este adevărată, respectiv valorile expresiilor logice `exp` și `exp1` sunt false și `exp2` are valoare adevărată.

De exemplu secvența de instrucțiuni ce urmează stabilește dacă elementele unui vector `v` cu `m` componente întregi sunt toate numere pare, toate impare sau sunt și pare și impare:

```
if rem(v,2)==0
    p=0;
elseif rem(v,2)==1
    p=1;
else
    p=2;
end
p
```

Rezultatul execuției acestei secvențe atribuie lui `p` una din valorile 0, 1, 2, în funcție de faptul că vectorul `v` are toate componentele numere pare, toate impare, respectiv și pare și impare.

1.8.2 Instrucțiunea `switch`

Această instrucțiune are ca rezultat trecerea execuției, în funcție de valoarea unei expresii, în diferite puncte ce urmează.

Forma generală a instrucțiunii `switch` este:

```
switch exp
case val1
    secv1
case val2
    secv2
. . . . .
otherwise
    secv
end
```

unde `exp` este o expresie aritmetică sau de tip caracter, iar `secv1`, `secv2`, ..., `secv` sunt secvențe de instrucțiuni Matlab. Dacă valoarea expresiei `exp` coincide cu una din valorile `val1`, `val2`, ... atunci se execută secvența de instrucțiuni corespunzătoare `secv1`, `secv2`, ..., iar în caz contrar se execută secvența de instrucțiuni `secv`.

De exemplu, secvența de instrucțiuni


```

m=input('m=');
switch expr
case 'ones'
    A=ones(m)
case 'eye'
    A=eye(m)
case 'magic'
    A=magic(m)
otherwise
    A=zeros(m)
end

```

generează matricea pătratică A de ordin m, care poate fi de la caz la caz, respectiv matrice cu toate elementele 1, matrice unitate, matrice magică, matrice cu toate elementele zero.

1.8.3 Instrucțiunea while

Instrucțiunea while este o instrucțiune de ciclare prin care se repetă execuția unei secvențe de instrucțiuni Matlab până când o condiție este satisfăcută. Forma generală a instrucțiunii este:

```

while exp
    secv
end

```

și are ca efect executarea secvenței de instrucțiuni secv atâta timp cât valoarea expresiei logice exp este adevărată.

De exemplu, dacă vrem să calculăm valoarea aproximativă a soluției ecuației $x = \cos(x)$, folosind metoda iterativă, impunând condiția de oprire prin aceea ca numărul iterațiilor să nu depășească o valoare dată n și distanța dintre două iterații consecutive să nu fie mai mare decât un $\epsilon > 0$ dat, atunci avem programul

```

n=input('n=');
epsilon=input('epsilon=');
xv=pi/4; k=1; d=1;
while d>epsilon & k<n
    k=k+1; xn=cos(xv);
    d=abs(xn-xv); xv=xn;
end
fprintf('k=%2d\n',k)
fprintf('d=%10.4f\n',d)
fprintf('x=%10.4f\n',xn)

```

Astfel, pentru valorile la intrare $n=50$ și $\epsilon=1e-5$, sunt afișate rezultatele:

```

k=25
d=    0.0000
x=    0.7391

```

1.8.4 Instrucțiunea for

Față de instrucțiunea `while`, instrucțiunea de ciclare `for` are fixat numărul repetărilor unei secvențe de instrucțiuni.

Forma generală a instrucțiunii este:

```
for v=vi:pas:vf
    secv
end
```

unde v este variabila de ciclare, care ia toate valorile, începând cu valoarea expresiei aritmetice vi până la valoarea expresiei vf , folosind pasul dat prin expresia aritmetică pas . Pentru fiecare din aceste valori ale lui v se execută succesiv secvența `secv` de instrucțiuni Matlab.

Dacă $pas=1$, atunci putem folosi următoarea formă a instrucției `for`:

```
for v=vi:vf
    secv
end
```

Mai remarcăm faptul că în interiorul unui ciclu `for` poate exista un alt ciclu `for` ș.a.m.d.

De exemplu, secvența următoare va genera matricea h de tipul (m, n) , numită *matricea lui Hilbert*:

```
m=input('m=');
n=input('n=');
h=[];
for i=1:m
    for j=1:n
        h(i,j)=1/(i+j-1);
    end
end
format rat
disp(h)
```

În urma executării acestui program, cu valorile $m=3$ și $n=6$, se obțin rezultatele

```
m=3
n=4
      1      1/2      1/3      1/4
1/2      1/3      1/4      1/5
1/3      1/4      1/5      1/6
```

Mai considerăm un exemplu de folosire a instrucțiunii `for`. Vrem să construim *matricea Vandermonde* V , corespunzătoare numerelor $a_i, i = \overline{1, m}$, care sunt date

prin vectorul a . Reamintim că forma generală a acestei matrice este

$$V = \begin{pmatrix} 1 & 1 & \dots & 1 \\ a_1 & a_2 & \dots & a_m \\ a_1^2 & a_2^2 & \dots & a_m^2 \\ \dots & \dots & \dots & \dots \\ a_1^{m-1} & a_2^{m-1} & \dots & a_m^{m-1} \end{pmatrix}.$$

Matricea V poate fi generată prin secvența de program următoare:

```
m=input('m=');
for i=1:m
    a(i)=input('a(i)=');
end
V=[];
for i=0:m-1
    V=[V;a.^i];
end
disp(V)
```

Se observă că matricea V a fost inițializată cu *matricea vidă*, care este specificată prin `[]`, iar construcția efectivă s-a făcut prin operația de concatenare.

Executând programul pentru $m=4$ și $a(i) = i, i = \overline{1,4}$, se obține

```
m=4
a(i)=1
a(i)=2
a(i)=3
a(i)=4
```

1	1	1	1
1	2	3	4
1	4	9	16
1	8	27	64

1.9 Instrucțiuni de întrerupere

1.9.1 Instrucțiunea `try...catch`

Forma generală a instrucțiunii este:

```
try
    secv
catch
    secv1
end
```

și are ca efect execuția primei secvențe de instrucțiuni, `secv`, până la apariția unei erori, după care se execută a doua secvență de instrucțiuni, `secv1`. Dacă nici o

eroare nu apare în execuția secvenței `sevc`, după executarea acesteia, se trece la prima instrucțiune după linia `end`. Dacă la execuția secvenței de instrucțiuni `sevc1` apare o eroare, execuția programului este oprită, dacă o altă instrucțiune de tipul `try...catch` nu există.

1.9.2 Instrucțiunea `pause`

Instrucțiunea `pause` are ca efect întreruperi temporare ale execuției programelor.

Instrucțiunea `pause` are următoarele forme:

```
pause(n)
pause
pause off
pause on
```

Prima formă are ca efect întreruperea executării programului timp de n secunde (n este număr pozitiv), când se ajunge cu executarea în acest punct.

A doua formă, `pause` fără argument, întrerupe executarea programului, reluarea executării făcându-se prin apăsarea oricărei taste.

Instrucțiunea `pause` cu argumentul `off` are ca efect inhibarea tuturor instrucțiunilor `pause` și `pause(n)` ce urmează, iar dacă are argumentul `on`, atunci se produce reactivarea acestora.

1.9.3 Instrucțiunea `return`

Instrucțiunea

```
return
```

are ca efect întoarcerea în programul de pe nivelul imediat superior (programul care apelează unitatea în care se află instrucțiunea). Nivelul cel mai înalt este cel care corespunde modului interactiv de lucru, prin urmare instrucțiunea `return`, la acest nivel nu are niciun efect.

1.9.4 Instrucțiunea `break`

Instrucțiunea

```
break
```

are efecte diferite în funcție de locul unde este poziționată.

Dacă instrucțiunea `break` se află într-un ciclu `while`, se produce ieșirea forțată din ciclul respectiv, iar dacă se află într-un ciclu `for` se produce la fel o ieșire forțată din ciclu, cu observația că dacă există un ciclu `for` superior, se continuă cu executarea acestuia.

Dacă instrucțiunea `break` se află într-o instrucțiune de tipul `if`, `switch` sau `try...catch`, atunci se abandonează execuția acestei instrucțiuni, iar dacă este poziționată în oricare altă parte, se încheie executarea programului.

1.9.5 Instrucțiunea `error`

Instrucțiunea

```
error('mesaj')
```

are ca efect abandonarea execuției și afișarea mesajului precizat prin `mesaj`.

1.10 Funcții (proceduri) în Matlab

Toate programele scrise până aici, numite *fișiere script*, ar putea fi considerate ca fiind *programe principale*. Un astfel de program principal poate să apeleze alte programe, pe care le numim *funcții* sau *proceduri*, care la rândul lor pot face apel la alte proceduri.

Un program principal (fișier script), nu poate fi apelat de o altă unitate de program. Lansarea execuției unui astfel de program, se poate face numai când sistemul Matlab este în mod interactiv, adică avem pe ecran cursorul `>>` al sistemului Matlab. Lansarea se face în două moduri:

- prin tastarea, după cursorul `>>`, a numelui său, fără extensia `.m`;
- prin copierea, după cursorul `>>`, a conținutului fișierului ce conține programul, un exemplu fiind cuplul de comenzi `Copy – Paste`,

după care se tastează comanda `Enter`.

Remarcăm faptul că toată istoria unei sesiuni se află pe ecran, chiar dacă unele secvențe mai vechi nu sunt vizibile la un moment dat, dar cu ajutorul *mouse*-ului, folosind *bara* din dreapta a ecranului se poate aduce secvența dorită în partea vizibilă a ecranului. Mai mult, cu ajutorul sistemului de săgeți pot fi readuse comenzi mai vechi sau mai noi, care pot fi reeditate (modificate) și lansate după o astfel de reeditare. Există chiar o căutare mai rapidă a unor instrucțiuni din lista de instrucțiuni ce au fost utilizate într-o sesiune Matlab, anume se tastează începutul instrucțiunii urmată de săgeată în sus.

Încă ceva, toate variabilele folosite în timpul unei sesiuni Matlab, fie în mod interactiv, fie prin execuția unui program principal (*script-file*) au un caracter global și sunt păstrate într-un spațiu de lucru (*workspace*), care la încheierea sesiunii Matlab sunt pierdute. De aceea, avem la dispoziție comenzile pe care le-am prezentat, `diary` și `save` care permit, prin lansarea lor, salvarea acestor informații.

1.10.1 Definirea și structura unei funcții Matlab

Colecția de programe ale sistemului Matlab este constituită, în mare parte, din funcții (proceduri) scrise în limbajul propriu. Acestea sunt fișiere cu extensia `.m` având următoarea structură:

```
function [lista1] = numef(lista2)
% linia H1 (help line)
% comentariu
%
corp
```

unde `lista1` și `lista2` reprezintă respectiv lista variabilelor de ieșire (despărțite prin spații sau virgule) și lista parametrilor de intrare (despărțiți prin spații sau virgule), `corp` este secvența de instrucțiuni ale funcției cu numele `numef`.

Variabilele din corpul procedurii au un caracter local, prin urmare la întoarcerea în unitatea care face apelul valorile acestora se pierd, desigur cu excepția valorilor din `lista1`. Există posibilitatea schimbării caracterului local, prin instrucțiunea

```
global lista
```

care are ca efect obținerea caracterului global pentru variabilele din `lista`.

În corpul unei funcții Matlab variabilele `nargin` și `nargout` specifică numărul parametrilor de intrare, adică lungimea listei `lista2`, respectiv a parametrilor de ieșire, adică lungimea listei `lista1`.

Dacă într-un fișier script sau de la tastatură se dă una din comenzile

```
nargin('numef')
nargout('numef')
```

atunci se obțin numărul parametrilor de intrare și respectiv numărul parametrilor de ieșire pentru funcția Matlab cu numele `numef`.

Liniile dinaintea corpului procedurii, care încep cu `%` sunt linii de comentariu și vor fi listate când se lansează comanda

```
help numef
```

Din acest motiv, partea de comentariu de la începutul procedurii este bine să conțină informații, care să descrie pe scurt cum trebuie să fie apelată și folosită funcția respectivă.

Mai remarcăm prima linie de comentariu, notată cu `H1`, care are o însușire aparte. Când se lansează comanda

```
lookfor key
```

care caută cuvântul cheie `key`, aceasta se face numai în liniile `H1` ale fișierelor cu extensia `.m`, după care se listează numele fișierelor ce conțin cuvântul cheie precizat.

În cadrul liniilor ce alcătuiesc `corp`-ul funcției pot exista de asemenea linii de comentariu, care încep cu simbolul `%`, dar care nu vor fi afișate odată cu lansarea comenzii `help`. Mai mult apariția simbolului `%` în interiorul unei linii face ca tot ce urmează pe linia respectivă să fie considerat comentariu.

Numele unei funcții Matlab, `numef`, scrisă de utilizator, este de dorit să fie diferit de numele funcțiilor deja existente în sistemul Matlab, iar numele fișierului, cu extensia `.m`, ce va conține funcția să fie `numef.m`.

Dacă lista parametrilor de intrare, `lista2`, este vidă, atunci linia de definiție a funcției este de forma

```
function [lista1] = numef
```

Dacă lista variabilelor de ieșire, `lista1`, este formată dintr-o singură variabilă, atunci se poate renunța la parantezele drepte, iar dacă este vidă se poate elimina complet.

Cu aceste precizări, cea mai simplă linie de definiție a unei funcții (proceduri) este

```
function numef
```

1.10.2 Apelul unei funcții (proceduri) Matlab

Apelarea unei funcții se poate face în două moduri.

Primul mod este cel interactiv, de exemplu, prin una din formele

```
[apel1] = numef(apel2)
numef(apel2)
```

unde `apel2` este lista parametrilor actuali, care poate fi vidă sau formată din cel mult atâtea elemente cât are lista parametrilor formali, `lista2`, din linia de definiție a funcției `numef`, de fiecare dată făcându-se legătura de tipul: primul parametru din `apel2` corespunde primului parametru din `lista2`, al doilea element din `apel2` corespunde celui de al doilea parametru din `lista2`,..., corespondența încheindu-se când se termină lista `apel2`. Desigur, dacă `apel2` are mai puține elemente decât `lista2`, atunci trebuie gestionată această situație în corpul funcției `numef` cu ajutorul variabilei `nargin`.

Rezultatele executării funcției `numef`, se păstrează în lista `apel1` în primul caz, respectiv în `ans` în al doilea caz. O discuție analogă asupra corespondenței dintre `apel2` și `lista2`, se face și în cazul `apel1` și `lista1`.

O altă cale de apel a unei funcții este dintr-un alt program sau procedură Matlab, prin apariția într-o instrucțiune Matlab a numelui funcției, împreună cu lista parametrilor actuali de tipul listei `apel2`, mai sus precizată.

1.10.3 Subfuncții

Funcțiile pot avea în corp-ul lor apeluri la alte funcții. Totdeauna întoarcerea în unitatea care cheamă se face, când se ajunge la linia finală a corp-ului unității chemate sau când se întâlnește instrucțiunea `return`.

Există două situații mai speciale în construcția funcțiilor în Matlab.

În primul rând, există posibilitatea ca după `corp`-ul unei funcții să fie definite una sau mai multe funcții, pe care le numim *subfuncții* și care pot fi apelate din `corp`-ul funcției. Forma unei astfel de funcții ar fi:

```
function [lista1] = numef(lista2)
%
% comentariu
%
corp
function [list1] = numef1(list2)
%
% comentariul
%
corp1
```

în secvența de instrucțiuni `corp` făcându-se apel la funcția `numef1`. Funcțiile de tipul funcției `numef1` nu se pot apela din exteriorul fișierului `numef.m`. Drept consecință nici `comentariul` nu este vizibil din exteriorul funcției `numef`.

Funcția 1.10.1. Să scriem o procedură, care reprezintă grafic pe intervalul $[-1, 1]$ o mulțime finită precizată de polinoame ale lui Cebîșev de speța I. Dacă numărul acestora este par, atunci să fie reprezentate pe două coloane, iar în caz contrar pe o singură coloană.

```
function cebisev(lista)
% Polinomul Cebisev
% lista -- contine gradele polinoamelor Cebisev
%         care vor fi reprezentate grafic;
% daca lista are un numar par de elemente,
% polinoamele se vor reprezenta pe doua coloane;
% daca lista are un numar impar de elemente,
% polinoamele se vor reprezenta pe o coloana;
%
clf, x=-1:0.001:1; r=length(lista);
if rem(r,2)==0
    m=fix(r/2); n=2;
else
    m=r; n=1;
end
for i=1:r
    y=cos(lista(i)*acos(x));
    splots(m,n,i,x,y,lista(i))
end
function splots(m,n,k,u,v,gr)
subplot(m,n,k), plot(u,v,'k-')
title(['Polinom Cebisev de gradul n=',num2str(gr)])
```

Dacă se apelează funcția `cebisev` prin

```
cebisev([2:5])
```

se obțin graficele din Figura 1.16.

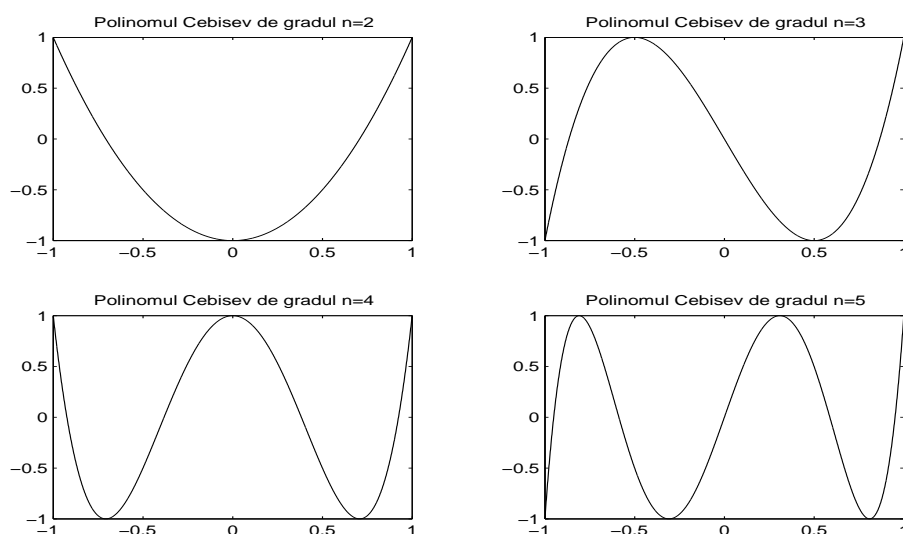


Figura 1.16: Polinoamele lui Cebîșev de grad $n = 2, 3, 4, 5$

1.10.4 Funcția feval

Un alt caz special este acela când printre parametrii formali de intrare din `lista2` există și parametri ce pot lua valori nume ale unor funcții. La apelul unei astfel de funcții pe poziția parametrului actual corespunzător din lista `apel2` trebuie să se afle un șir de caractere ce denumește funcția actuală.

Funcția 1.10.2. Funcția cu numele `opmatrix` operează prin operatorii Matlab cunoscuți asupra a două matrice:

```
function Z = opmatrix(op,X,Y)
%
% Efectueaza operatia definita prin op
% asupra matricelor X si Y.
% Rezultatul este returnat in Z
%
Z = feval(op,X,Y);
```

Folosind această funcție, de exemplu prin comenzile

```
opmatrix('plus',A,B)
opmatrix('minus',A,B)
opmatrix('and',A,B)
```

sau

```
opmatrix('+',A,B)
opmatrix('-',A,B)
opmatrix('&',A,B)
```

se obțin suma, diferența și respectiv conjuncția matricelor A și B.

1.10.5 Comanda **echo**

Instrucțiunea **echo** are un comportament puțin diferit pentru cele două tipuri de fișiere cu extensia **.m**.

Dacă suntem în cazul unui *fișier script*, atunci instrucțiunile

```
echo on
echo off
echo
```

au respectiv următoarele efecte: prima are ca efect afișarea pe ecran a fișierului în derulare, a doua anulează acest efect, iar a treia are efectul primei forme, dacă sistemul se află în mod **off**, iar dacă se află în mod **off** se va face trecerea la modul **on**.

Dacă suntem în cazul funcțiilor Matlab cu extensia **.m**, atunci avem formele

```
echo numef on
echo numef off
echo numef
echo on all
echo off all
```

Primele trei forme sunt analoge celor trei din cazul *script*, numai că se mai specifică și numele funcției la care se referă. Celelalte două acționează asupra tuturor funcțiilor.

1.11 Instrucțiuni de evaluare a eficienței

Compararea eficienței algoritmului este posibilă în sistemul Matlab fie prin numărarea operațiilor în virgulă flotantă, fie prin determinarea duratei, în secunde, a execuției unei secvențe de instrucțiuni, CPU (Central Processor Unit).

1.11.1 Instrucțiunea **flops**

Instrucțiunea **flops** este încorporată în Matlab 6 prin pachetul LAPACK. Prin urmare este funcțională numai dacă sistemul conține și LAPACK.

Obținerea numărului operațiilor în virgulă flotantă efectuate pentru execuția unei secvențe **secv** de instrucțiuni Matlab se realizează prin:

```
flops(0)
secv
flops
```

La încheierea execuției acestei succesiuni de instrucțiuni în variabila **flops** se află numărul acestor operații.

1.11.2 Instrucțiunile `tic` și `toc`

Pentru obținerea duratei execuției unei secvențe `secv` de instrucțiuni Matlab avem următoarea succesiune de instrucțiuni:

```
tic
    secv
toc
```

unde `tic` marchează începutul contorizării timpului, iar `toc` anunță încheierea contorizării, variabila `toc` conținând timpul măsurat în secunde.

1.12 Grafică tridimensională

Grafica tridimensională sau altfel numită 3D este axată în sistemul Matlab pe două direcții: reprezentarea curbelor în spațiu și reprezentarea suprafețelor.

1.12.1 Instrucțiunea `plot3`

Instrucțiunea `plot3` este utilizată pentru reprezentarea curbelor în spațiu și în principiu nu se deosebește mult de instrucțiunea `plot` din grafica bidimensională.

Formele întâlnite pentru această instrucțiune sunt:

```
plot3(x,y,z)
plot3(x,y,z,s)
plot3(x1,y1,z1,x2,y2,z2,...)
plot3(x1,y1,z1,s1,x2,y2,z2,s2,...)
```

Dacă x, y, z sunt vectori de aceeași lungime m , atunci se unesc punctele de coordonate (x_i, y_i, z_i) , $i = \overline{1, m}$, printr-o linie poligonală de tipul precizat prin parametrul s , de la instrucțiunea `plot`, iar în lipsa acestui parametru, sistemul Matlab efectuează o alegere implicită a tipului curbei. Desigur că netezimea curbei este cu atât mai bună cu cât numărul punctelor este mai mare.

Dacă x, y, z sunt matrice de același tip (m, n) , atunci pentru fiecare coloană se obține câte o linie poligonală, adică n linii poligonale. Acestea sunt generate respectiv de punctele de coordonate $(x_{ij}, y_{ij}, z_{ij})_{i=1}^m$, pentru fiecare $j = \overline{1, n}$, cu aceeași observație privind parametrul s , care precizează tipul curbei. În acest caz pe aceeași figură vor fi obținute mai multe grafice.

Ultimele două forme sunt extinderi ale primelor două, pentru a putea controla mai bine în special tipul curbelor cu ajutorul parametrilor de tip s .

Remarcăm de asemenea extinderea acțiunilor comenzilor ce controlează un grafic din cazul bidimensional pentru cazul tridimensional.

Programul 1.12.1. Prezentăm un program care ilustrează mișcarea unei particule din punctul A pe o spirală până în punctul B.

```

pas=input('pas='); h=input('h=');
t=0:pas:h; r=exp(-0.2*t); th=pi*t*0.5;
x=r.*cos(th); y=r.*sin(th); z=t;
plot3(x,y,z), hold on
plot3([1,1],[-0.5,0],[0,0])
text(1,-0.7,0,'A'), n=length(t);
text(x(n),y(n),h+2,'B')
xlabel('x'), ylabel('y'), zlabel('z')

```

Să explicăm pe scurt o parte a instrucțiunilor acestui program.

Figura este obținută în două etape, din cauza aceasta a fost introdusă instrucțiunea `hold on`. Prima instrucțiune `plot3` trasează efectiv graficul spiralei, care se mai completează prin segmentul din planul xoy delimitat de punctele de coordonate $(1, -0.5, 0)$ și $(1, 0, 0)$, trasat cu a doua instrucțiune `plot3`. Cele două instrucțiuni au același rezultat ca și instrucțiunea

```
plot3(x,y,z,[1,1],[-0.5,0],[0,0])
```

doar că spirala va fi trasată cu o anumită culoare, iar segmentul din planul xoy cu alta. Dacă vrem ca întreaga curbă să fie trasată cu aceeași culoare se poate apela la parametrul de tip `s`.

Se mai observă că poziția inițială A a punctului și poziția finală B sunt afișate cu instrucțiunile `text`, iar etichetele celor trei axe de coordonate sunt date prin instrucțiunile de tip `label`.

După executarea programului, cu $h=20$ și $pas=0.01$, s-a obținut graficul din Figura 1.17.

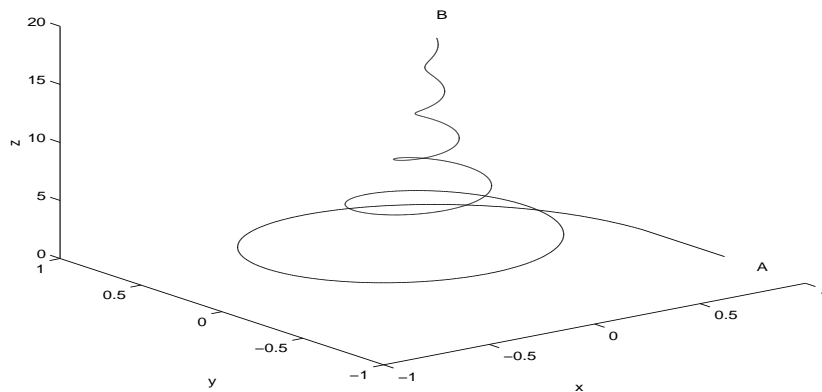


Figura 1.17: Mișcarea pe spirală

1.12.2 Instrucțiunea `ezplot3`

Pentru reprezentarea grafică a curbelor în spațiu, poate fi folosită instrucțiunea `ezplot3`, care se lansează prin una din comenzile:

```
ezplot3('x','y','z')
ezplot3('x','y','z',[a,b])
ezplot3('x','y','z','animate')
ezplot3('x','y','z',[a,b],'animate')
```

Parametrii x , y și z , conțin expresii algebrice ale reprezentărilor curbei în funcție de un parametru t , din $[a, b]$. Valoarea implicită a acestui interval este $[0, 2\pi]$.

Prezența parametrului `animate` produce pe figură un buton cu numele `repeat`, care prin activare prezintă modul de deplasare a unui punct pe curba reprezentată grafic.

De exemplu, dacă se consideră comanda

```
ezplot3('sin(t)','cos(t)','t',[0,10*pi],'animate')
```

se va obține graficul din Figura 1.18, în care se observă și butonul `repeat`.

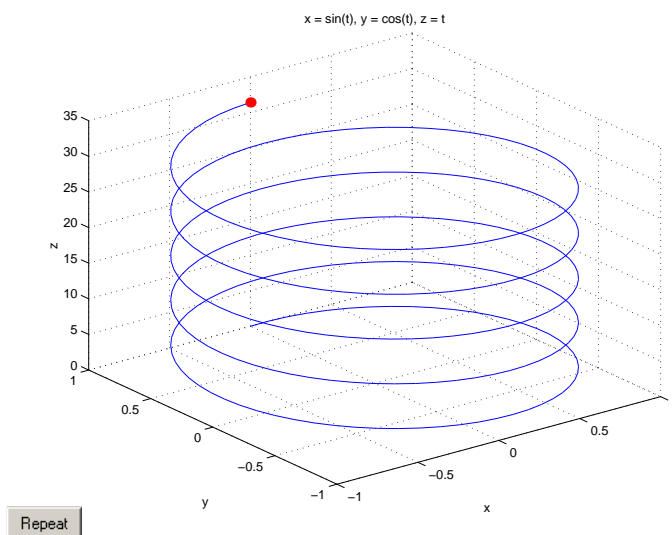


Figura 1.18: $x(t) = \sin(t)$, $y(t) = \cos(t)$, $z(t) = t$, $t \in [0, 10\pi]$

1.12.3 Instrucțiunile `meshgrid`, `mesh` și `surf`

Reprezentarea grafică a unei suprafețe, dată prin funcția $z = f(x, y)$, pe domeniul $D \subset \mathbb{R}^2$, cuprinde două etape.

Prima dată se realizează o rețea rectangulară de puncte, cu ajutorul instrucțiunii `meshgrid`, care are sintaxa

```
[X,Y]=meshgrid(x,y)
```

unde x și y sunt vectori, în general de lungimi diferite, m și n , care conțin puncte de pe axele ox și oy și care sunt folosite la generarea rețelei rectangulare (x_i, y_j) $i = \overline{1, m}, j = \overline{1, n}$.

Rezultatul executării instrucțiunii `meshgrid`, are ca efect generarea matricelor X și Y având același tip (n, m) . Fiecare linie a matricei X este formată din componentele vectorului x , iar fiecare coloană a matricei Y este formată din componentele vectorului y . În acest fel punctul de coordonate (x_i, y_j) al rețelei rectangulare, se poate exprima, cu ajutorul matricelor X și Y , prin perechea $(X(j, i), Y(j, i))$.

Înainte de a trece la a doua etapă, să facem observația că dacă x coincide cu y , atunci se poate utiliza

```
[X,Y]=meshgrid(x)
```

Reprezentarea grafică a funcției $z = f(x, y)$, se face cu ajutorul instrucțiunii `mesh`, care are una din formele:

```
mesh(Z)
mesh(Z,C)
mesh(x,y,Z)
mesh(x,y,Z,C)
mesh(X,Y,Z)
mesh(X,Y,Z,C)
```

De la început, să remarcăm faptul că parametrul C precizează scara culorilor. Când lipsește acest parametru, scara culorilor este dată de parametrul Z , care reprezintă valorile funcției $z = f(x, y)$ pe o rețeaua rectangulară. În acest ultim caz culoarea va fi proporțională cu înălțimea z din punctul (x, y) .

De asemenea, trebuie să remarcăm faptul că reprezentarea grafică a suprafeței este de forma unei rețele curbilinie sprijinită pe înălțimile date prin matricea Z .

Parametrii X , Y și Z sunt matrice de aceleași dimensiuni, dacă avem în vedere cum au fost obținute matricele X și Y , atunci toate aceste trei matrice sunt de tipul (n, m) .

Parametrii x și y sunt vectori de tipul celor ce generează matricele X și Y prin instrucțiunea `meshgrid`, adică dacă aceștia au lungimile m și respectiv n , matricea Z va fi de tipul (n, m) , iar punctele pe care se sprijină rețeaua curbilinie, care reprezintă graficul funcției, au coordonatele $(x(j), y(i), Z(i, j))$.

Dacă primii doi parametri lipsesc, adică suntem în cazul primelor două forme ale instrucțiunii `mesh`, acestea sunt echivalente respectiv cu următoarele două, prin considerarea valorilor implicite $x=1:m$ și $y=1:n$.

Instrucțiunea `surf` are aceeași sintaxă ca instrucțiunea `mesh`, adică

```
surf(Z)
surf(Z,C)
```

```
surf(x,y,z)
surf(x,y,z,c)
surf(X,Y,Z)
surf(X,Y,Z,C)
```

unde parametrii au aceleași caracteristici. Deosebirea în reprezentarea grafică este că prin instrucțiunea `surf` se produce o colorare a patruleterelor obținute prin curbele care generează suprafața corespunzătoare funcției $z = f(x, y)$.

Programul 1.12.2. Vom scrie un program care să reprezinte grafic funcția

$$z = f(x, y) = \frac{1}{2\pi} e^{-\frac{1}{2}(x^2+y^2)},$$

pe domeniul $D = [-1, 1] \times [-1, 1]$, în două figuri alăturate, odată cu instrucțiunea `mesh` și apoi cu instrucțiunea `surf`.

```
clear, clf
m=input('m=');
h=1/m; x=-1:h:1;
[X,Y]=meshgrid(x);
Z=1/(2*pi)*exp(-0.5*(X.^2+Y.^2));
subplot(1,2,1), mesh(X,Y,Z)
title('Grafica prin MESH')
xlabel('x'), ylabel('y'), zlabel('z')
subplot(1,2,2), surf(X,Y,Z)
title('Grafica prin SURF')
xlabel('x'), ylabel('y'), zlabel('z')
```

Executarea programului, cu $m = 6$, conduce la reprezentările grafice din Figura 1.19.

În încheierea acestui paragraf, amintim că există posibilitatea de umbrire a unei suprafețe cu ajutorul comenzii

```
shading tip
```

unde `tip` poate fi:

- `faceted` (implicit) – umbrește fiecare patruleter de pe suprafața trasată cu o intensitate fixă și trasează laturile patruleterelor;
- `flat` – umbrește fiecare patruleter de pe suprafața trasată cu o intensitate fixă, fără trasarea laturilor patruleterelor;
- `interp` – umbrirea fiecărui patruleter de pe suprafața trasată se face în mod gradat, printr-un procedeu de interpolare, care asigură o trecere netedă de la un patruleter la altul, fără trasarea laturilor patruleterelor.

Programul 1.12.3. Să scriem un program care efectuează umbrirea suprafeței reprezentată grafic prin Programul 1.12.2:

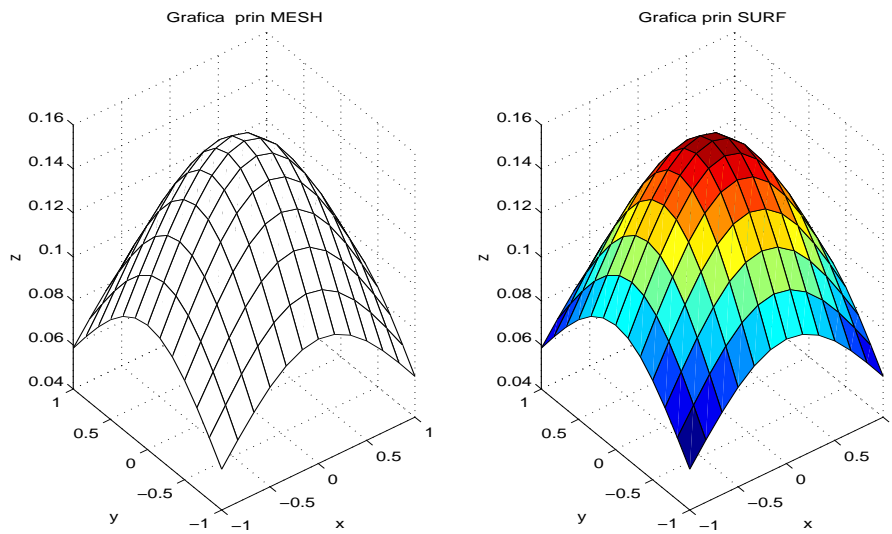


Figura 1.19: $z = \frac{1}{2\pi}e^{-\frac{1}{2}(x^2+y^2)}$

```
m=input('m=');
h=1/m; x=-1:h:1;
[X,Y]=meshgrid(x);
Z=1/(2*pi)*exp(-0.5*(X.^2+Y.^2))
subplot(1,3,1), surf(X,Y,Z),
shading faceted
title('Umbrire - faceted')
subplot(1,3,2), surf(X,Y,Z),
shading flat
title('Umbrire - flat')
xlabel('x'), ylabel('y'), zlabel('z')
subplot(1,3,3), surf(X,Y,Z),
shading interp
title('Umbrire - interp')
```

Rezultatul executării programului este în Figura 1.20.

1.12.4 Instrucțiunile `contour` și `contourf`

Reprezentarea curbelor de nivel ale suprafeței, dată prin funcția $z = f(x, y)$, se poate face cu una din următoarele forme ale instrucțiunii `contour`:

```
contour(Z)
contour(Z,n)
contour(Z,v)
contour(X,Y,Z)
```

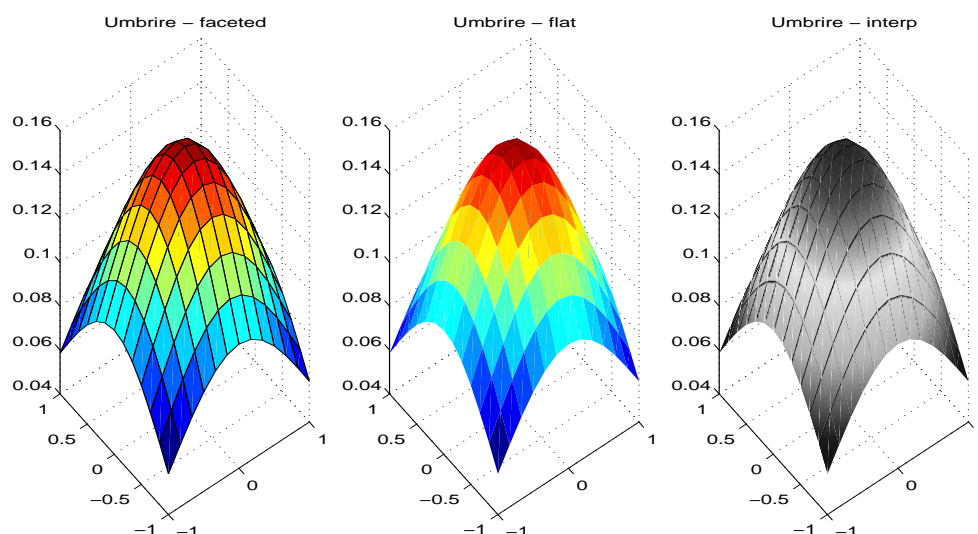



Figura 1.20: Umbrirea unei suprafețe

```

contour(X,Y,Z,v)
contour(Z,[w w])
contour(X,Y,Z,[w w])

```

unde parametrii X , Y și Z au aceleași interpretări ca și în cazul instrucțiunilor `mesh` și `surf`.

Parametrul v este un vector prin care se specifică nivelurile curbelor de nivel ce urmează a fi reprezentate grafic, iar în absența sa, valorile acestui parametru sunt automat calculate. Numărul curbelor de nivel se poate preciza prin parametrul n .

Formele instrucțiunii `contour` cu parametrul de tipul `[w w]` are ca efect trasarea curbei de nivel precizată prin w . Aceste forme permit reprezentarea grafică a funcțiilor date sub formă implicită. Astfel, dacă vrem să reprezentăm grafic funcția $y = y(x)$ dată prin $f(x, y) = 0$, atunci graficul funcției $y = y(x)$ va fi dat de curba de nivel pentru funcția $z = f(x, y)$, considerând $w=0$.

Instrucțiunea `contourf` diferă de `contour` doar prin faptul că ariile delimitate de curbe de nivel sunt umbrite.

Programul 1.12.4. Să reprezentăm n curbe de nivel pentru funcția considerată în Programul 1.12.2. Reprezentările acestor curbe de nivel să fie făcute atât cu instrucțiunea `contour`, cât și cu `contourf`.

```

clear, clf
m=input('m='); n=input('n=');
h=1/m; x=-1:h:1;

```

```
[X,Y]=meshgrid(x);
Z=1/(2*pi)*exp(-0.5*(X.^2+Y.^2));
subplot(1,2,1), contour(Z,n)
title(['n=',num2str(n),' curbe de nivel'])
subplot(1,2,2), contourf(Z,n)
title(['n=',num2str(n),' curbe de nivel umbrite'])
```

Executarea programului, cu $m = 6$ și $n=5$, conduce la reprezentările grafice din Figura 1.21.

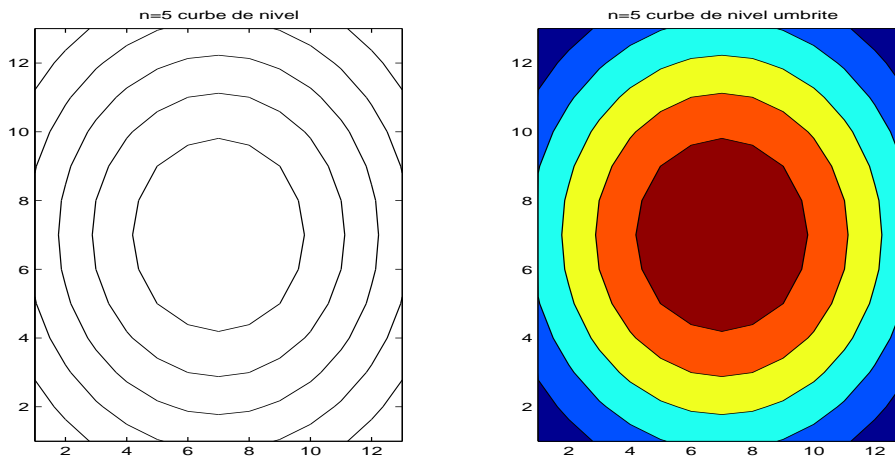


Figura 1.21: Curbe de nivel

1.12.5 Instrucțiunile ezcontour și ezcontourf

O formă simplificată, pentru reprezentarea curbilor de nivel, pentru o funcție dată prin $z = f(x, y)$, se obține cu ajutorul instrucțiunii ezcontour, având una din formele:

```
ezcontour(f)
ezcontour(f,lim)
ezcontour(f,n)
```

unde f este expresia matematică a unei funcții de două variabile, privită ca un șir de caractere.

Parametrul `lim` este fie un vector de forma `[xmin,xmax]`, fie de forma `[xmin,xmax,ymin,ymax]`. Dacă lipsește, se ia implicit $[-2\pi, 2\pi, -2\pi, 2\pi]$, iar dacă este de prima formă se ia `[xmin,xmax,xmin,xmax]`.

Parametrul `n` specifică faptul că în reprezentarea grafică se folosește o rețea de $n \times n$ puncte (implicit se consideră $n=60$).

De exemplu,

```
f='sqrt(x^2+y^2)'  
ezcontour(f,[-2,2])
```

are același efect cu

```
ezcontour('sqrt(x^2+y^2)',[-2,2])
```

și produc curbele de nivel pentru funcția $z = f(x, y) = \sqrt{x^2 + y^2}$, pe domeniul $D = [-2, 2] \times [-2, 2]$.

Instrucțiunea `ezcontourf`, față de `ezcontour` execută și umbrirea ariilor generate de curbele de nivel.

1.12.6 Instrucțiunile `ezmesh`, `ezsurf`, `ezmeshc` și `ezsurf c`

Sintaxa apelurilor funcțiilor `ezmesh`, `ezsurf`, `ezmeshc` și `ezsurf c` este aceeași, iar efectul executării lor este reprezentarea grafică a unei suprafețe, respectiv a unei suprafețe, împreună cu sau fără curbele de nivel corespunzătoare. De aceea prezentăm modurile de apelare, numai pentru una dintre ele.

```
ezmesh('f')  
ezmesh('f',lim)  
ezmesh('x','y','z')  
ezmesh('x','y','z',lim)
```

Parametrul `f` reprezintă o expresie algebrică, ce definește suprafața, care urmează să fie reprezentată grafic, pe domeniul definit prin `lim`. Implicit, în absența parametrului `lim`, domeniul este $[-2\pi, 2\pi] \times [-2\pi, 2\pi]$. Dacă `lim=[a,b]`, atunci domeniul va fi $[a,b] \times [a,b]$, iar dacă `lim=[a,b,c,d]`, atunci domeniul va fi $[a,b] \times [c,d]$.

Parametrii `x`, `y`, `z`, reprezintă expresii algebrice ce definesc suprafața în modul parametric. În acest caz, corespunzător, parametrul `lim`, precizează domeniul în care iau valori parametrii cu ajutorul cărora se definește în mod parametric suprafața ce urmează a fi reprezentată.

Remarcăm faptul că formele mai sus prezentate, mai pot avea doi parametri. Unul, `n`, de tip întreg pozitiv, care precizează numărul $n \times n$ al punctelor rețelei. Implicit se consideră `n=60`. Un altul este `'circ'`, care prin prezența sa atrage reprezentarea lui `f` pe un disc centrat în domeniul mai sus precizat.

Pentru exemplificare, dacă se execută programul Matlab format din următoarele instrucțiuni:

```
subplot(1,2,1)  
ezmesh('sin(x)*sin(y)',[-pi,pi], 'circ')  
subplot(1,2,2)  
ezmeshc('sin(x)*sin(y)',[-pi,pi], 'circ')  
colormap spring
```

se obțin graficele din Figura 1.22.

Mai remarcăm faptul că `f`, `x`, `y`, `z` pot fi numele unor funcții predefinite sau scrise de utilizator.

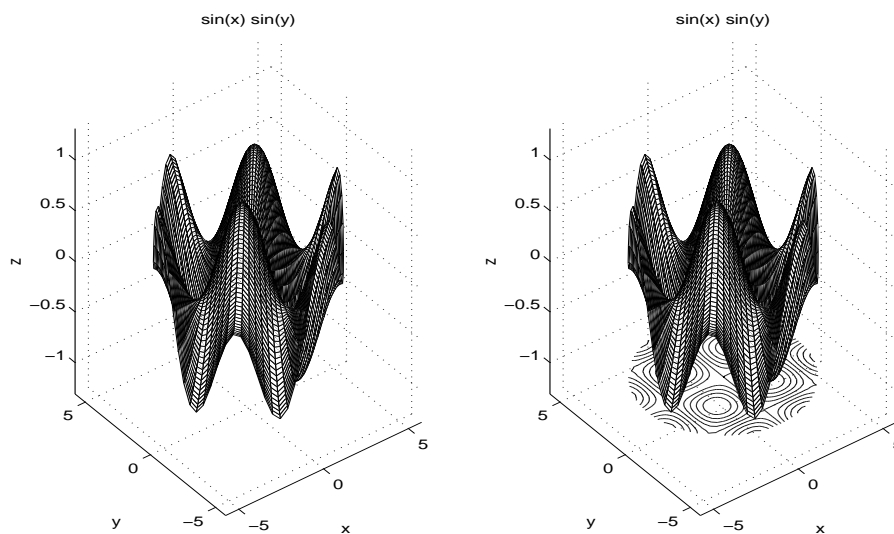


Figura 1.22: Suprafață fără și cu curbe de nivel

Pentru a ilustra că cele două instrucțiuni de tipul `ezmesh` și `ezsurf` au de regulă același efect, programul următor reprezintă grafic aceeași suprafață, respectiv cu `ezmesh` și `ezsurf`.

```
clf
subplot(1,2,1),
ezsurf('x*exp(-x^2 - y^2)',[-2,2],20)
title('Reprezentare cu ezsurf')
subplot(1,2,2)
ezsurf('x*exp(-x^2 - y^2)',[-2,2],20)
colormap autumn
title('Reprezentare cu ezmesh')
```

Graficele sunt prezentate în Figura 1.23.

1.12.7 Instrucțiunile `bar3` și `bar3h`

Instrucțiunile `bar3` și `bar3h` corespund instrucțiunilor `bar` și respectiv `barh` din grafica bidimensională.

Forme ale acestor instrucțiuni sunt:

```
bar3(z,w,'tip','color')
bar3(y,z,w,'tip','color')
bar3h(y,w,'tip','color')
bar3h(y,z,w,'tip','color')
```

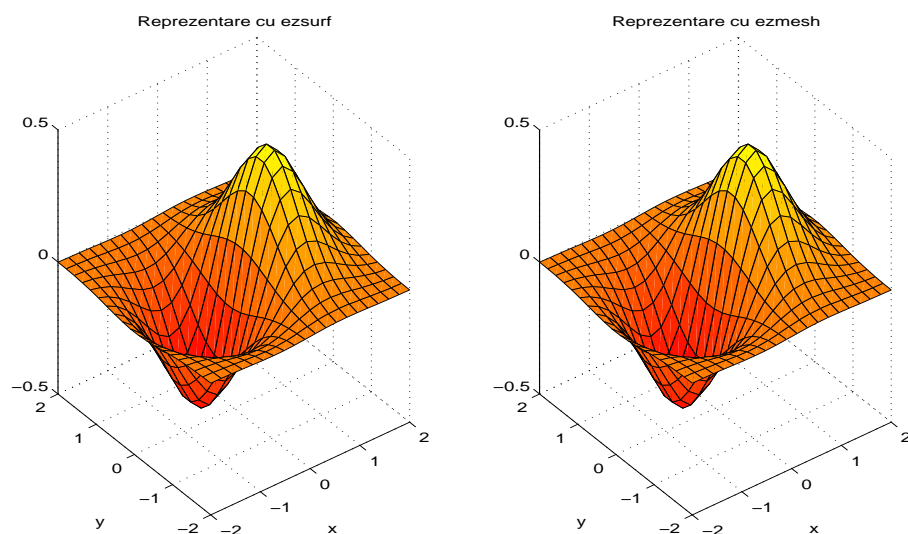


Figura 1.23: Suprafață reprezentată $z = f(x, y) = xe^{-(x^2+y^2)}$

argumentele `w`, `'tip'` și `'color'` sunt opționale, dar când sunt specificate, trebuie să păstreze această ordine.

Parametrul `y` este un vector de lungime `m`, având componentele ordonate crescător sau descrescător. Dacă `y` lipsește, atunci se consideră implicit `y=1:m`.

Dacă `z` este un vector de aceeași lungime cu `y`, atunci vor fi reprezentate grafic `m` bare de lungimi z_i în dreptul punctelor y_i , $i = \overline{1, m}$, de pe axa oy . Dacă `z` este o matrice de tipul (m, n) , reprezentările sunt efectuate pentru fiecare din cele `n` coloane ale matricei, adică în fiecare punct y_i , $i = \overline{1, m}$, de pe axa oy , vor fi reprezentate câte `n` bare (batoane).

Parametrul `w` specifică grosimea barelor, implicit fiind `w = 0.8`, iar pentru `w > 1` barele se suprapun.

Dacă parametrul `tip=detached`, care este valoarea implicită, atunci barele din punctele y_i de pe axa oy vor fi detașate în adâncime după axa ox . Dacă `tip=grouped`, barele (batoanele) vor fi grupate în punctele y_i , de pe axa oy , iar dacă `tip=stacked`, atunci barele vor fi stivuite în aceste puncte.

Parametrul `color` specifică culoarea barelor (batoanelor) și are una din următoarele valori: `r` (roșu), `g` (verde), `b` (albastru), `y` (galben), `m` (violet), `c` (ciclamen), `k` (negru), `w` (alb).

Deosebirea între `bar3` și `bar3h` este că prima instrucțiune acționează pe verticală, iar a doua pe orizontală.

Programul 1.12.5. Să scriem un program care generează o matrice magică de ordin $m > 2$. Folosind primele trei coloane ale matricei magice, să reprezentăm prin bare verticale valorile acestora cu fiecare din tipurile `detached`, `grouped` și `stacked`. Folosind ultimele trei coloane să se facă astfel de reprezentări cu ajutorul barelor orizontale.

```
m=input('m:'); A=magic(m);
subplot(2,3,1), bar3(A(:,1:3),.6,'detached')
title('Bar3 - detached')
subplot(2,3,2), bar3(A(:,1:3),'grouped')
title('Bar3 - grouped')
subplot(2,3,3), bar3(A(:,m-2:m),.6,'stacked')
title('Bar3 - stacked')
subplot(2,3,4), bar3h(A(:,m-2:m),.6,'detached')
title('Bar3h - detached')
subplot(2,3,5), bar3h(A(:,m-2:m),'grouped')
title('Bar3h - grouped')
subplot(2,3,6), bar3h(A(:,1:3),.6,'stacked')
title('Bar3h - stacked')
colormap spring
```

Executarea acestui program, pentru $m=4$, produce graficele din Figura 1.24.

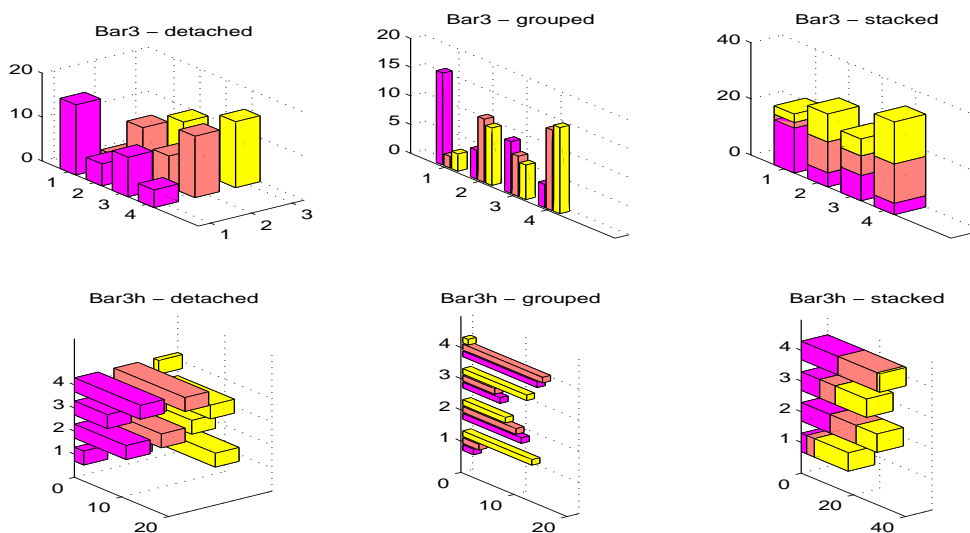


Figura 1.24: Bare verticale și bare orizontale în 3D

Capitolul 2

Elemente de teoria probabilităților

Teoria probabilităților are ca obiect de studiu fenomene nedeterminate sau aleatoare (întâmplătoare). Un fenomen nedeterminist \mathcal{E} este rezultatul imposibilității cunoașterii exacte și complete a tuturor elementelor care concură la realizarea acestuia. Astfel de fenomene aleatoare ar putea fi apelurile dintr-o centrală telefonică, accidentele de circulație dintr-o intersecție, extragerile numerelor loto, rezultatele de la ruletă, noii născuți la o maternitate, numărul particulelor dezintegrate dintr-o substanță radioactivă, mișcarea browniană a particulelor unei substanțe dizolvată într-un fluid, sosirea clienților într-o stație de servire, răspândirea erorilor de tipar dintr-o carte, vânzarea pâinii la o unitate alimentară și în fine răspândirea stafidelor în prăjitura cumpărată de la o cofetărie. Toate aceste fenomene sunt fenomene aleatoare prin faptul că nu sunt cunoscute toate condițiile în care sunt realizate, sau dacă o parte din aceste condiții pot fi precizate, ele nu sunt cunoscute cu exactitate. Factorii necunoscuți, care se regăsesc în fenomenul real, conduc la astfel de fenomene aleatoare (nedeterminate).

Teoria probabilităților vine să dea metode unitare pentru analiza și studiul unor astfel de fenomene nedeterminate și care provin dintr-o gamă largă a activității și cunoașterii umane. Pentru aceasta este necesar să fie formalizat limbajul prin care se tratează fenomenele aleatoare, ca apoi prin particularizări să fie obținute rezultate specifice fenomenului considerat.

Vom presupune că avem anumite cunoștințe de relative la teoria probabilităților. Reamintim în acest capitol doar o parte din noțiunilor de bază ale teoriei probabilităților, pentru stabilirea unui limbaj comun în abordarea capitolelor următoare. Ținând seama de această motivație recomandăm celor neinițiați în această disciplină a teoriei probabilităților lucrările [10], [6], [9], [8].

2.1 Câmp de probabilitate

Dacă se consideră fenomenul aleator (*experimentul*) \mathcal{E} , vom numi *probe* rezultatele experimentului, iar mulțimea probelor relative la experimentul \mathcal{E} o numim *spațiul probelor* și pe care îl notăm prin Ω .

Definiția 2.1.1. Spunem că submulțimea nevidă $\mathcal{K} \subset \mathcal{P}(\Omega)$ este un σ -corp (corp borelian, σ -algebră) dacă satisface următoarele condiții:

$$(2.1.1) \quad \forall A \in \mathcal{K} \implies \bar{A} \in \mathcal{K}; \quad \forall A_i \in \mathcal{K}, i \in I \implies \bigcup_{i \in I} A_i \in \mathcal{K},$$

iar perechea (Ω, \mathcal{K}) se numește câmp de evenimente.

Mulțimea de indici I este o mulțime cel mult numărabilă.

Definiția 2.1.2. Fiind dat câmpul de evenimente (Ω, \mathcal{K}) , se numește probabilitate o aplicație $P: \mathcal{K} \longrightarrow \mathbb{R}$, care verifică următoarele trei axiome:

$$(i) \quad \forall A \in \mathcal{K} \implies P(A) \geq 0,$$

$$(ii) \quad P(\Omega) = 1, \quad \Omega \text{ fiind evenimentul sigur (cert),}$$

$$(iii) \quad \forall A_i \in \mathcal{K}, i \in I, A_i \cap A_j = \emptyset, i \neq j \implies P\left(\bigcup_{i \in I} A_i\right) = \sum_{i \in I} P(A_i),$$

iar tripletul (Ω, \mathcal{K}, P) se numește câmp de probabilitate.

2.2 Variabile aleatoare

Fie câmpul de probabilitate (Ω, \mathcal{K}, P) și σ -algebrele mulțimilor Borel $\mathcal{B}(\mathbb{R})$, $\mathcal{B}(\mathbb{R}^n)$, respectiv pe \mathbb{R} și \mathbb{R}^n .

Definiția 2.2.1. Numim variabilă aleatoare, aplicația $X: \Omega \longrightarrow \mathbb{R}$, dacă este \mathcal{K} -măsurabilă, adică

$$(2.2.1) \quad \forall B \in \mathcal{B}(\mathbb{R}) \implies X^{-1}(B) = \{\omega \in \Omega \mid X(\omega) \in B\} \in \mathcal{K}.$$

Observația 2.2.2. Având în vedere că $\mathcal{B}(\mathbb{R})$ poate fi generat de familia de intervale $\{(-\infty, x]\}_{x \in \mathbb{R}}$, condiția (2.2.1) de \mathcal{K} -măsurabilitate poate fi înlocuită prin

$$(2.2.2) \quad X^{-1}((-\infty, x]) = (X \leq x) = \{\omega \in \Omega \mid X(\omega) \leq x\} \in \mathcal{K}, \quad \forall x \in \mathbb{R}.$$

Observația 2.2.3. Cardinalul mulțimii $X(\Omega)$ a valorilor unei variabile aleatoare X , conduce la o clasificarea a mulțimii variabilelor aleatoare:

- *variabile aleatoare de tip discret*, dacă $|X(\Omega)| \leq \aleph_0$, adică X are o mulțime cel mult numărabilă de valori;
- *variabile aleatoare simple*, dacă $|X(\Omega)| < \aleph_0$, adică X are o mulțime finită de valori;
- *variabile aleatoare de tip continuu*, dacă $|X(\Omega)| = \aleph$, adică X are o mulțime valori de puterea continuumului.

Definiția 2.2.4. Numim vector aleator n -(dimensional), aplicația $\mathbf{X}: \Omega \rightarrow \mathbb{R}^n$, dacă este \mathcal{K} -măsurabilă, adică

$$\forall \mathbf{B} \in \mathcal{B}(\mathbb{R}^n) \implies \mathbf{X}^{-1}(\mathbf{B}) = \{\omega \in \Omega \mid \mathbf{X}(\omega) \in \mathbf{B}\} \in \mathcal{K}$$

sau având în vedere Observația 2.2.2

$$\begin{aligned} (\mathbf{X} \leq \mathbf{x}) &= (X_1 \leq x_1, \dots, X_n \leq x_n) \\ &= \{\omega \in \Omega \mid X_1(\omega) \leq x_1, \dots, X_n(\omega) \leq x_n\} \in \mathcal{K}, \quad \forall \mathbf{x} \in \mathbb{R}^n, \end{aligned}$$

unde X_i și x_i , $i = \overline{1, n}$, reprezintă respectiv componentele vectorului aleator \mathbf{X} și ale vectorului \mathbf{x} .

2.3 Funcție de repartiție

Fie câmpul de probabilitate (Ω, \mathcal{K}, P) și variabila aleatoare $X: \Omega \rightarrow \mathbb{R}$.

Definiția 2.3.1. Numim funcție de repartiție atașată variabilei aleatoare X , aplicația $F: \mathbb{R} \rightarrow \mathbb{R}$, definită prin

$$(2.3.1) \quad F(x) = P(X \leq x), \quad \forall x \in \mathbb{R}.$$

Observația 2.3.2. Marea majoritate a tratatelor de teoria probabilităților definesc funcția de repartiție prin formula

$$F(x) = P(X < x), \quad \forall x \in \mathbb{R},$$

dar având în vedere că sistemul Matlab folosește formula (2.3.1), în definirea funcției de repartiție, vom considera această abordare în cele ce urmează.

Definiția 2.3.3. Numim funcție de repartiție atașată vectorului aleator n -dimensional \mathbf{X} , aplicația $F: \mathbb{R}^n \rightarrow \mathbb{R}$, definită prin

$$F(\mathbf{x}) = F(x_1, \dots, x_n) = P(X_1 \leq x_1, \dots, X_n \leq x_n), \quad \forall \mathbf{x} \in \mathbb{R}^n.$$

Fie vectorul aleator \mathbf{X} cu funcția de repartiție F și fie notată respectiv prin F_i funcția de repartiție a variabilei aleatoare X_i , care este componenta i a vectorului aleator, pentru fiecare $i = \overline{1, n}$.

Definiția 2.3.4. Spunem că variabilele aleatoare X_i , $i = \overline{1, n}$, sunt independente dacă

$$F(\mathbf{x}) = F(x_1, \dots, x_n) = F_1(x_1) \dots F_n(x_n), \quad \forall \mathbf{x} \in \mathbb{R}^n.$$

2.4 Legi de probabilitate de tip discret

Dacă variabila aleatoare X este de tip discret, adică are un număr cel mult numărabil de valori, fie acestea $x_i \in \mathbb{R}$, $i \in I$, atunci funcția de repartiție F atașată este o funcție în scară și este dată prin:

$$(2.4.1) \quad F(x) = \sum_{\substack{i \in I \\ x_i \leq x}} p_i, \quad \forall x \in \mathbb{R},$$

unde $p_i = P(X = x_i)$.

Definiția 2.4.1. Numim distribuția sau repartiția variabilei aleatoare X de tip discret, tabloul

$$X \begin{pmatrix} x_i \\ p_i \end{pmatrix}_{i \in I},$$

unde $x_i \in \mathbb{R}$, $i \in I$, sunt valorile pe care le ia variabila aleatoare X , iar p_i este probabilitatea cu care variabila aleatoare X ia valoarea x_i , adică $p_i = P(X = x_i)$, pentru fiecare $i \in I$.

Deoarece evenimentele $(X = x_i)$, $i \in I$, formează un sistem complet de evenimente, avem că $\sum_{i \in I} p_i = 1$. Remarcăm de asemenea că p_i , reprezintă mărimea saltului funcției de repartiție în punctul de discontinuitate x_i , pentru fiecare $i \in I$.

În teoria probabilităților și în aplicațiile sale, se întâlnesc clase de variabile aleatoare de tip discret. Forma cea mai generală a unei variabile aleatoare aparținând unei clase se numește *lege de probabilitate de tip discret*.

2.4.1 Funcțiile Matlab pdf și cdf

Distribuția variabilei aleatoare X poate fi precizată prin ceea ce numim *funcție de probabilitate* (pdf – probability distribution function, în Matlab) definită prin:

$$f(x_i) = p_i, \quad i \in I,$$

sau prin *funcția de repartiție* (cdf – cumulative distribution function).

Dacă valorile p_i , $i \in I$, sunt calculate, atunci folosind instrucțiunile `plot`, `bar` și `stairs`, se pot reprezenta grafic funcția de probabilitate (pdf) și funcția de repartiție (cdf). Anume, dacă vectorul x conține valorile variabilei aleatoare X , iar p probabilitățile corespunzătoare, atunci instrucțiunile

```
plot(x,p,'s')
bar(x,p)
stairs(x,p)
```

vor reprezenta grafic respectiv funcția de probabilitate prin simbolul precizat prin s , funcția de probabilitate prin bare și funcția de repartiție (funcție în scară).

Remarcăm faptul că dacă X ia o infinitate numărabilă de valori, atunci trebuie să ne limităm la un număr finit de valori ale variabilei aleatoare, iar reprezentările grafice se vor face pe domeniul cuprins între valorile minimă și maximă ale acestora.

Programul 2.4.2. Să considerăm variabila aleatoare X care are distribuția

$$X \begin{pmatrix} -1 & 0 & 1 \\ \frac{1}{3} & \frac{1}{6} & \frac{1}{2} \end{pmatrix}$$

și vrem să reprezentăm grafic funcția de probabilitate (prin puncte și bare) și funcția de repartiție pe aceeași figură.

Am putea presupune că variabila aleatoare X se referă la aruncarea unui zar. Anume, dacă în urma aruncării zarului se obține un număr compus (4 sau 6), atunci se pierde o miză ($X = -1$), dacă se obține un număr prim (2, 3 sau 5) se câștigă o miză ($X = 1$), altfel nu se câștigă și nu se pierde nimic ($X = 0$).

Următorul program

```
x = [-1:1]; p = [1/3,1/6,1/2];
pc = [1/3,1/2,1];
subplot(1,3,1), plot(x,p,'o')
axis([-1.5 1.5 0 1])
title('Funcția de probabilitate')
subplot(1,3,2), bar(x,p)
axis([-1.5 1.5 0 1])
title('Funcția de probabilitate')
subplot(1,3,3), stairs(x,pc)
title('Funcția de repartiție')
```

produce reprezentările grafice din Figura 2.1.

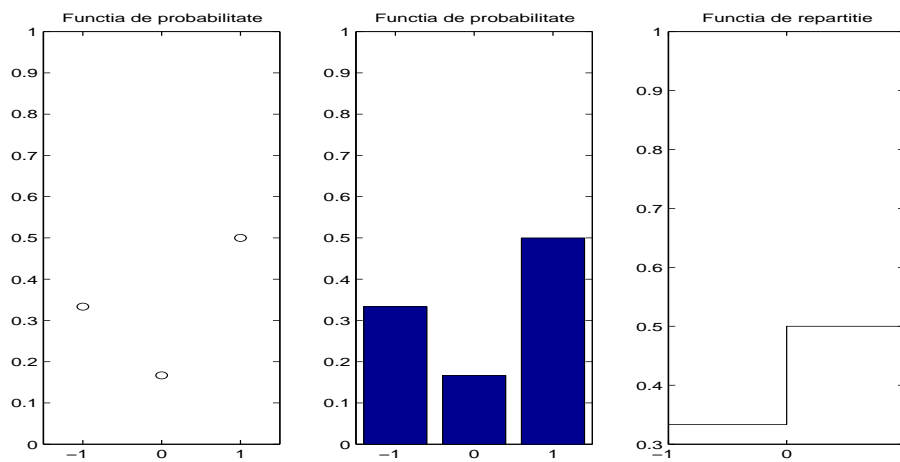


Figura 2.1: Funcția de frecvență și funcția de repartiție

Sistemul Matlab dispune de funcții (instrucțiuni) care calculează valorile funcției de probabilitate și ale funcției de repartiție pentru legile de probabilitate implementate prin *Statistics toolbox*.

Funcția pdf

Două moduri de apel pentru calculul valorilor funcției pdf avem:

```
p=pdf('legea',x,par1,par2,...)
p=numef(x,par1,par2,...)
```

unde *legea* este un șir de caractere predefinit pentru fiecare din legile de probabilitate disponibile în *Statistics toolbox*, *numef* este un șir de caractere din care ultimele trei sunt pdf, iar cele ce le preced sunt cele care dau numele predefinit al legii de probabilitate (ca și cele din parametrul *legea*).

În urma executării uneia din cele două instrucțiuni, se calculează matricea *p* a probabilităților legii precizată prin parametrii *legea*, respectiv *numef*, corespunzătoare valorilor date prin matricea *x* și având parametrii dați prin matricele *par1*, *par2*,... Aceste matrice trebuie să fie de aceeași dimensiuni, cu excepția că dacă unele sunt scalari, aceștia se extind la matricele constante de aceeași dimensiuni cu celelalte și care iau valorile scalarilor corespunzători.

Funcția cdf

Și pentru calculul valorilor funcției de repartiție funcția *cdf* avem două forme de apel

```
cp=cdf('legea',x,par1,par2,...)
cp=numef(x,par1,par2,...)
```

unde *legea* este un șir de caractere predefinit pentru fiecare din legile de probabilitate disponibile în *Statistics toolbox*, *numef* este un șir de caractere din care ultimele trei sunt *cdf*, iar cele ce le preced sunt cele care dau numele predefinit al legii de probabilitate (ca și cele din parametrul *legea*).

În urma executării uneia din cele două instrucțiuni, se calculează matricea *cp* a probabilităților cumulate (valorile funcției de repartiție) a legii precizată prin *legea* respectiv *numef*, corespunzătoare valorilor date prin matricea *x* și având parametrii dați prin matricele *par1*, *par2*,... Aceste matrice trebuie să fie de aceleași dimensiuni, cu excepția că dacă unele sunt scalari, aceștia se extind la matricele constante de aceleași dimensiuni cu celelalte și care iau valorile scalarilor corespunzători.

2.4.2 Legi de probabilitate de tip discret clasice

Vom prezenta în continuare legile de probabilitate de tip discret recunoscute de sistemul Matlab, prin pachetul de programe *Statistics toolbox*. În paranteză, pentru fiecare lege de probabilitate sunt trecute și denumirile acceptate de sistemul Matlab. Facem observația că pentru unele funcții sunt acceptate și alte denumiri.

Legea lui Bernoulli

Variabila aleatoare X urmează *legea lui Bernoulli*, pe care o notăm $Be(p)$, dacă are distribuția

$$(2.4.2) \quad X \begin{pmatrix} 1 & 0 \\ p & q \end{pmatrix}, \quad \text{unde } p \in (0, 1), q = 1 - p.$$

Legea lui Bernoulli modelează un fenomen aleator în desfășurarea căruia suntem interesați în apariția unui eveniment fixat A , pe care îl vom numi *succes*, respectiv a lui \bar{A} numit *insucces*. Probabilitatea obținerii succesului este p , iar valoarea corespunzătoare a variabilei aleatoare X este 1, respectiv $q = 1 - p$ este probabilitatea producerii insuccesului, cu valoarea 0 corespunzătoare pentru X .

Funcția de probabilitate este

$$f(x|p) = p^x q^{1-x}, \quad x = 0, 1.$$

Legea uniformă discretă (unid)

Variabila aleatoare X urmează *legea uniformă discretă*, notăm $\mathcal{U}(N)$, dacă are distribuția

$$(2.4.3) \quad X \left(\frac{k}{N} \right)_{k=\overline{1, N}}, \quad \text{unde } N \in \mathbb{N}.$$

Funcția de probabilitate este

$$f(x | N) = \frac{1}{N}, \quad x = \overline{1, N}.$$

Legea binomială (bino)

Spunem că variabila aleatoare X urmează *legea binomială*, notăm $\mathcal{B}(n, p)$, dacă are distribuția

$$(2.4.4) \quad X \left(\frac{k}{P(n, k)} \right)_{k=\overline{0, n}}, \quad \text{unde } P(n, k) = \binom{n}{k} p^k q^{n-k},$$

iar $p \in (0, 1)$ și $q = 1 - p$.

Variabila aleatoare X reprezintă numărul succeselor obținute în n repetări independente ale unui experiment.

Descoperirea legii binomială este atribuită lui James Bernoulli și care se află în cartea *Ars Conjectandi* (1713). Pascal a considerat cazul particular $p = \frac{1}{2}$.

Funcția de probabilitate este

$$f(x | n, p) = \binom{n}{x} p^x q^{n-x}, \quad x = \overline{0, n}.$$

Remarcăm faptul că variabila aleatoare X se obține ca suma a n variabile aleatoare independente, ce urmează fiecare legea lui Bernoulli $\mathcal{B}e(p)$. Prin urmare, putem spune și faptul că legea lui Bernoulli $\mathcal{B}e(p)$ se obține din legea binomială $\mathcal{B}(n, p)$, pentru $n = 1$.

Programul 2.4.3. Să scriem un program Matlab care să reprezinte grafic funcția de probabilitate (prin puncte și prin bare) și funcția de repartiție ale legii de probabilitate $\mathcal{B}(n, p)$.

Executarea programului

```

n = input('n=');
p = input('p='); x = 0:n;
f = pdf('bino',x,n,p);
subplot(1,3,1), plot(x,f,'o')
axis([-0.5 n+0.5 0 max(p)+0.02])
title('Functia de probabilitate')
subplot(1,3,2), bar(x,f)
axis([-0.5 n+0.5 0 max(p)+0.02])
title('Functia de probabilitate')
f = cdf('bino',x,n,p);
subplot(1,3,3), stairs(x,f)
title('Functia de repartitie')
axis([0 n 0 1])

```

pentru $n=7$ și $p=0.3$, realizează graficele din Figura 2.2.

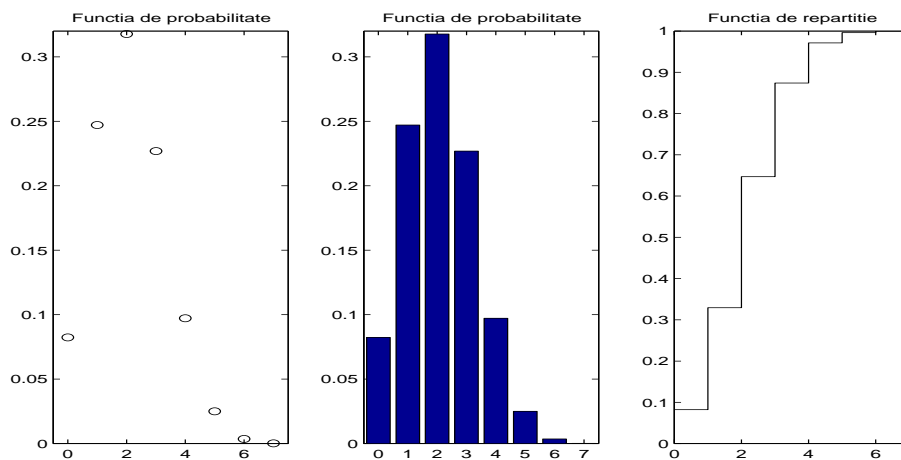


Figura 2.2: Legea binomială $\mathcal{B}(7, 0.3)$

Legea hipergeometrică (hyge)

Spunem că variabila aleatoare X urmează *legea hipergeometrică* și vom nota prin $\mathcal{H}(n, M, K)$, dacă are distribuția

$$(2.4.5) \quad X \left(\begin{matrix} k \\ P(n, k) \end{matrix} \right)_{k=0, \overline{n}}, \text{ unde } P(n, k) = \frac{\binom{K}{k} \binom{M-K}{n-k}}{\binom{M}{n}}, \text{ iar } n \leq K \leq M.$$

Variabila aleatoare X reprezintă numărul succeselor obținute în n extrageri dintr-o populație de volum M , fără întoarcere, dacă numărul indivizilor cu proprietatea cercetată este K .

Funcția de probabilitate corespunzătoare este:

$$f(x | n, M, K) = \frac{\binom{K}{x} \binom{M-K}{n-x}}{\binom{M}{n}}, \quad x = \overline{0, n}.$$

Observația 2.4.4. Dacă se notează $p = \frac{K}{M}$ și $q = \frac{M-K}{M}$, adică probabilitățile ca la prima extragere să se obțină succes, respectiv insucces, iar $M \rightarrow \infty$, atunci

$$P(n, k) \rightarrow \binom{n}{k} p^k q^{n-k},$$

adică se obține distribuția binomială.

Programul 2.4.5. Pentru a ilustra afirmația precedentă, să scriem un program Matlab care să reprezinte grafic prin bare funcțiile de probabilitate pentru $\mathcal{H}(n, M, K)$ și $\mathcal{B}(n, p)$, cu $p = \frac{K}{M}$.

Executarea programului

```
clf;
M = input('M: ');
K = input('K(K<=M): ');
n = input('n(n<=K): ');
x = 0:n; p=K/M;
fh = pdf('hyge', x, M, K, n);
fb = pdf('bino', x, n, p);
bar(x', [fh', fb'])
colormap winter
```

pentru $M=100$, $K=40$ și $n=10$, realizează graficul din Figura 2.3.

Legea lui Poisson (**poiss**)

Variabila aleatoare X urmează *legea lui Poisson*, vom nota $\mathcal{Po}(\lambda)$, dacă are distribuția

$$(2.4.6) \quad X \left(\begin{matrix} k \\ p_k(\lambda) \end{matrix} \right)_{k=0,1,2,\dots}, \text{ unde } p_k(\lambda) = \frac{\lambda^k}{k!} e^{-\lambda}, \text{ iar } \lambda > 0.$$

Funcția de probabilitate corespunzătoare este

$$f(x | \lambda) = \frac{\lambda^x}{x!} e^{-\lambda}, \quad x = 0, 1, \dots$$

Variabila aleatoare X , care urmează legea lui Poisson, numără de câte ori apare un anumit eveniment într-un interval de tip, pe o distanță, pe o suprafață etc. Poisson

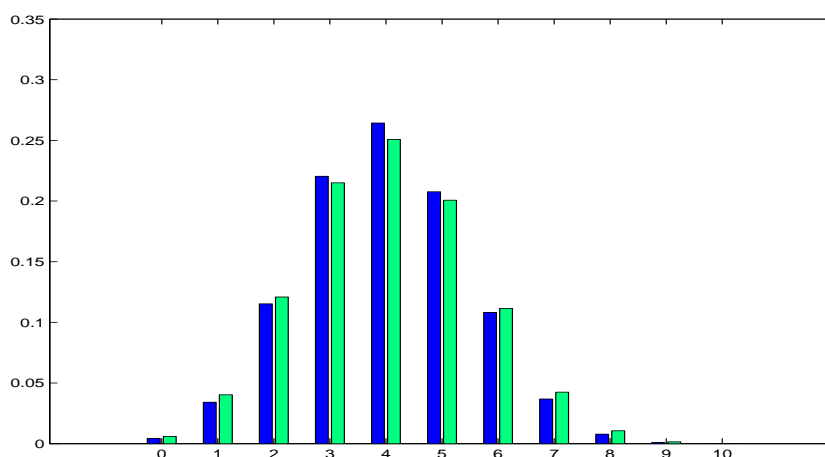


Figura 2.3: Legea hipergeometrică $\mathcal{H}(10, 100, 40)$ și legea binomială $\mathcal{B}(10, 0.4)$

(1837) a arătat că legea, care-i poartă numele, este un caz limită a legii binomiale, anume când $np \rightarrow \lambda$, pentru $n \rightarrow \infty$. Altfel spus, când n este mare, p trebuie să fie mic, adică probabilitatea producerii evenimentului considerat este mică, de aici denumirea de *lege a evenimentelor rare*.

Mai remarcăm un rezultat cunoscut, care stabilește legătura dintre legea lui Poisson și *legea exponențială*, care va fi prezentată mai încolo. Anume, pe când legea lui Poisson numără evenimentele ce apar într-un interval de timp, legea exponențială dă lungimea intervalului dintre două apariții consecutive ale evenimentelor.

Programul 2.4.6. Vom scrie un program Matlab, care să reprezinte grafic prin bare funcțiile de probabilitate pentru $\mathcal{B}(n, p)$ și $\mathcal{Po}(\lambda)$ cu $\lambda = np$.

Având în vedere comportarea legii $\mathcal{B}(n, p)$, când $n \rightarrow \infty$, graficele funcțiilor de probabilitate pentru cele două legi de probabilitate trebuie să se apropie.

Mai remarcăm faptul că legea $\mathcal{Po}(\lambda)$ are o infinitate (numărabilă) de valori, prin urmare în reprezentarea grafică pentru această lege de probabilitate ne vom limita la un număr finit de valori. Anume, vom considera numai valorile întregi din intervalul $[\lambda - 3\sqrt{\lambda}, \lambda + 3\sqrt{\lambda}]$, alegere care va fi justificată mai târziu.

Programul

```
clf;
n = input('n='); p = input('p=');
lambda = n*p;
vi = fix(lambda-3*sqrt(lambda));
vf = fix(lambda+3*sqrt(lambda));
x = vi:vf;
```

```
fb = binopdf(x,n,p);
fl = poisspdf(x,lambda);
colormap autumn
bar(x',[fb',fl'])
```

pentru $n=100$ și $p=0.05$, realizează graficul din Figura 2.4.

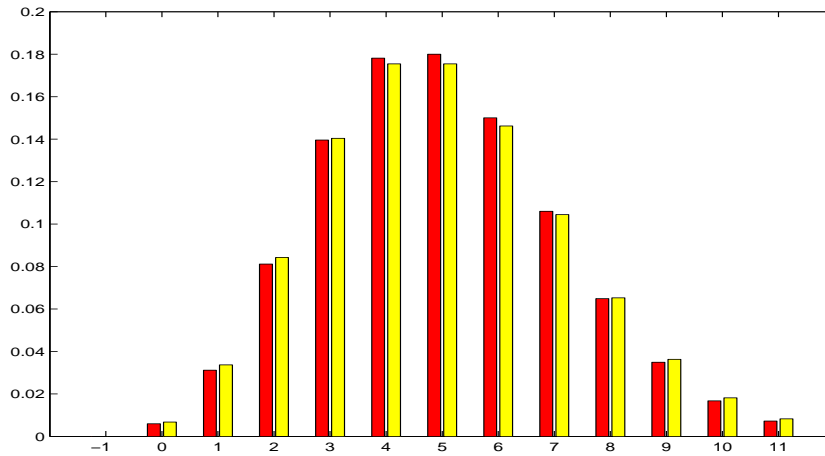


Figura 2.4: Legea binomială $\mathcal{B}(100, 0.05)$ și legea Poisson $\mathcal{Po}(5)$

Legea binomială negativă (nbin)

Spunem că variabila aleatoare X urmează *legea binomială negativă* sau *legea lui Pascal*, notată $\mathcal{BN}(r, p)$, dacă are distribuția

$$X \left(\begin{matrix} k \\ P(r, k) \end{matrix} \right)_{k=0,1,2,\dots}, \text{ cu } P(r, k) = \binom{r+k-1}{k} p^r q^k,$$

iar $p \in (0, 1)$ și $q = 1 - p$. Variabila aleatoare X ce urmează legea $\mathcal{BN}(r, p)$ reprezintă numărul insucceselor până la apariția succesului de rang r , dacă repetările sunt considerate independente.

Funcția de probabilitate corespunzătoare legii $\mathcal{BN}(r, p)$ este

$$f(x|r, p) = \binom{r+x-1}{x} p^r q^x, \quad x = 0, 1, \dots$$

Observația 2.4.7. Unele tratate consideră că o variabilă aleatoare X ce urmează legea lui Pascal reprezintă rangul repetării la care se produce succesul cu numărul r .

Dacă se are în vedere această definiție se obține distribuția

$$(2.4.7) \quad X \left(\begin{matrix} k \\ P(r, k) \end{matrix} \right)_{k=1,2,\dots}, \text{ cu } P(r, k) = \binom{k-1}{r-1} p^r q^{k-r}.$$

Legea geometrică (geo)

Legea geometrică, notată $Ge(p)$, se obține ca și caz particular, $r = 1$, al legii $BN(r, p)$. Prin urmare variabila aleatoare X urmează legea $Ge(p)$, dacă are distribuția

$$X \left(\begin{matrix} k \\ pq^k \end{matrix} \right)_{k=0,1,2,\dots}, \quad p \in (0, 1), \quad p + q = 1,$$

și reprezintă numărul insucceselor până la apariția primului succes.

Funcția de probabilitate pentru legea $Ge(p)$ este

$$f(x|p) = pq^x, \quad x = 0, 1, \dots$$

De remarcat faptul că dacă se adună r variabile aleatoare independente, fiecare urmând legea $Ge(p)$, atunci se obține o variabilă aleatoare ce urmează legea $BN(r, p)$.

Observația 2.4.8. Corespunzător observației de la legea lui Pascal, avem că unele tratate consideră că o variabilă aleatoare X ce urmează legea geometrică reprezintă rangul repetării la care se produce primul succes. Dacă se are în vedere această definiție se obține distribuția

$$(2.4.8) \quad X \left(\begin{matrix} k \\ pq^{k-1} \end{matrix} \right)_{k=1,2,\dots}$$

2.5 Legi de probabilitate continue

Definiția 2.5.1. Fie variabila aleatoare X având funcția de repartiție F . Vom spune că X este variabilă aleatoare (absolut) continuă, dacă funcția de repartiție F este absolut continuă, sau echivalent se poate reprezenta sub forma

$$F(x) = \int_{-\infty}^x f(t) dt, \quad \text{pentru orice } x \in \mathbb{R},$$

funcția $f: \mathbb{R} \rightarrow \mathbb{R}$, numindu-se densitatea de probabilitate a variabilei aleatoare X .

Definiția 2.5.2. Fie vectorul aleator \mathbf{X} având funcția de repartiție F . Spunem că \mathbf{X} este vector aleator de tip continuu, dacă funcția de repartiție F se poate reprezenta sub forma

$$F(\mathbf{x}) = F(x_1, \dots, x_n) = \int_{-\infty}^{x_1} \dots \int_{-\infty}^{x_n} f(t_1, \dots, t_n) dt_1 \dots dt_n, \quad \forall \mathbf{x} \in \mathbb{R}^n,$$

funcția $f: \mathbb{R}^n \rightarrow \mathbb{R}$, numindu-se densitatea de probabilitate a vectorului aleator \mathbf{X} .

Ca la variabilele aleatoare de tip discret și la cele de tip continuu există clase de variabile aleatoare. Forma cea mai generală a densității de probabilitate a unei variabile aleatoare din clasa respectivă, ne dă o *lege de probabilitate de tip continuu*.

Funcțiile Matlab `pdf` și `cdf`

Distribuția variabilei aleatoare X de tip continuu poate fi precizată prin *funcție densitate de probabilitate* (`pdf` – probability density function) sau prin *funcția de repartiție* (`cdf` – cumulative distribution function).

Reprezentările grafice ale densității de probabilitate și funcției de repartiție în Matlab se fac folosind instrucțiunea de tip `plot`.

Valorile densității de probabilitate f și ale funcției de repartiție F , pentru legile de probabilitate implementate prin *Statistics toolbox*, ca și în cazul discret, se obțin cu ajutorul funcțiilor Matlab `pdf` și respectiv `cdf`. Trebuie însă să ținem seama că f este continuă (eventual cu o mulțime cel mult numărabilă de puncte de discontinuitate), iar F nu mai este o funcție în scară. Prin urmare, pentru asigurarea netezimii graficelor se impune să fie considerat un număr suficient de puncte ale graficelor acestor funcții.

2.5.1 Legi de probabilitate continue clasice

Prezentăm și în acest caz legile de probabilitate de tip continue recunoscute de sistemul Matlab, prin pachetul de programe *Statistics toolbox*.

Legea uniformă (`unif`)

Spunem că variabila aleatoare X urmează *legea uniformă* pe intervalul $[a, b]$, și vom nota $\mathcal{U}(a, b)$, dacă are densitatea de probabilitate

$$f(x | a, b) = \begin{cases} \frac{1}{b-a}, & \text{dacă } x \in [a, b], \\ 0, & \text{dacă } x \notin [a, b]. \end{cases}$$

Funcția de repartiție pentru variabila aleatoare X , ce urmează legea uniformă pe intervalul $[a, b]$, este

$$(2.5.1) \quad F(x|a, b) = \int_{-\infty}^x f(t|a, b) dt = \begin{cases} 0, & \text{dacă } x < a, \\ \frac{x-a}{b-a}, & \text{dacă } a \leq x \leq b, \\ 1, & \text{dacă } x > b. \end{cases}$$

Denumirea de lege uniformă este legată de faptul că dacă se consideră subintervale ale intervalului $[a, b]$ de lungimi egale, să zicem că au lungimea ℓ , atunci probabilitatea ca variabila aleatoare X să ia valori într-un astfel de interval este $\ell/(b-a)$. Adică, această probabilitate nu depinde de poziția subintervalului, ci numai de lungimea lui. Intuitiv, am putea spune că valorile întâmplătoare (aleatoare) ale funcției X sunt “uniform” răspândite pe intervalul $[a, b]$. Acest lucru nu se întâmplă la alte legi de probabilitate, când aceste valori sunt “grupate” în jurul uneia sau mai multor valori pe care le ia variabila aleatoare X .

Programul 2.5.3. Să reprezentăm grafic densitatea de probabilitate f și funcția de repartiție F pentru legea $\mathcal{U}(a, b)$.

Programul Matlab care urmează

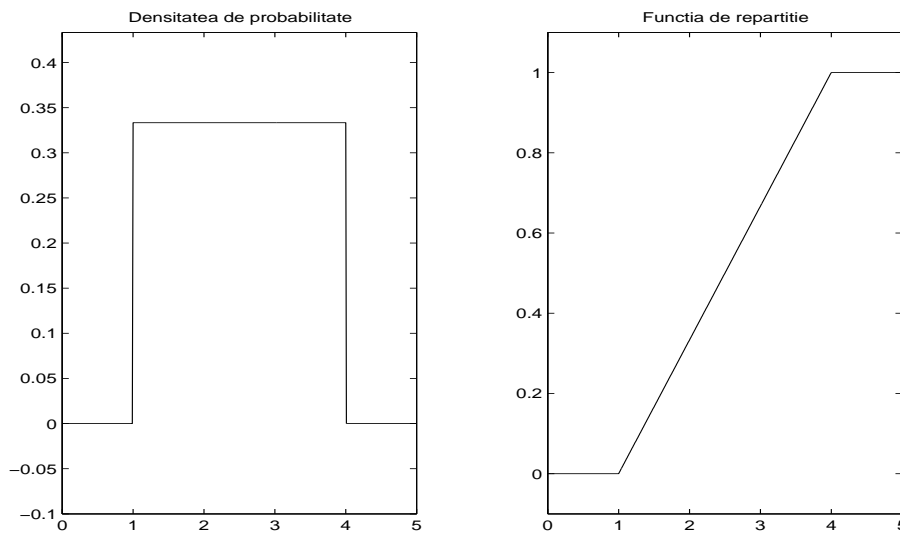
```
clf;
a = input('a='); b = input('b=');
x = a-1:0.01:b+1;
f = pdf('unif',x,a,b);
F = cdf('unif',x,a,b);
subplot(1,2,1), plot(x,f,'k-')
axis([a-1 b+1 -0.1 1/(b-a)+0.1])
title('Densitatea de probabilitate')
subplot(1,2,2), plot(x,F,'k-')
axis([a-1 b+1 -0.1 1.1])
title('Funcția de repartiție')
```

pentru $a=1$ și $b=4$, realizează graficele din Figura 2.5.

Putem remarca faptul că funcția de repartiție este continuă. Astfel valorile argumentului $x = a$ și $x = b$ în formula (2.5.1) pot fi atașate respectiv la cazurile $F(x) = 0$ și $F(x) = 1$. Prin urmare intervalul pe care densitatea de probabilitate este diferită de zero poate fi interval închis la ambele capete, închis la un capăt și deschis la celălalt, respectiv deschis la ambele capete. În toate aceste cazuri vom considera că avem legea uniformă de probabilitate $\mathcal{U}(a, b)$.

Proprietatea 2.5.4. Dacă variabila aleatoare U urmează legea uniformă pe intervalul $[0, 1]$, adică $\mathcal{U}(0, 1)$, numită și lege uniformă standard, și dacă $a < b$, atunci variabila aleatoare

$$X = (b-a)U + a,$$

Figura 2.5: Legea uniformă pe $[1, 4]$

urmează legea uniformă $\mathcal{U}(a, b)$.

Legea normală (norm)

Spunem că variabila aleatoare X urmează *legea normală* (*legea Gauss–Laplace*) de parametri $\mu \in \mathbb{R}$ și $\sigma > 0$, notăm aceasta prin $\mathcal{N}(\mu, \sigma)$, dacă are densitatea de probabilitate

$$f(x | \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \quad \text{pentru orice } x \in \mathbb{R}.$$

Funcția de repartiție pentru *legea normală standard* $\mathcal{N}(0, 1)$ este

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{t^2}{2}} dt, \quad x \in \mathbb{R},$$

și se numește *funcția lui Laplace*.

În Matlab (și nu numai) există o funcție `erf` (*funcția eroare*):

$$\operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt,$$

care verifică relația

$$\operatorname{erf}(x) = 2\Phi(x\sqrt{2}) - 1.$$

Remarcăm că funcția

$$\tilde{\Phi}(x) = \frac{1}{\sqrt{2\pi}} \int_0^x e^{-\frac{t^2}{2}} dt, \quad x \in \mathbb{R},$$

se numește de asemenea *funcția lui Laplace*. Între cele două funcții Laplace există relația

$$\Phi(x) = \frac{1}{2} + \tilde{\Phi}(x).$$

Folosind cele două funcții ale lui Laplace avem:

$$\begin{aligned} F(x|\mu, \sigma) &= \int_{-\infty}^x f(t|\mu, \sigma) dt = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{(t-\mu)^2}{2\sigma^2}} dt \\ &= \Phi\left(\frac{x-\mu}{\sigma}\right) = \frac{1}{2} + \tilde{\Phi}\left(\frac{x-\mu}{\sigma}\right). \end{aligned}$$

Programul 2.5.5. Vom reprezenta grafic densitatea de probabilitate f pentru legea $\mathcal{N}(\mu, \sigma)$ și funcția de repartiție F corespunzătoare.

Execuția programului

```
clf;
m = input('mu='); s = input('sigma=');
x = m-3*s:0.01:m+3*s;
f = pdf('norm',x,m,s);
F = cdf('norm',x,m,s);
subplot(1,2,1), plot(x,f,'k-')
title('Densitatea de probabilitate')
subplot(1,2,2), plot(x,F,'k-')
title('Funcția de repartiție')
```

pentru $\mu=0$ și $\sigma=1$, produce graficele din Figura 2.6.

Observația 2.5.6. Dacă se consideră variabilele aleatoare X_k , $k = \overline{1, n}$, independente, fiecare urmând respectiv legile normale $\mathcal{N}(\mu_k, \sigma_k)$, atunci combinația liniară

$$Z = \sum_{k=1}^n a_k X_k, \quad a_k \in \mathbb{R}, \quad \sum_{k=1}^n a_k^2 > 0,$$

este o variabilă aleatoare ce urmează legea normală $\mathcal{N}(\mu, \sigma)$, unde

$$\mu = \sum_{k=1}^n a_k \mu_k, \quad \sigma^2 = \sum_{k=1}^n a_k^2 \sigma_k^2.$$

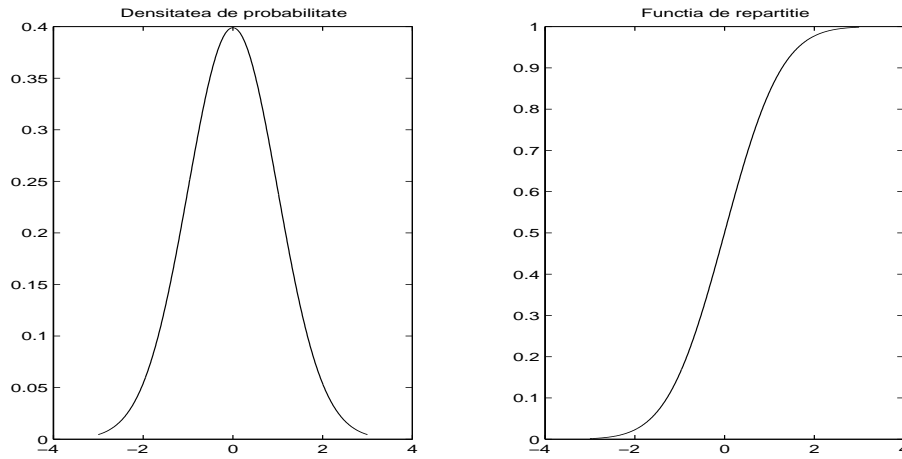


Figura 2.6: Legea normală standard

Legea lognormală (logn)

Variabila aleatoare X urmează *legea lognormală de probabilitate* și vom nota aceasta prin $\mathcal{LN}(\mu, \sigma)$, dacă are densitatea de probabilitate

$$f(x | \mu, \sigma) = \frac{1}{x\sigma\sqrt{2\pi}} e^{-\frac{(\ln x - \mu)^2}{2\sigma^2}}, \quad x > 0; \quad \mu \in \mathbb{R}, \quad \sigma > 0.$$

Dacă X urmează legea lognormală de parametri μ și σ , atunci variabila aleatoare $Y = \ln X$ urmează legea normală ($X = e^Y$).

Legea gamma (gam)

Variabila aleatoare X urmează *legea gamma*, notăm aceasta prin $\mathcal{Ga}(a, b)$, dacă are densitatea de probabilitate

$$f(x | a, b) = \frac{1}{b^a \Gamma(a)} x^{a-1} e^{-\frac{x}{b}}, \quad x > 0; \quad a, b > 0,$$

unde $\Gamma(a)$ este *funcția lui Euler de speța a doua*, adică

$$\Gamma(a) = \int_0^\infty x^{a-1} e^{-x} dx, \quad a > 0.$$

Când $a \rightarrow \infty$, legea gamma tinde la legea normală. Adică, dacă X urmează legea $\mathcal{Ga}(a, b)$, atunci pentru $a \rightarrow \infty$, variabila aleatoare poate fi considerată ca urmând legea $\mathcal{N}(\mu, \sigma)$, cu $\mu = ab$ și $\sigma = b\sqrt{a}$.

Programul 2.5.7. Pentru a ilustra această ultimă afirmație, să reprezentăm grafic, pe aceeași figură, densitățile de probabilitate pentru cele două legi, parametrii acestora satisfăcând relațiile precizate mai sus.

Executând programul

```
clf;
a = input('a='); b = input('b=');
m = a*b; s = b*sqrt(a);
x = m-3*s:0.01:m+3*s;
fn = pdf('norm',x,m,s);
fg = pdf('gam',x,a,b);
plot(x,fn,'k-.',x,fg,'k-')
legend('Legea normala','Legea gamma')
```

pentru $a=100$ și $b=10$, acesta produce graficele din Figura 2.7.

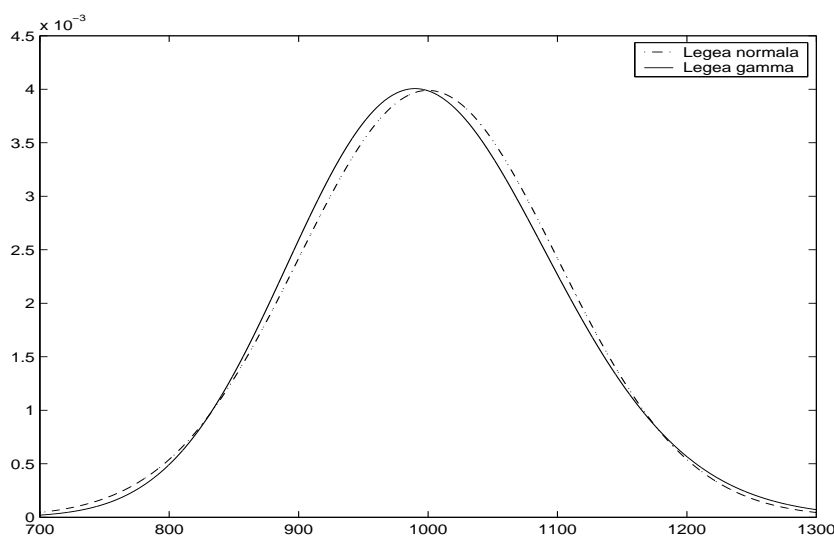


Figura 2.7: Legile $\mathcal{G}a(100, 10)$ și $\mathcal{N}(1000, 100)$

Legea exponențială (**exp**)

Variabila aleatoare X urmează *legea exponențială*, notăm $\mathcal{Exp}(\mu)$, dacă are densitatea de probabilitate

$$f(x|\mu) = \frac{1}{\mu} e^{-\frac{x}{\mu}}, \quad x > 0; \quad \mu > 0.$$

Observăm că se obține ca și caz particular al legii $\mathcal{G}a(a, b)$, când se ia $a = 1$ și $b = \mu$.

Legea exponențială se aplică la modelarea evenimentelor ce apar aleator în timp. De asemenea are aplicații în studiile privind durata de viață.

Legea beta (beta)

Variabila aleatoare X urmează *legea beta*, notată $Beta(a, b)$, dacă are densitatea de probabilitate:

$$f(x|a, b) = \frac{1}{B(a, b)} x^{a-1} (1-x)^{b-1}, \quad x \in (0, 1); \quad a, b > 0,$$

unde

$$B(a, b) = \int_0^1 x^{a-1} (1-x)^{b-1} dx,$$

este funcția lui Euler de speța întâi.

Este cunoscut faptul că dacă două variabile aleatoare independente U și V urmează legile $\mathcal{G}_a(a, 1)$ și $\mathcal{G}_a(b, 1)$, atunci variabila aleatoare $X = U/(U + V)$ urmează legea $Beta(a, b)$.

Dacă se consideră legea $Beta(a, b)$, cu $a = b = \frac{1}{2}$, se obține *legea de probabilitate arcsin*, care are densitatea de probabilitate

$$f(x) = \frac{1}{\pi \sqrt{x(1-x)}}, \quad x \in (0, 1).$$

Legea Weibull (weib)

Variabila aleatoare X urmează *legea lui Weibull*, notăm $\mathcal{W}(a, b)$, dacă are densitatea de probabilitate

$$f(x|a, b) = abx^{b-1}e^{-ax^b}, \quad x > 0; \quad a, b > 0.$$

Legea a fost descoperită de Weibull (1939) și modelează rezistența la rupere a materialelor. De asemenea se pare că modelează durata de viață mai bine decât legea exponențială. Se observă că legea exponențială este un caz particular, anume se obține când $a = \frac{1}{\mu}$ și $b = 1$.

Legea Rayleigh (ray1)

Variabila aleatoare X urmează *legea lui Rayleigh de probabilitate*, notată prin $\mathcal{R}(b)$, dacă are densitatea de probabilitate:

$$f(x|b) = \frac{x}{b^2} e^{-\frac{x^2}{2b^2}}, \quad x > 0; \quad b > 0.$$

Se observă că este caz particular al legii Weibull, $b = 2$, $a := 1/(2b^2)$.

Dacă viteza unei particule în direcțiile x și y sunt două variabile aleatoare independente, care urmează legi normale de probabilitate, cu mediile zero și dispersiile egale, atunci distanța parcursă de particulă pe unitate de timp urmează legea lui Rayleigh.

2.5.2 Legi de probabilitate continue statistice

Prezentăm legile de probabilitate de tip continuu, denumite astfel în sistemul Matlab, prin pachetul de programe *Statistics toolbox*, deoarece utilizarea lor este strâns legată de statistică.

Legea t (Student) (\mathfrak{t})

Variabila aleatoare X urmează *legea t (Student) de probabilitate*, notată $\mathcal{T}(n)$, dacă are densitatea de probabilitate

$$f(x|n) = \frac{\Gamma\left(\frac{n+1}{2}\right)}{\sqrt{n\pi}\Gamma\left(\frac{n}{2}\right)} \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}}, \quad x \in \mathbb{R},$$

$n \in \mathbb{N}$ (numărul gradelor de libertate).

Gosset (1908) a descoperit această lege de probabilitate. Nu i-a fost permis să publice rezultatul, drept urmare a publicat-o sub pseudonimul Student.

Pentru $n = 1$, se obține *legea de probabilitate a lui Cauchy*.

Dacă X_0, X_1, \dots, X_n sunt independente, fiecare urmând legea normală de parametri $\mu = 0$ și $\sigma = 1$, atunci

$$T = \frac{X_0}{\sqrt{\frac{1}{n} \sum_{k=1}^n X_k^2}},$$

urmează legea $\mathcal{T}(n)$.

Dacă Y urmează legea $\mathcal{T}(n)$, atunci

$$X = \frac{1}{2} + \frac{1}{2} \frac{Y}{\sqrt{n + Y^2}}$$

urmează legea $\text{Beta}\left(\frac{n}{2}, \frac{n}{2}\right)$.

Când $n \rightarrow \infty$, se ajunge la legea normală standard $\mathcal{N}(0, 1)$.

Programul 2.5.8. Pentru ilustrarea ultimei afirmații, vom scrie un program care să reprezinte pe aceeași figură graficele densităților de probabilitate pentru legile $\mathcal{T}(n)$ și $\mathcal{N}(0, 1)$.

Prin executarea programului Matlab

```

clf;
n = input('n='); x = -3:0.01:3;
fn = normpdf(x,0,1); fs = tpdf(x,n);
plot(x,fn,'k-.',x,fs,'k-')
legend('Legea normala', 'Legea Student')

```

pentru $n=10$, acesta produce graficele din Figura 2.8.

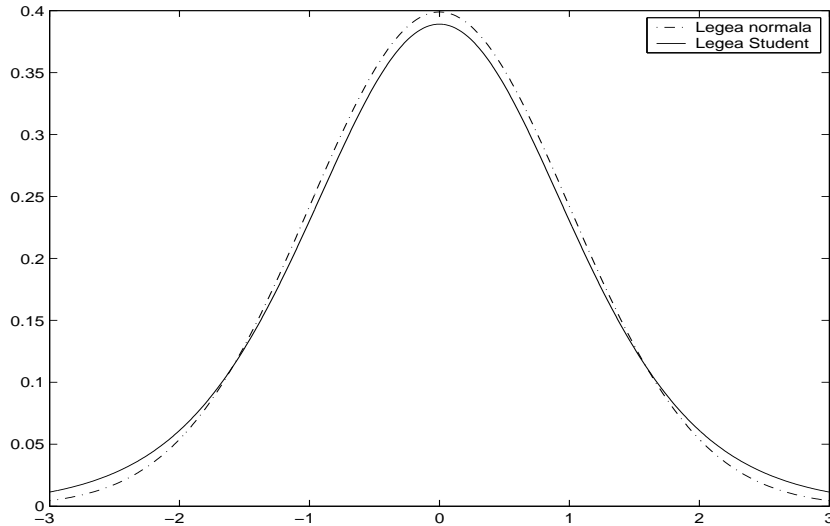


Figura 2.8: Legea $\mathcal{T}(10)$ și legea $\mathcal{N}(0, 1)$

Legea t (Student) necentrată (nct)

Legea t (Student) necentrată, notată $\mathcal{T}nc(n, \delta)$, generalizează legea $\mathcal{T}(n)$ și are densitatea de probabilitate

$$f(x | n, \delta) = \frac{\Gamma\left(\frac{n+1}{2}\right)}{\sqrt{n\pi}\Gamma\left(\frac{n}{2}\right)} e^{-\frac{\delta^2}{2}} \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}} h(x | n, \delta), \quad x \in \mathbb{R},$$

unde

$$h(x | n\delta) = \sum_{k=0}^{\infty} \left(\frac{\delta t \sqrt{2}}{n}\right)^k \frac{\Gamma\left(\frac{n+k+1}{2}\right)}{k! \Gamma\left(\frac{n+1}{2}\right)} \left(1 + \frac{x^2}{n}\right)^{-\frac{k}{2}}.$$

Programul 2.5.9. Programul Matlab ce urmează reprezintă grafic funcțiile de repartiție pentru legile $\mathcal{T}(n)$ și $\mathcal{T}nc(n, \delta)$:

```

clf;
n = input('n='); d = input('delta=');
x = -5:0.01:5;
Fc = tcdf(x,n); Fnc = nctcdf(x,n,d);
plot(x,Fc,'k-.',x,Fnc,'k-')
legend('Legea Student',...
      'Legea Student necentrata',2)

```

Pentru $n=10$ și $d=1$, se obțin graficele din Figura 2.9.

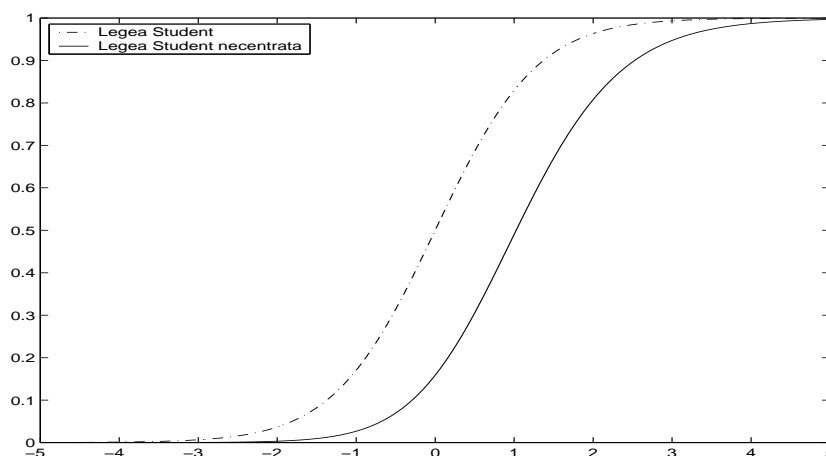


Figura 2.9: Legea $\mathcal{T}_{nc}(10, 1)$ și legea $\mathcal{T}(10)$

Legea χ^2 (chi2)

Variabila aleatoare X urmează *legea χ^2* sau *legea Helmer–Pearson*, vom nota prin $\chi^2(n)$, dacă are densitatea de probabilitate:

$$f(x|n) = \frac{1}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})} x^{\frac{n}{2}-1} e^{-\frac{x}{2}}, \quad x > 0,$$

$n \in \mathbb{N}$ (numărul gradelor de libertate).

Legea $\chi^2(n)$ este un caz particular al legii $\mathcal{Ga}(a, b)$, când $a = \frac{n}{2}$, iar $b = 2$.

Dacă X_1, X_2, \dots, X_n sunt independente, fiecare urmând legea normală de parametri $\mu = 0$ și $\sigma = 1$, atunci variabila aleatoare

$$\chi^2 = \sum_{k=1}^n X_k^2,$$

urmează legea $\chi^2(n)$, iar dacă variabilele aleatoare independente X și Y urmează respectiv legile $\mathcal{N}(0, 1)$ și $\chi^2(n)$, atunci variabila aleatoare

$$T = \frac{X}{\sqrt{\frac{Y}{n}}}$$

urmează legea $\mathcal{T}(n)$.

Dacă variabila aleatoare X urmează legea $\chi^2(n)$, atunci, pentru $n \rightarrow \infty$, variabila aleatoare urmează legea $\mathcal{N}(\mu, \sigma)$, cu $\mu = n$ și $\sigma = \sqrt{2n}$.

Programul 2.5.10. Pentru ilustrarea afirmației, scriem un program care să reprezinte pe aceeași figură graficele densităților de probabilitate pentru legile $\chi^2(n)$ și $\mathcal{N}(n, \sqrt{2n})$.

Prin executarea programului Matlab

```
clf;
n = input('n='); s = sqrt(2*n);
x = n-3*s:0.01:n+3*s;
fn = normpdf(x,n,s); fchi = chi2pdf(x,n);
plot(x,fn,'k-.',x,fchi,'k-')
legend('Legea normala', 'Legea chi2',2)
```

pentru $n=10$, se obțin graficele din Figura 2.10.

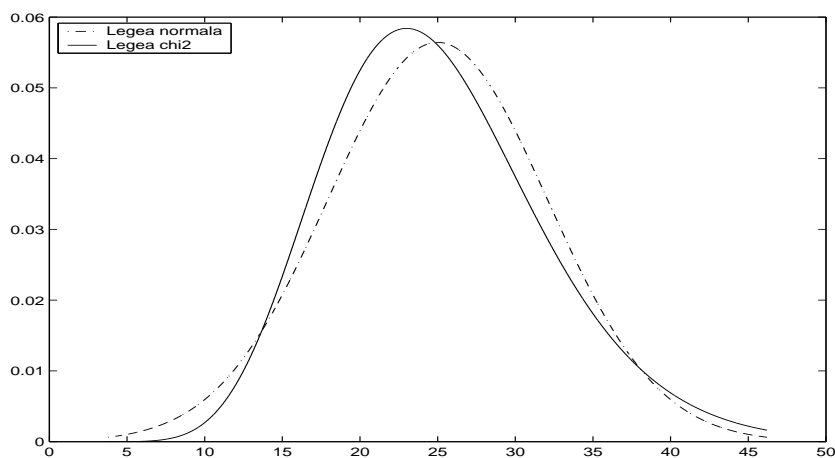


Figura 2.10: Legea $\chi^2(10)$ și legea $\mathcal{N}(10, 2\sqrt{5})$

Legea χ^2 necentrată (ncx2)

Să notăm prin χ_n^2 o variabilă aleatoare ce urmează legea $\chi^2(n)$. Spunem că variabila aleatoare X urmează *legea χ^2 necentrată*, dacă are funcția de repartiție:

$$F(x|n, \delta) = \sum_{k=0}^{\infty} \frac{\left(\frac{\delta}{2}\right)^k e^{-\frac{\delta}{2}}}{k!} P(\chi_{n+2k}^2 \leq x).$$

Programul 2.5.11. Programul Matlab ce urmează reprezintă grafic densitățile de probabilitate pentru legea $\chi^2(n)$ și legea $\chi^2(n, \delta)$ necentrată corespunzătoare:

```
clf;
n = input('n='); d = input('delta=');
x = 0:0.01:15;
c = chi2pdf(x,n); fnc = ncx2pdf(x,n,d);
plot(x,c,'k-.',x,fnc,'k-')
legend('Legea chi2', 'Legea chi2 necentrata')
```

Pentru $n=4$ și $d=2$, se obțin graficele din Figura 2.11.

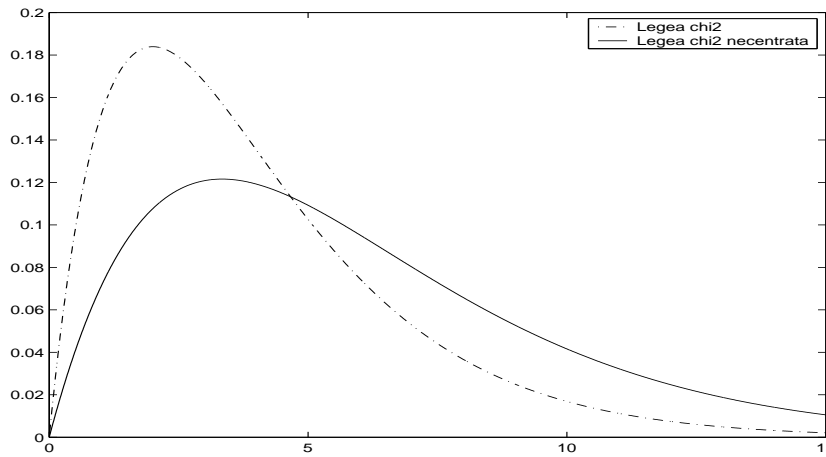


Figura 2.11: Legea $\chi^2(4, 2)$ și legea $\chi^2(4)$

Legea F (Fisher–Snedecor) (f)

Variabila aleatoare X urmează legea F (Fisher–Snedecor), notată $\mathcal{F}(m, n)$, dacă are densitatea de probabilitate:

$$f(x|m, n) = \frac{\Gamma\left(\frac{m+n}{2}\right)}{\Gamma\left(\frac{m}{2}\right)\Gamma\left(\frac{n}{2}\right)} \left(\frac{m}{n}\right)^{\frac{m}{2}} x^{\frac{m}{2}-1} \left(1 + \frac{m}{n}x\right)^{-\frac{m+n}{2}}, \quad x > 0;$$

$m, n \in \mathbb{N}$ (numărul gradelor de libertate)

Dacă variabilele aleatoare independente X și Y urmează legile $\chi^2(m)$ și $\chi^2(n)$, atunci

$$F = \frac{X}{m} \bigg/ \frac{Y}{n}$$

urmează legea $\mathcal{F}(m, n)$.

Dacă variabila aleatoare T urmează legea $\mathcal{T}(n)$, atunci T^2 va urma legea $\mathcal{F}(1, n)$.

Dacă variabila aleatoare X urmează legea $\mathcal{F}(m, n)$, atunci transformata lui Fisher, adică $Z = \frac{\log X}{2}$, urmează legea $\mathcal{N}(\mu, \sigma)$, când $m, n \rightarrow \infty$, unde

$$(2.5.2) \quad \mu = \frac{1}{2} \left(\frac{1}{n} - \frac{1}{m} \right), \quad \sigma = \sqrt{\frac{1}{2} \left(\frac{1}{n} + \frac{1}{m} \right)}.$$

Programul 2.5.12. Programul care urmează reprezintă pe aceeași figură graficele densităților de probabilitate pentru transformata lui Fisher Z și legea $\mathcal{N}(\mu, \sigma)$, unde μ și σ au valorile date prin (2.5.2).

Remarcăm faptul că densitatea de probabilitate a variabilei aleatoare Z este

$$f_Z(x) = 2e^{2x} f_X(e^{2x} | m, n).$$

```
clf
m = input('m='); n = input('n=');
mu = (1/n-1/m)/2; s = sqrt((1/n+1/m)/2);
x = mu-3*s:0.01:mu+3*s;
f1 = normpdf(x,mu,s);
f2 = 2*exp(2*x).*fpdf(exp(2*x),m,n);
plot(x,f1,'k-.',x,f2,'k-')
legend('Legea normala','Transformata Fisher',2)
```

Pentru $m=20$ și $n=30$, se obțin graficele din Figura 2.12.

Legea F (Fisher–Snedecor) necentrată (ncf)

Dacă variabilele aleatoare independente X și Y urmează respectiv legile $\chi^2(m, \delta)$ necentrată și $\chi^2(n, \delta)$ necentrată, atunci

$$F = \frac{X}{m} \bigg/ \frac{Y}{n}$$

urmează legea F (Fisher–Snedecor) necentrată, notată prin $\mathcal{Fnc}(m, n, \delta)$.

Funcția de repartiție pentru variabila aleatoare legea $\mathcal{Fnc}(m, n, \delta)$ este

$$F(x | m, n, \delta) = \sum_{k=0}^{\infty} \frac{\left(\frac{\delta}{2}\right)^k e^{-\frac{\delta}{2}}}{k!} I\left(\frac{mx}{n+mx} \mid \frac{m}{2} + k, \frac{n}{2}\right),$$

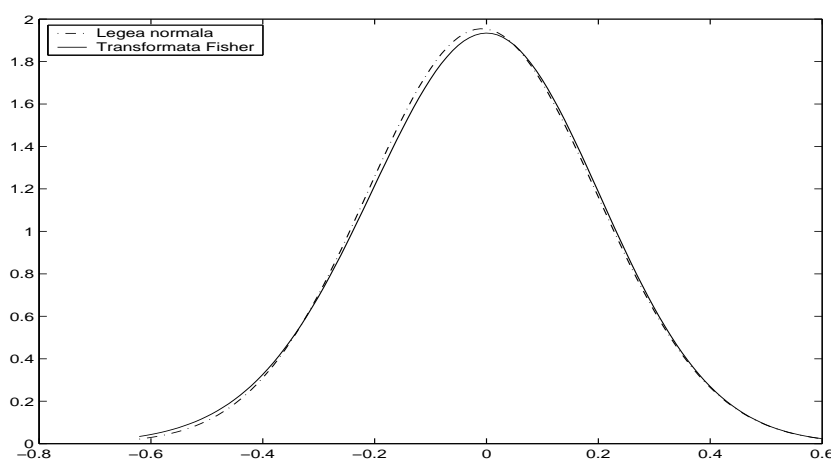


Figura 2.12: Transformata Fisher și legea normală corespunzătoare

unde $I(x|a, b)$ este funcția Beta incompletă de parametri $a, b > 0$.

Programul 2.5.13. Programul Matlab ce urmează reprezintă grafic densitățile de probabilitate pentru legile $\mathcal{F}(m, n)$ și $\mathcal{Fnc}(m, n, \delta)$:

```
clf;
m = input('m='); n = input('n=');
d = input('delta=');
x = 0:0.01:10.01;
Fc = fpdf(x,m,n); Fnc = ncfdpdf(x,m,n,d);
plot(x,Fc,'k-.',x,Fnc,'k-')
legend('Legea F', 'Legea F necentrata')
```

Pentru $m=5$, $n=20$ și $d=10$, se obțin graficele din Figura 2.13.

2.5.3 Funcția Matlab normspec

Având în vedere că legea normală joacă un rol important în teoria probabilităților și în statistică, sistemul Matlab îi acordă la rândul său o atenție specială. Astfel, în sistemul Matlab de bază există funcția `normspec`, cu următoarele moduri de apel:

```
normspec(sp,mu,sigma)
p = normspec(sp,mu,sigma)
```

Funcția `normspec` reprezintă grafic densitatea de probabilitate a legii normale de parametri μ și σ și umbrește aria mărginită de dreptele perpendiculare pe axa absciselor ce trec prin punctele axei precizate prin cele două componente ale vectorului `sp`, curba ce reprezintă graficul densității de probabilitate și axa absciselor.

Parametrul `p`, după execuția funcției `normspec`, conține valoarea ariei umbrite, adică probabilitatea ca variabila aleatoare să ia valori din intervalul precizat prin `sp`.

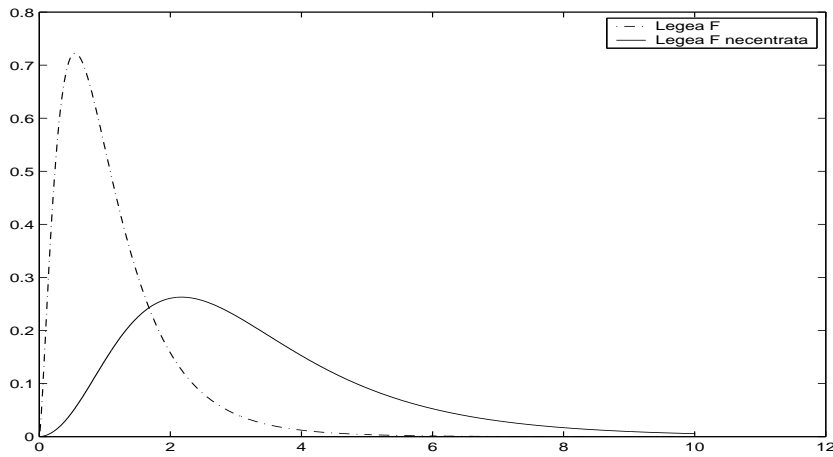


Figura 2.13: Legea $\mathcal{F}nc(5, 20, 10)$ și legea $\mathcal{F}(5, 20)$

Programul 2.5.14. Programul Matlab ce urmează va reprezenta grafic densitatea de probabilitate a legii $\mathcal{N}(\mu, \sigma)$ și va umbri aria cuprinsă între dreptele de ecuații $x = a$ și $x = b$ ($a < b$), curba ce reprezintă graficul densității de probabilitate și axa absciselor.

```
m = input('mu='); s = input('sigma=');
a = input('a:'); b = input('b (b>a):');
p = normspec([a,b],m,s);
xlabel(['a= ', num2str(a), '      b= ', num2str(b)])
ylabel('Densitatea de probabilitate')
title(['Probabilitatea intre a si b : ', num2str(p)])
```

Executând acest program cu valorile de intrare $\mu=0$, $\sigma=1$, $a=-1.5$ și $b=2$, se obține graficul din Figura 2.14.

Se observă că instrucțiunea `title` conține parametrul `num2str(p)`, care transformă valoarea numerică a lui `p` într-un șir de caractere.

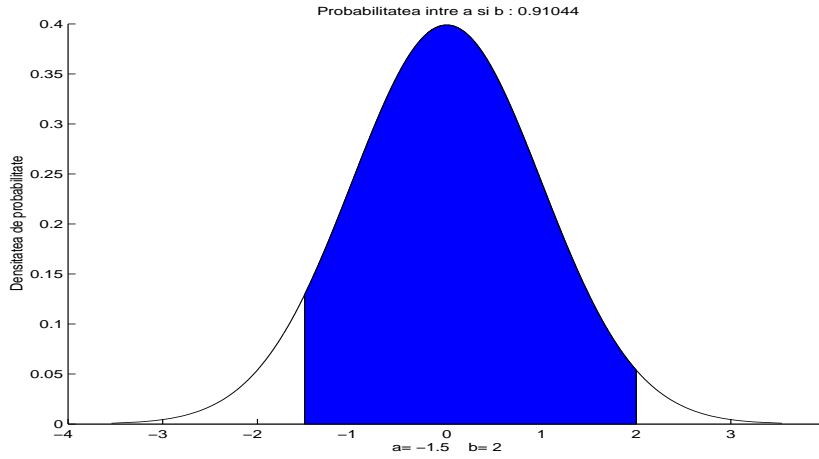
2.5.4 Distribuție marginală

Fie vectorul aleator $\mathbf{X} = (X_1, \dots, X_n)$ și vectorul aleator $\tilde{\mathbf{X}} = (X_{i_1}, \dots, X_{i_k})$, unde $1 \leq k \leq n$, iar $1 \leq i_1 < \dots < i_k \leq n$.

Definiția 2.5.15. Numim distribuție marginală a vectorului distribuția vectorului $\tilde{\mathbf{X}}$.

Dacă F este funcția de repartiție a vectorului \mathbf{X} , iar $F_{i_1 \dots i_k}$ este funcția vectorului aleator $\tilde{\mathbf{X}}$, numindu-se *funcție de repartiție marginală*, atunci

$$F_{i_1 \dots i_k}(x_{i_1}, \dots, x_{i_k}) = F(+\infty, \dots, x_{i_1}, \dots, x_{i_k}, \dots, +\infty),$$

Figura 2.14: Legea $\mathcal{N}(0, 1)$

adică, pentru obținerea funcției de repartiție a vectorului $\tilde{\mathbf{X}}$ se fac să tindă la ∞ toate argumentele funcției de repartiție F , cu excepția argumentelor x_{i_1}, \dots, x_{i_k} .

Dacă vectorul aleator \mathbf{X} este de tip continuu, având densitatea de probabilitate f , iar vectorul aleator $\tilde{\mathbf{X}}$ are densitatea de probabilitate $f_{i_1 \dots i_k}$, care se numește *densitate de probabilitate marginală*, atunci

$$f_{i_1 \dots i_k}(x_{i_1}, \dots, x_{i_k}) = \int_{-\infty}^{+\infty} \dots \int_{-\infty}^{\infty} f(x_1, \dots, x_n) dx_1 \dots / \dots / \dots dx_n,$$

adică $f_{i_1 \dots i_k}$ se obține integrând densitatea de probabilitate f pe \mathbb{R}^{n-k} , în raport cu toate variabilele independente, cu excepția variabilelor x_{i_1}, \dots, x_{i_k} .

2.5.5 Funcție de repartiție condiționată

Fie vectorul aleator bidimensional (X, Y) , având funcția de repartiție F .

Definiția 2.5.16. Numim funcție de repartiție condiționată a variabilei aleatoare X de către variabila aleatoare Y , funcția $F_{X|Y}: \mathbb{R} \rightarrow \mathbb{R}$, dată prin

$$F_{X|Y}(x|y) = P(X \leq x | Y = y), \quad \forall x \in \mathbb{R},$$

pentru fiecare $y \in \mathbb{R}$ fixat.

Trebuie să remarcăm faptul că formula de definiție se poate scrie sub forma

$$F_{X|Y}(x|y) = \lim_{d \searrow 0} P(X \leq x | y - d < Y \leq y)$$

sau

$$F_{X|Y}(x|y) = \lim_{d \searrow 0} \frac{P(X \leq x, y-d < Y \leq y)}{P(y-d < Y \leq y)} = \lim_{d \searrow 0} \frac{F(x, y) - F(x, y-d)}{F_Y(y) - F_Y(y-d)}.$$

Dacă vectorul aleator (X, Y) este de tip discret, adică are distribuția dată prin tabloul bidimensional

$X \setminus Y$	\dots	y_j	\dots
\vdots		\vdots	
x_i	\dots	p_{ij}	\dots
\vdots		\vdots	

atunci distribuția variabilei aleatoare X condiționată de $(Y = y_j)$ este dată prin

$$P(X = x_i | Y = y_j) = \frac{P(X = x_i, Y = y_j)}{P(Y = y_j)} = \frac{p_{ij}}{p_{\cdot j}},$$

unde $p_{\cdot j} = P(Y = y_j) = \sum_{i \in I} p_{ij}$.

Remarcăm de asemenea că are loc formula

$$(2.5.3) \quad p_{i\cdot} = P(X = x_i) = \sum_{j \in J} p_{\cdot j} p_{i|j},$$

unde $p_{i|j} = P(X = x_i | Y = y_j)$.

Dacă vectorul aleator (X, Y) este de tip continuu, adică are densitatea de probabilitate f , atunci densitatea de probabilitate a variabilei aleatoare X condiționată de $(Y = y)$ este dată prin

$$f_{X|Y}(x|y) = \frac{f(x, y)}{f_Y(y)}, \quad \text{pentru } f_Y(y) \neq 0.$$

Are loc formula (2.5.3) în variantă continuă:

$$f_X(x) = \int_{-\infty}^{+\infty} f_Y(y) f_{X|Y}(x|y) dy.$$

Exemplul 2.5.17 (Legea normală multidimensională). Vectorul aleator, notat prin \mathbf{X} , urmează *legea normală multidimensională*, notăm aceasta prin $\mathcal{N}(\boldsymbol{\mu}, \mathbf{V})$, dacă are densitatea de probabilitate

$$(2.5.4) \quad f(\mathbf{x} | \boldsymbol{\mu}, \mathbf{V}) = \frac{[\det(\mathbf{A})]^{\frac{1}{2}}}{(2\pi)^{\frac{n}{2}}} \times \exp \left[-\frac{1}{2} \sum_{i,j=1}^n a_{ij} (x_i - \mu_i) (x_j - \mu_j) \right], \quad \mathbf{x} \in \mathbb{R}^n, \quad \boldsymbol{\mu} \in \mathbb{R}^n,$$

unde V este o matrice pătratică de ordin n pozitiv definită, iar A este inversa matricei V .

Dacă $n = 2$ și notăm prin (X, Y) vectorul aleator bidimensional corespunzător, densitatea de probabilitate devine

$$(2.5.5) \quad f(x, y) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-r^2}} \times \\ \times \exp\left\{-\frac{1}{2(1-r^2)}\left[\frac{(x-\mu_1)^2}{\sigma_1^2} - 2r\frac{(x-\mu_1)(y-\mu_2)}{\sigma_1\sigma_2} + \frac{(y-\mu_2)^2}{\sigma_2^2}\right]\right\},$$

unde $\sigma_1, \sigma_2 > 0$ și $|r| < 1$.

Se cunoaște că fiecare din componentele X și Y ale vectorului aleator urmează legea normală respectiv $\mathcal{N}(\mu_1, \sigma_1)$ și $\mathcal{N}(\mu_2, \sigma_2)$.

Densitatea de probabilitate condiționată se exprimă în acest fel prin:

$$(2.5.6) \quad f_{X|Y}(x|y) = \frac{1}{\sigma_1\sqrt{2\pi(1-r^2)}} e^{-\frac{(x-m(y))^2}{2(1-r^2)\sigma_1^2}},$$

unde $m(y) = \mu_1 + r\frac{\sigma_1}{\sigma_2}(y - \mu_2)$, adică se obține densitatea de probabilitate a legii normale $\mathcal{N}(m(y), \sigma_1\sqrt{1-r^2})$.

Programul 2.5.18. Scriem în cele ce urmează un program Matlab, care reprezintă grafic densitatea de probabilitate a legii normale bidimensionale.

```
clf, clear
m1 = input('mu1='); m2 = input('mu2=');
s1 = input('sigma1='); s2 = input('sigma2=');
r = input('r(-1<r<1):');
x = m1-3*s1:.2:m1+3*s1; y = m2-3*s2:.2:m2+3*s2;
[X,Y] = meshgrid(x,y);
Z = 1/(2*pi*s1*s2*sqrt(1-r^2))*exp(-1/(2*(1-r^2))*...
    *((X-m1).^2/s1^2-2*r*(X-m1).*(Y-m2)/(s1*s2)...
    +(Y-m2).^2/s2^2));
mesh(X,Y,Z)
title('Legea normala bidimensionala')
```

Prin executarea acestui program, pentru $m1=m2=0$, $s1=1$, $s2=2$ și $r=-0.5$, se obține graficul din Figura 2.15.

Programul 2.5.19. Programul Matlab care urmează reprezintă grafic densitatea de probabilitate condiționată în cazul unui vector aleator ce urmează legea normală bi-dimensională.

```
clf, clear
m1 = input('mu1='); m2 = input('mu2=');
s1 = input('sigma1='); s2 = input('sigma2=');
```

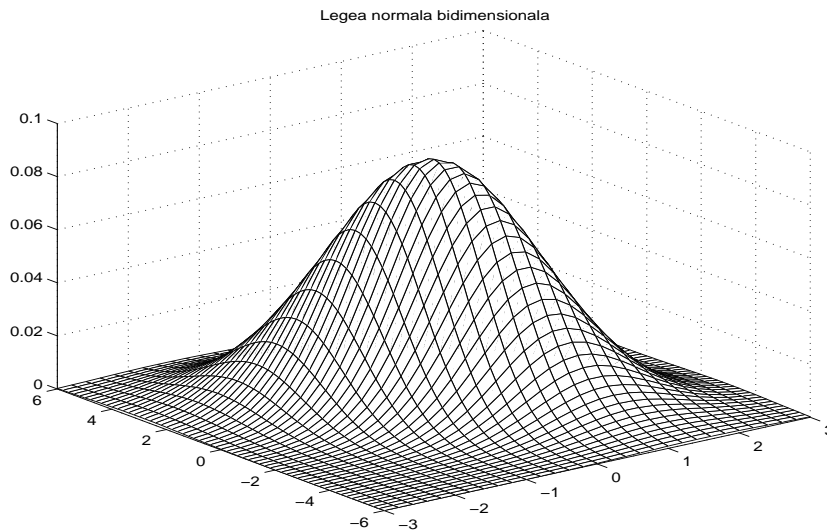


Figura 2.15: Legea $\mathcal{N}(0, 0; 1, 2; -0.5)$

```

r = input('r(-1<r<1):');
Y(1:3) = input('Y(trei valori):');
m = m1+s1*r*(Y-m2)/s2; s = sqrt(1-r^2)*s1;
MM = max(m); mm = min(m);
X = []; X = mm-3*s:0.01:MM+3*s;
hold on
for i=1:3
    P = []; P = normpdf(X,m(i),s);
    if i==1 plot(X,P,'k-.'), end
    if i==2 plot(X,P,'k-'), end
    if i==3 plot(X,P,'k--'), end
end
title('Legea normala conditionata')

```

Prin executarea acestui program, pentru $m_1=1$, $m_2=0$, $s_1=1$, $s_2=2$, $r=-0.5$ și $Y=[-2, 0, 2]$, se obțin graficele din Figura 2.16. Graficul marcat prin punct–linie este pentru prima valoare a lui Y , cel cu linie continuă, pentru a doua valoare, iar prin linie întreruptă, pentru a treia.

2.5.6 Funcție de supraviețuire. Funcție hazard

În aplicații tehnice și nu numai, sunt utilizate și alte funcții atașate variabilelor aleatoare, pentru exprimarea mai potrivită a fenomenului aleator cercetat.

Definiția 2.5.20. Dacă variabila aleatoare X are funcția de repartiție F , se numește

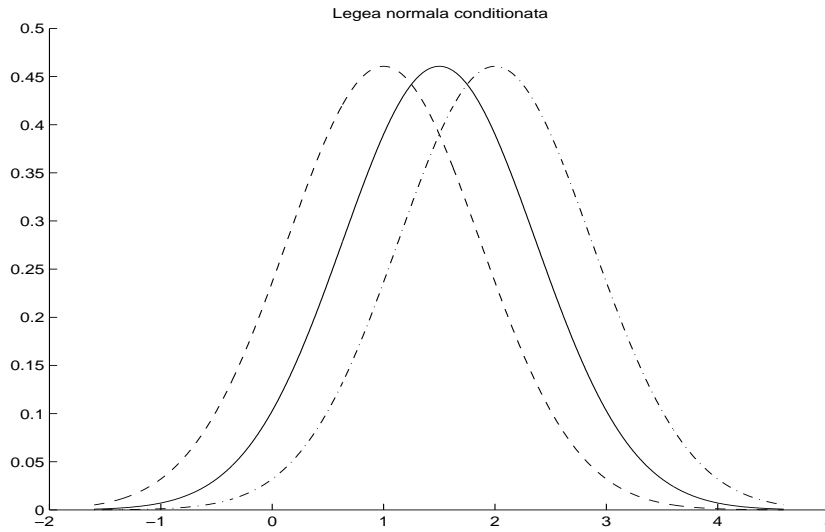


Figura 2.16: Legea normală condiționată pentru $y = -2, 0, 2$

funcție de supraviețuire *atașată variabilei aleatoare X funcția definită prin:*

$$S(x) = 1 - F(x) = 1 - P(X > x), \quad \forall x \in \mathbb{R}.$$

Se observă că dacă X reprezintă durata de viață a unui individ, atunci $S(x)$ reprezintă probabilitatea ca durată de viață să depășească un prag x dat.

Definiția 2.5.21. Dacă variabila aleatoare de tip continuu X are funcția de repartiție F și densitatea de probabilitate f , se numește funcție hazard sau funcție de risc instantaneu *atașată variabilei aleatoare X funcția definită prin:*

$$h(x) = \frac{f(x)}{S(x)} = -\frac{S'(x)}{S(x)}, \quad \forall x \in \mathbb{R}.$$

Având în vedere definiția densității de probabilitate putem scrie

$$f(x) dx \cong P(x \leq X < x + dx),$$

de unde

$$h(x) dx \cong \frac{P(x \leq X < x + dx)}{P(X > x)},$$

adică

$$h(x) dx \cong P(x \leq X < x + dx | X > x).$$

Prin urmare $h(x)$ ar reprezenta probabilitatea ca valoarea variabilei aleatoare X să rămână neschimbată într-un interval imediat următor, de lungime mică, dacă valoarea ei a depășit pragul x .

2.6 Caracteristici numerice

2.6.1 Valoare medie. Dispersie (varianță). Covarianță

Fie variabila aleatoare X având funcția de repartiție F .

Definiția 2.6.1. Numim valoare medie sau speranță matematică, *caracteristica numerică*

$$(2.6.1) \quad E(X) = \int_{-\infty}^{+\infty} x dF(x),$$

unde integrala Stieltjes se impune să fie absolut convergentă pentru ca valoarea medie să existe.

Observația 2.6.2. Dacă variabila aleatoare este de tip discret, adică are distribuția $X \left(\begin{smallmatrix} x_i \\ p_i \end{smallmatrix} \right)_{i \in I}$, atunci integrala Stieltjes se reduce la o sumă, anume

$$E(X) = \sum_{i \in I} x_i p_i,$$

care există pentru cazul când X este variabilă aleatoare simplă, altfel se impune ca seria ce apare în această formulă să fie absolut convergentă.

Dacă variabila aleatoare este de tip continuu, adică are densitatea de probabilitate f , atunci

$$E(X) = \int_{-\infty}^{+\infty} x f(x) dx.$$

Mai remarcăm faptul că dacă \mathbf{X} este un vector aleator n -dimensional, iar funcția $h: \mathbb{R}^n \rightarrow \mathbb{R}$ este astfel încât $Y = h(\mathbf{X})$ să fie o variabilă aleatoare, atunci

$$E(Y) = \int_{-\infty}^{+\infty} x dF_Y(x) = \int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty} h(\mathbf{x}) dF_{\mathbf{X}}(\mathbf{x}),$$

care în cazul continuu devine

$$E(Y) = \int_{-\infty}^{+\infty} x f_Y(x) dx = \int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty} h(\mathbf{x}) f_{\mathbf{X}}(\mathbf{x}).$$

Definiția 2.6.3. Numim dispersie sau varianță caracteristica numerică atașată variabilei aleatoare X definită prin:

$$(2.6.2) \quad Var(X) = E[(X - E(X))^2].$$

Remarcăm formula de calcul $Var(X) = E(X^2) - [E(X)]^2$.

Definiția 2.6.4. Dacă se consideră vectorul aleator bidimensional (X, Y) , numim covarianță sau corelație, caracteristica numerică

$$(2.6.3) \quad Cov(X, Y) = E[(X - E(X))(Y - E(Y))],$$

respectiv coeficient de corelație, raportul

$$(2.6.4) \quad r(X, Y) = \frac{Cov(X, Y)}{\sqrt{Var(X)Var(Y)}}.$$

Observația 2.6.5. Pentru suma unui număr finit de variabile aleatoare avem:

$$Var\left(\sum_{k=1}^n X_k\right) = \sum_{k=1}^n Var(X_k) + 2 \sum_{i < k} Cov(X_i, X_k),$$

iar dacă variabilele aleatoare sunt independente două câte două

$$Var\left(\sum_{k=1}^n X_k\right) = \sum_{k=1}^n Var(X_k).$$

Definiția 2.6.6. Dacă avem vectorul aleator $\mathbf{X} = (X_1, \dots, X_n)^\top$, numim valoare medie, vectorul

$$E(\mathbf{X}) = (E(X_1), \dots, E(X_n))^\top,$$

respectiv matricea covarianțelor, matricea

$$V(\mathbf{X}) = (Cov(X_i, X_k))_{i,k=\overline{1,n}}.$$

Observația 2.6.7. Din cauza simetriei covarianței, avem că matricea covarianțelor este simetrică.

Dacă componentele vectorului aleator \mathbf{X} sunt independente două câte două, atunci matricea covarianțelor este o matrice diagonală, pe diagonala principală aflându-se dispersiile componentelor vectorului aleator.

Având în vedere extinderea naturală a valorii medii la matrice aleatoare, putem scrie:

$$V(\mathbf{X}) = E[(\mathbf{X} - E(\mathbf{X}))(\mathbf{X} - E(\mathbf{X}))^\top].$$

Teorema 2.6.8. Fie vectorul aleator $\mathbf{X} = (X_1, \dots, X_n)^\top$ și matricea \mathbf{A} , de tipul (m, n) , cu elemente numere reale. Dacă se consideră vectorul $\mathbf{Y} = \mathbf{A}\mathbf{X}$, atunci

(i) \mathbf{Y} este un vector aleator;

(ii) valoarea medie a vectorului aleator \mathbf{Y} este

$$E(\mathbf{Y}) = \mathbf{A}E(\mathbf{X}) \quad \text{sau} \quad E(\mathbf{Y}^\top) = E(\mathbf{X}^\top) \mathbf{A}^\top;$$

(iii) matricea covarianțelor vectorului aleator \mathbf{Y} este dată prin

$$(2.6.5) \quad V(\mathbf{Y}) = \mathbf{A} V(\mathbf{X}) \mathbf{A}^\top.$$

Demonstrație. (i) Transformarea liniară a unui vector aleator este vector aleator.

(ii) Având în vedere proprietăți ale valorii medii, pentru o componentă a vectorului $E(\mathbf{Y})$, se poate scrie

$$E(Y_i) = E\left(\sum_{k=1}^n a_{ik} X_k\right) = \sum_{k=1}^n a_{ik} E(X_k),$$

de unde rezultă $E(\mathbf{Y}) = \mathbf{A}E(\mathbf{X})$.

În mod analog se arată și a doua relație.

(iii) Se poate proceda ca și pentru valoarea medie, dacă se ține seama de proprietăți ale corelației. Se scrie succesiv:

$$\begin{aligned} Cov(Y_i, Y_k) &= Cov\left(\sum_{p=1}^n a_{ip} X_p, \sum_{q=1}^n a_{kq} X_q\right) \\ &= E\left\{\left[\sum_{p=1}^n a_{ip} (X_p - E(X_p))\right] \left[\sum_{q=1}^n a_{kq} (X_q - E(X_q))\right]\right\} \\ &= \sum_{p=1}^n \sum_{q=1}^n a_{ip} E[(X_p - E(X_p))(X_q - E(X_q))] a_{kq} \\ &= \sum_{p=1}^n \sum_{q=1}^n a_{ip} Cov(X_p, X_q) a_{kq}, \end{aligned}$$

de unde se obține relația (2.6.5).

Același lucru se poate obține și prin folosirea relației de la punctul (ii). Anume,

se poate scrie:

$$\begin{aligned}
 V(\mathbf{Y}) &= E[(\mathbf{Y} - E(\mathbf{Y}))(\mathbf{Y} - E(\mathbf{Y}))^\top] \\
 &= E[(\mathbf{A}\mathbf{X} - \mathbf{A}E(\mathbf{X}))(\mathbf{A}\mathbf{X} - \mathbf{A}E(\mathbf{X}))^\top] \\
 &= E[\mathbf{A}(\mathbf{X} - E(\mathbf{X}))(\mathbf{X} - E(\mathbf{X}))^\top \mathbf{A}^\top] \\
 &= \mathbf{A}E[(\mathbf{X} - E(\mathbf{X}))(\mathbf{X} - E(\mathbf{X}))^\top] \mathbf{A}^\top \\
 &= \mathbf{A}V(\mathbf{X}) \mathbf{A}^\top,
 \end{aligned}$$

de unde se obține din nou relația (2.6.5). \square

Proprietatea 2.6.9. Dacă $V(\mathbf{X})$ este matricea covarianțelor vectorului aleator \mathbf{X} , atunci aceasta este o matrice pozitiv semi-definită, adică

$$\mathbf{a}^\top V(\mathbf{X}) \mathbf{a} \geq 0, \quad \forall \mathbf{a} \in \mathbb{R}^n.$$

Demonstrație. Pornind din membrul drept al relației ce se dorește a fi demonstrată, putem scrie succesiv:

$$\begin{aligned}
 \mathbf{a}^\top V(\mathbf{X}) \mathbf{a} &= \sum_{i=1}^n \sum_{k=1}^n a_i a_k \text{Cov}(X_i, X_k) = \text{Cov}\left(\sum_{i=1}^n a_i X_i, \sum_{k=1}^n a_k X_k\right) \\
 &= \text{Var}\left(\sum_{i=1}^n a_i X_i\right) \geq 0,
 \end{aligned}$$

de unde se obține relația din enunț. \square

Proprietatea 2.6.10. Dacă matricea \mathbf{W} este o matrice pătratică pozitiv definită, adică

$$\mathbf{a}^\top V(\mathbf{X}) \mathbf{a} > 0, \quad \forall \mathbf{a} \in \mathbb{R}^n \setminus \{\mathbf{0}\},$$

atunci \mathbf{W} poate fi considerată matricea covarianțelor unui vector aleator.

Demonstrație. Din algebra liniară este cunoscut faptul că o matrice pozitiv definită admite decompunerea

$$\mathbf{W} = \mathbf{U}^\top \mathbf{A} \mathbf{U} \quad \text{sau} \quad \mathbf{A} = \mathbf{U} \mathbf{W} \mathbf{U}^\top,$$

unde $\mathbf{U}^\top \mathbf{U} = \mathbf{U} \mathbf{U}^\top = \mathbf{I}$, \mathbf{I} fiind matricea unitate, iar \mathbf{A} este o matrice diagonală cu elementele numere pozitive. Putem face precizarea că diagonală matricei \mathbf{A} conține valorile proprii ale matricei \mathbf{W} , iar matricea \mathbf{U} este formată din vectorii proprii ortonormați corespunzători.

Considerăm vectorul aleator \mathbf{Z} având matricea covarianțelor \mathbf{I} . Un astfel de vector aleator există, având în vedere de exemplu vectorul aleator care are ca și componente variabile aleatoare independente ce urmează fiecare legea normală $\mathcal{N}(0, 1)$.

Vectorul aleator $\mathbf{X} = \mathbf{U}^\top \mathbf{A}^{\frac{1}{2}} \mathbf{Z}$ va avea ca matrice a covarianțelor chiar matricea \mathbf{W} .

Întrădevăr, dacă se ține seama de relația (2.6.5), se obține:

$$V(\mathbf{X}) = \mathbf{U}^\top \mathbf{A}^{\frac{1}{2}} V(\mathbf{Z}) \left(\mathbf{U}^\top \mathbf{A}^{\frac{1}{2}} \right)^\top = \mathbf{U}^\top \mathbf{A}^{\frac{1}{2}} \mathbf{A}^{\frac{1}{2}} \mathbf{U} = \mathbf{U}^\top \mathbf{A} \mathbf{U} = \mathbf{W}.$$

Din această succesiune de egalități se obține afirmația făcută în enunț. \square

Corolarul 2.6.11. *Dacă matricea covarianțelor $V(\mathbf{X})$ a unui vector aleator \mathbf{X} este pozitiv definită, atunci există o transformare $\mathbf{X} \mapsto \mathbf{Y}$ astfel încât pentru noul vector aleator \mathbf{Y} să avem $E(\mathbf{Y}) = \mathbf{0}$ și $V(\mathbf{Y}) = \mathbf{I}$.*

Demonstrație. Folosind descompunerea $V(\mathbf{X}) = \mathbf{U}^\top \mathbf{A} \mathbf{U}$ și alegând $\mathbf{A} = \mathbf{A}^{-\frac{1}{2}} \mathbf{U}$, obținem că $\mathbf{Y} = \mathbf{A}[\mathbf{X} - E(\mathbf{X})]$ este vectorul aleator căutat.

Pe de o parte avem

$$E(\mathbf{Y}) = \mathbf{A}[E(\mathbf{X}) - E(\mathbf{X})] = \mathbf{A}\mathbf{0} = \mathbf{0}.$$

Pe de altă parte

$$\begin{aligned} V(\mathbf{Y}) &= \mathbf{A} V(\mathbf{X} - E(\mathbf{X})) \mathbf{A}^\top = \mathbf{A} V(\mathbf{X}) \mathbf{A}^\top \\ &= \mathbf{A}^{-\frac{1}{2}} \mathbf{U} V(\mathbf{X}) \mathbf{U}^\top \mathbf{A}^{-\frac{1}{2}} = \mathbf{A}^{-\frac{1}{2}} \mathbf{A} \mathbf{A}^{-\frac{1}{2}} = \mathbf{I}. \end{aligned}$$

Din cel două succesiuni de relații se obține că vectorul \mathbf{Y} satisface condițiile dorite. \square

Corolarul 2.6.12. *Dacă matricea covarianțelor $V(\mathbf{X})$ a unui vector aleator \mathbf{X} este pozitiv definită, atunci variabila aleatoare*

$$S^2 = [\mathbf{X} - E(\mathbf{X})]^\top [V(\mathbf{X})]^{-1} [\mathbf{X} - E(\mathbf{X})],$$

are valoarea medie $E(S^2) = n$.

Demonstrație. Cu notațiile din demonstrația corolarului precedent avem

$$\mathbf{A}^\top \mathbf{A} = \mathbf{U}^\top \mathbf{A}^{-\frac{1}{2}} \mathbf{A}^{-\frac{1}{2}} \mathbf{U} = \mathbf{U}^\top \mathbf{A}^{-1} \mathbf{U} = [V(\mathbf{X})]^{-1}.$$

Prin urmare, avem succesiv:

$$S^2 = [\mathbf{X} - E(\mathbf{X})]^\top \mathbf{A}^\top \mathbf{A} [\mathbf{X} - E(\mathbf{X})] = \mathbf{Y}^\top \mathbf{Y} = \sum_{i=1}^n Y_i^2.$$

Dar avem că $E(Y_i) = 0$ și $Var(Y_i) = 1$, astfel că $E(Y_i^2) = Var(Y_i) = 1$, pentru fiecare $i = \overline{1, n}$, ceea ce conduce la $E(S^2) = n$. \square

2.6.2 Funcții Matlab pentru valoare medie și dispersie

Sistemul Matlab dispune de proceduri pentru calculul valorilor medii și ale dispersiilor legilor de probabilitate implementate prin *Statistics toolbox*.

Pentru calculul valorilor acestor caracteristici numerice se folosește apelul

```
[m,v]=numef(x,par1,par2,...)
```

unde numef este un șir de caractere, care definesc legea de probabilitate, din care ultimele patru sunt stat, iar cele ce le preced sunt cele care dau numele legii.

În urma executării instrucțiunii, se calculează matricele m și v ale valorilor medii și ale dispersiilor pentru legea considerată, având parametrii precizați prin matricele par1, par2, ... Aceste matrice trebuie să fie de aceleași dimensiuni, cu excepția că dacă unele sunt scalari, aceștia se extind la matricele constante de aceleași dimensiuni cu celelalte și care iau valorile scalarilor corespunzători.

Valorile medii și dispersiile se calculează, pentru fiecare lege de probabilitate în parte, folosindu-se formulele din Tabelul 2.1.

2.6.3 Valoare medie condiționată. Dispersie (varianță) condiționată

Fie vectorul aleator bidimensional (X, Y) și $F_{X|Y}$ funcția de repartiție a lui X condiționată de Y .

Definiția 2.6.13. Numim valoare medie condiționată a variabilei aleatoare X de către Y , variabila aleatoare

$$E(X|Y) = \int_{-\infty}^{+\infty} x dF_{X|Y}(x|Y),$$

o realizare a variabilei aleatoare $E(X|Y)$ fiind

$$E(X|y) = \int_{-\infty}^{+\infty} x dF_{X|Y}(x|y), \quad y \in \mathbb{R}.$$

Dacă vectorul aleator (X, Y) este de tip discret, cu distribuția

$X \setminus Y$...	y_j	...
\vdots		\vdots	
x_i	...	p_{ij}	...
\vdots		\vdots	

atunci variabila aleatoare $E(X|Y)$ va avea distribuția

$$E(X|Y) \left(\begin{matrix} E(X|y_j) \\ p_{\cdot j} \end{matrix} \right)_{j \in J},$$

Legea	Denumirea	Valoarea medie și dispersia
$\mathcal{U}(N)$	unid	$E(X) = \frac{N+1}{2}, Var(X) = \frac{N^2-1}{12}$
$\mathcal{B}(n, p)$	bino	$E(X) = np, Var(X) = np(1-p)$
$\mathcal{H}(n, M, K)$	hyge	$E(X) = n \frac{K}{M}, Var(X) = n \frac{K}{M} \frac{M-K}{M} \frac{M-n}{M-1}$
$\mathcal{P}o(\lambda)$	poiss	$E(X) = \lambda, Var(X) = \lambda$
$\mathcal{BN}(r, p)$	nbins	$E(X) = \frac{r(1-p)}{p}, Var(X) = \frac{r(1-p)}{p^2}$
$\mathcal{G}e(p)$	geo	$E(X) = \frac{1-p}{p}, Var(X) = \frac{1-p}{p^2}$
$\mathcal{U}(a, b)$	unif	$E(X) = \frac{a+b}{2}, Var(X) = \frac{(b-a)^2}{12}$
$\mathcal{N}(\mu, \sigma)$	norm	$E(X) = \mu, Var(X) = \sigma^2$
$\mathcal{LN}(\mu, \sigma)$	logn	$E(X) = e^{\mu + \frac{\sigma^2}{2}}, Var(X) = e^{2\mu + 2\sigma^2} - e^{2\mu + \sigma^2}$
$\mathcal{G}a(a, b)$	gam	$E(X) = ab, Var(X) = ab^2$
$\mathcal{E}xp(\mu)$	exp	$E(X) = \mu, Var(X) = \mu^2$
$\mathcal{B}eta(a, b)$	beta	$E(X) = \frac{a}{a+b}, Var(X) = \frac{ab}{(a+b+1)(a+b)^2}$
$\mathcal{W}(a, b)$	weib	$E(X) = a^{-\frac{1}{b}} \Gamma\left(1 + \frac{1}{b}\right),$ $Var(X) = a^{-\frac{2}{b}} \left[\Gamma\left(1 + \frac{2}{b}\right) - \Gamma^2\left(1 + \frac{1}{b}\right) \right]$
$\mathcal{R}(b)$	rayl	$E(X) = b\sqrt{\frac{\pi}{2}}, Var(X) = \frac{4-\pi}{2} b^2$
$\mathcal{T}(n)$	t	$E(X) = 0, Var(X) = \frac{n}{n-2}, n > 2$
$\mathcal{T}nc(n, \delta)$	nct	$E(X) = \frac{\delta \sqrt{\frac{n}{2}} \Gamma\left(\frac{n-1}{2}\right)}{\Gamma\left(\frac{n}{2}\right)},$ $Var(X) = \frac{n(1+\delta^2)}{n-2} - \frac{n\delta^2}{2} \left[\frac{\Gamma\left(\frac{n-1}{2}\right)}{\Gamma\left(\frac{n}{2}\right)} \right]^2, n > 2$
$\chi^2(n)$	chi2	$E(X) = n, Var(X) = 2n$
$\chi^2(n, \delta)$	ncx2	$E(X) = n + \delta, Var(X) = 2(n + 2\delta)$
$\mathcal{F}(m, n)$	f	$E(X) = \frac{n}{n-2}, n > 2,$ $Var(X) = \frac{2n^2(m+n-2)}{m(n-2)^2(n-4)}, n > 4$
$\mathcal{F}nc(m, n, \delta)$	ncf	$E(X) = \frac{n(m+\delta)}{m(n-2)}, n > 2,$ $Var(X) = 2 \left(\frac{n}{m} \right)^2 \frac{(m+\delta)^2 + (m+2\delta)(n-2)}{(n-2)^2(n-4)}, n > 4$

Tabelul 2.1: Tabelul valorilor medii și dispersiilor

unde

$$E(X|y_j) = \sum_{i \in I} x_i P(X = x_i | Y = y_j) = \sum_{i \in I} x_i \frac{p_{ij}}{p_{\cdot j}}, \quad j \in J.$$

Dacă vectorul aleator (X, Y) este de tip continuu, având densitatea de probabilitate f , atunci o realizare a variabilei aleatoare $E(X|Y)$ este dată prin

$$E(X|y) = \int_{-\infty}^{+\infty} x f_{X|Y}(x|y) dx = \int_{-\infty}^{+\infty} x \frac{f(x, y)}{f_Y(y)} dy.$$

Proprietatea 2.6.14. $E[E(X|Y)] = E(X)$.

Demonstrație. Demonstrăm proprietatea în cazul continuu. Cazul discret se demonstrează analog.

Folosind definiția valorii medii avem succesiv:

$$\begin{aligned} E[E(X|Y)] &= \int_{-\infty}^{+\infty} E(X|y) f_Y(y) dy = \int_{-\infty}^{+\infty} f_Y(y) \left[\int_{-\infty}^{+\infty} x \frac{f(x, y)}{f_Y(y)} dx \right] dy \\ &= \int_{-\infty}^{+\infty} x \left[\int_{-\infty}^{+\infty} f(x, y) dy \right] dx = \int_{-\infty}^{+\infty} x f_X(dx) = E(X). \end{aligned}$$

Demonstrația se încheie dacă se rețin extremitățile acestui șir de egalități. \square

Definiția 2.6.15. Fie vectorul aleator (X, Y) . Numim dispersie condiționată sau varianță condiționată a variabilei aleatoare X de către Y , variabila aleatoare

$$Var(X|Y) = E[(X - E(X|Y))^2 | Y].$$

Proprietatea 2.6.16. $Var(X) = E[Var(X|Y)] + Var[E(X|Y)]$.

Demonstrație. Pornind de la definiția varianței și folosind proprietăți ale ei, putem scrie:

$$\begin{aligned} Var(X) &= E[(X - E(X))^2] = E[(X - E(X|Y) + E(X|Y) - E(X))^2] \\ &= E[(X - E(X|Y))^2] + E[(E(X|Y) - E(X))^2] \\ &\quad + 2E[(X - E(X|Y))(E(X|Y) - E(X))]. \end{aligned}$$

Vom calcula pe rând cei trei termeni din partea dreaptă a acestei relații.

Pentru primul termen, dacă se notează $Z = (X - E(X|Y))^2$ și se aplică proprietatea valorii medii condiționate avem:

$$E(Z) = E[E(Z|Y)] = E\left\{E[(X - E(X|Y))^2 | Y]\right\},$$

de unde se obține că

$$E \left[(X - E(X|Y))^2 \right] = E[Var(X|Y)].$$

Pentru al doilea termen, dacă se notează $Z = E(X|Y)$, se poate scrie:

$$\begin{aligned} E \left[(E(X|Y) - E(X))^2 \right] &= E \left[(E(X|Y) - E(E(X|Y)))^2 \right] \\ &= E \left[(Z - E(Z))^2 \right] = Var(Z) = Var[E(X|Y)], \end{aligned}$$

de unde

$$E \left[(E(X|Y) - E(X))^2 \right] = Var[E(X|Y)].$$

A mai rămas să arătăm ca al treilea termen este zero.

Fixăm $Y = y$ și având în vedere că $E(X|y) - E(X) = \text{const.}$, putem scrie:

$$\begin{aligned} E[(X - E(X|y))(E(X|y) - E(X))] &= [E(X|y) - E(X)] E[X - E(X|y)] \\ &= [E(X|y) - E(X)] [E(X) - E(X)] = 0. \end{aligned}$$

Ceea ce a mai rămas de arătat. □

Exemplul 2.6.17. Să considerăm vectorul aleator (X, Y) , care urmează legea normală, având densitatea dată prin (2.5.5). Este cunoscut că parametrii acestei legi de probabilitate au următoarele interpretări probabilistice:

$$\mu_1 = E(X), \quad \mu_2 = E(Y), \quad \sigma_1^2 = Var(X), \quad \sigma_2^2 = Var(Y),$$

iar $r \in (-1, 1)$ este coeficientul de corelație dintre X și Y .

Mai mult, pentru legea normală multidimensională, având densitatea de probabilitate dată prin (2.5.4), se arată că μ este valoarea medie a vectorului aleator ce urmează legea normală multidimensională, iar V este matricea covarianțelor.

Să revenim la cazul bidimensional și să constatăm că

$$\begin{aligned} E(X|Y = y) &= m(y) = \mu_1 + \frac{\sigma_1}{\sigma_2} r (y - \mu_2), \\ Var(X|Y = y) &= \sigma_1^2 (1 - r^2). \end{aligned}$$

În mod analog, avem că

$$\begin{aligned} E(Y|X = x) &= m(x) = \mu_2 + \frac{\sigma_2}{\sigma_1} r (x - \mu_1), \\ Var(Y|X = x) &= \sigma_2^2 (1 - r^2). \end{aligned}$$

Curbele de ecuații

$$\begin{aligned}x(y) &= \mu_1 + \frac{\sigma_1}{\sigma_2} r(y - \mu_2), \\y(x) &= \mu_2 + \frac{\sigma_2}{\sigma_1} r(x - \mu_1),\end{aligned}$$

se numesc *curbele de regresie* a lui X în raport cu Y , respectiv a lui Y în raport cu X .

2.6.4 Funcție caracteristică

Fie variabila aleatoare X și vectorul aleator n -dimensional \mathbf{X} .

Definiția 2.6.18. Numim funcție caracteristică atașată variabilei aleatoare X , respectiv vectorului aleator \mathbf{X} , funcția definită prin

$$\begin{aligned}\varphi(t) &= M(e^{itX}), \quad \forall t \in \mathbb{R}, \\ \varphi(\mathbf{t}) &= M\left(e^{i \sum_{k=1}^n t_k X_k}\right), \quad \forall \mathbf{t} \in \mathbb{R}^n.\end{aligned}$$

Utilitatea funcției caracteristice este dată de netezimea ei, care este superioară funcției de repartiție, această din urmă putând fi discontinuă, pe când funcția caracteristică este uniform continuă. În plus, pe baza formulei de inversiune, există o caracterizare completă a distribuției unei variabile aleatoare sau a unui vector aleator cu ajutorul funcției caracteristice.

2.6.5 Momente

Definiția 2.6.19. Fie variabila aleatoare X . Numim moment (inițial) și respectiv moment centrat de ordin n , caracteristicile numerice

$$\nu_n = E(X^n), \quad \mu_n = E[(X - E(X))^n] = E[(X - \nu_1)^n].$$

Remarcăm că în statistică sunt folosite în primul rând primele patru momente.

Astfel momentele centrate de ordinele doi și trei se folosesc pentru definirea *asimetriei* (*skewness*)

$$a = \frac{\mu_3}{\mu_2^{\frac{3}{2}}},$$

iar momentele centrate de ordinele doi și patru se folosesc pentru definirea *excesului* (*kurtosis*)

$$e = \frac{\mu_4}{\mu_2^2} \quad \text{sau} \quad e = \frac{\mu_4}{\mu_2^2} - 3.$$

În sistemul Matlab, excesul este definit cu prima formulă.

2.6.6 Mediana. Cuartile. Cuantile

Definiția 2.6.20. Fie variabila aleatoare X , care are funcția de repartiție F . Numim mediană, caracteristica numerică m , care satisface condițiile

$$(2.6.6) \quad P(X \geq m) \geq \frac{1}{2} \leq P(X \leq m).$$

Având în vedere proprietăți ale funcției de repartiție, condițiile (2.6.6) se pot scrie echivalent sub forma

$$F(m-0) \leq \frac{1}{2} \leq F(m).$$

Dacă variabila aleatoare este de tip continuu, relațiile (2.6.6) se reduc la relația $F(m) = \frac{1}{2}$, adică $m = F^{-1}(\frac{1}{2})$.

În cazul discret, folosind (2.6.6) s-ar putea ca mediana să nu fie determinată în mod unic, de aceea se alege m ca fiind cel mai mic număr pentru care $F(m) \geq \frac{1}{2}$.

Definiția 2.6.21. Fie variabila aleatoare X , care are funcția de repartiție F . Numim cuantilă de ordin $\gamma \in (0, 1)$, caracteristica numerică x_γ , care satisface condițiile

$$(2.6.7) \quad P(X \geq x_\gamma) \geq 1 - \gamma \quad \text{și} \quad P(X \leq m) \geq \gamma.$$

Ca și în cazul medianei, aceste condiții se pot scrie echivalent sub forma

$$F(x_\gamma - 0) \leq \gamma \leq F(x_\gamma),$$

care se reduce la relația $F(x_\gamma) = \gamma$, adică $x_\gamma = F^{-1}(\gamma)$, în cazul continuu.

În cazul discret, cuantila x_γ se ia ca fiind cel mai mic număr pentru care are loc inegalitatea $F(x_\gamma) \geq \gamma$.

Unele cuantile au denumiri speciale. Iată de exemplu, cuantila de ordin $\gamma = \frac{1}{2}$ este chiar mediana.

Dacă $\gamma = \frac{1}{4}, \frac{2}{4}, \frac{3}{4}$, se obțin ceea ce se numesc *cuartile*, respectiv *cuartila inferioară*, *mediana* și *cuartila superioară*.

Dacă $\gamma = \frac{i}{10}, i = \overline{1, 9}$, $\gamma = \frac{k}{100}, k = \overline{1, 99}$, se obțin ceea ce poartă denumirile *decile* și respectiv *centile*.

2.6.7 Funcția Matlab `icdf`

Am văzut că pentru calculul cuantilelor este necesară inversarea funcției de repartiție. Sistemul Matlab prin *Statistics toolbox* dispune de funcții pentru inversarea funcțiilor de repartiție ale legilor de probabilitate implementate.

Apelarea acestor funcții se face cu una din următoarele forme:

```
x = icdf('legea', P, p1, p2, ...)  
x = numef(P, p1, p2, ...)
```

unde `legea` este un șir de caractere predefinit pentru fiecare din legile de probabilitate disponibile în *Statistics toolbox*, `numef` este un șir de caractere din care ultimele trei sunt `inv`, iar cele ce le preced sunt cele care dau numele predefinit al legii de probabilitate (ca și cele din parametrul `legea`).

În urma executării uneia din cele două instrucțiuni, se calculează matricea x a cuantilelor legii precizată prin parametrii `legea`, respectiv `numef`, corespunzătoare valorilor date prin matricea P și având parametrii dați prin matricele `par1`, `par2`,... Aceste matrice trebuie să fie de aceeași dimensiuni, cu excepția că dacă unele sunt scalari, aceștia se extind la matricele constante de aceeași dimensiuni cu celelalte și care iau valorile scalarilor corespunzători.

Programul 2.6.22. Pentru a ilustra grafic modul de calcul al cuantilelor, vom da un program care calculează mediana pentru legea uniformă discretă și reprezintă grafic acest lucru pentru două valori distincte ale parametrului N , precum și cuartilele legii normale.

```
clf, clear
N1 = input('N1='); N2 = input('N2=');
m = input('mu='); s = input('sigma=');
xu1 = 0:N1+1; yu1 = unidcdf(xu1,N1);
xu2 = 0:N2+1; yu2 = unidcdf(xu2,N2);
xn = m-3*s:0.01:m+3*s; yn = normcdf(xn,m,s);
me1 = icdf('unid',1/2,N1);
me2 = icdf('unid',1/2,N2);
Q = icdf('norm',[1/4,2/4,3/4],m,s);
subplot(3,1,1), stairs(xu1,yu1)
set(gca,'Xlim',[0,N1+1]), set(gca,'xtick',[me1])
set(gca,'xticklabel',[me1]), hold on
plot([0,me1],[1/2,1/2],'k:',me1,1/2,'o')
plot([me1,me1],[0,1/2],'k-.')
subplot(3,1,2), stairs(xu2,yu2)
set(gca,'Xlim',[0,N2+1]), set(gca,'xtick',[me2])
set(gca,'xticklabel',[me2]), hold on
plot([0,me2],[1/2,1/2],'k:',me2,1/2,'o')
plot([me2,me2],[0,1/2],'k-.')
subplot(3,1,3),
plot(xn,yn,[Q(1),Q(2),Q(3)],[1/4,2/4,3/4],'o')
set(gca,'xtick',[Q(1),Q(2),Q(3)])
set(gca,'xticklabel',[Q(1),Q(2),Q(3)])
hold on
X = [m-3*s,Q(1);m-3*s,Q(2);m-3*s,Q(3)];
plot(X',[1/4,1/4;1/2,1/2;3/4,3/4]','k:')
plot([Q(1),Q(1)],[0,1/4],'k-.')
plot([Q(2),Q(2)],[0,2/4],'k-.')
plot([Q(3),Q(3)],[0,3/4],'k-.')
```

Executarea programului pentru $N=5$ și $N=6$, respectiv pentru $\mu=0$ și $\sigma=2$, are ca rezultat graficele din Figura 2.17.

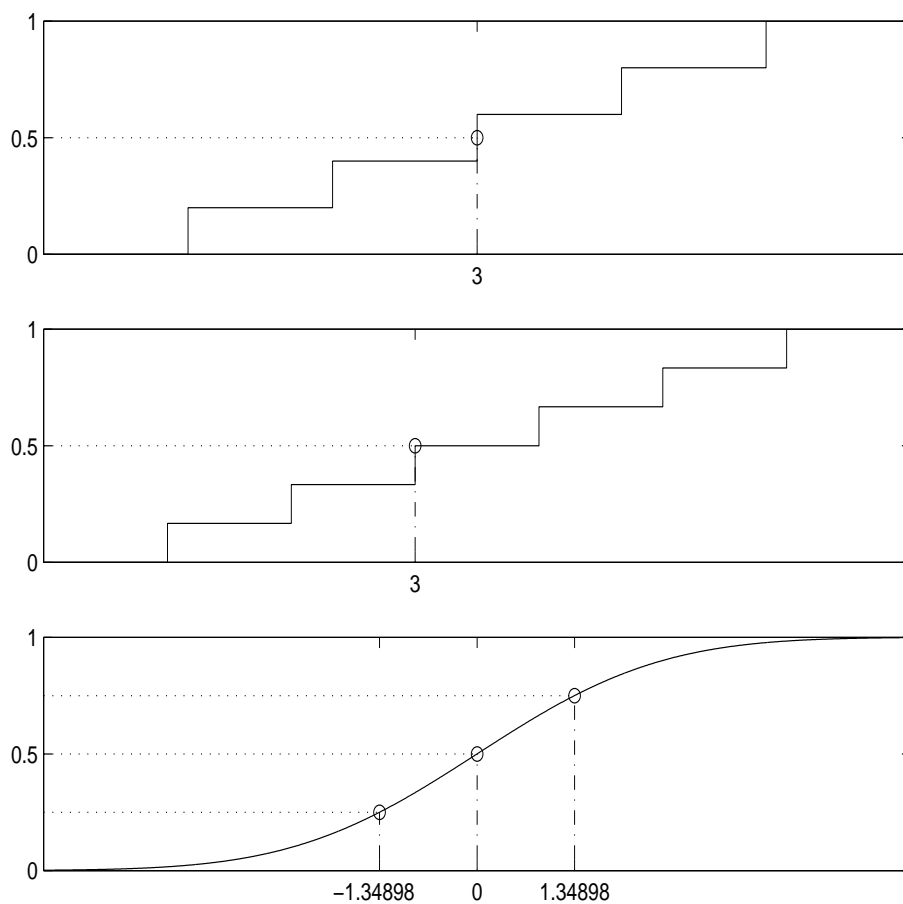


Figura 2.17: Medianele pentru $\mathcal{U}(5)$ și $\mathcal{U}(6)$ și cuartilele pentru $\mathcal{N}(0, 2)$

2.6.8 Funcția Matlab `disttool`

Funcția (comanda) `disttool` este un program de demonstrativ. Lansarea acestuia se face prin:

```
>>disttool
```

în urma căreia se produce o fereastră grafică interactivă (demonstrativă) privind funcțiile de repartiție (`cdf`) și funcțiile de probabilitate (`pdf`).

Fixarea (stabilirea) legii de probabilitate în scop demonstrativ se face prin alegerea din meniul legilor de probabilitate situat în partea stângă sus a ferestrei, iar una din alternativele `pdf` și `cdf` se face folosind meniul din partea dreaptă sus a ferestrei. Pentru stabilirea valorilor parametrilor legii de probabilitate considerate se poate proceda în două moduri. Fie prin introducerea în ferestrele corespunzătoare ale valorilor dorite, fie prin deplasarea barelor atașate acestora. În plus, limite pentru parametrii legii de probabilitate considerate pot fi precizate prin introducerea acestora în ferestrele considerate.

Determinarea valorii funcției `pdf` respectiv `cdf` într-un punct, se poate obține prin introducerea argumentului funcției în fereastra de pe axa absciselor sau prin deplasarea drepte verticale afișată pe grafic, cu ajutorul *mouse*-lui, până când aceasta trece prin punctul respectiv. Inversa funcției de repartiție (`icdf`) se poate obține de asemenea, prin introducerea valorii funcției de repartiție (`cdf`) în fereastra de pe axa ordonatelor sau prin deplasarea drepte orizontale afișată pe grafic, cu ajutorul *mouse*-lui, până când aceasta trece prin punctul respectiv.

2.7 Șiruri de variabile aleatoare

2.7.1 Inegalitatea lui Cebîșev

Fie variabila aleatoare X , pentru care există valoarea medie și dispersia.

Inegalitatea lui Cebîșev este cunoscută sub una din formele

$$P(|X - E(X)| < \varepsilon) \geq 1 - \frac{Var(X)}{\varepsilon^2},$$

$$P(|X - E(X)| \geq \varepsilon) \leq \frac{Var(X)}{\varepsilon^2}, \quad \forall \varepsilon > 0.$$

Observația 2.7.1. Dacă în inegalitatea lui Cebîșev se ia $\varepsilon = k\sigma$, unde s-a notat $\sigma = \sqrt{Var(X)}$ abaterea standard, atunci aceasta devine

$$P(|X - E(X)| < k\sigma) \geq \frac{k^2 - 1}{k^2}.$$

Astfel, pentru $k = 3$ se obține $P(|X - E(X)| < 3\sigma) \geq \frac{8}{9} = 0,88\dots$, ceea ce înseamnă că probabilitatea ca valorile lui X să se abată de la valoarea medie $E(X)$ mai puțin de trei ori abaterea standard este mai mare decât $\frac{8}{9}$.

Observația 2.7.2. Fie variabilele aleatoare X_k , $k = \overline{1, n}$, independente și identic repartizate, având valorile medii și dispersiile $E(X_k) = \mu$, $Var(X_k) = \sigma^2$, pentru fiecare $k = \overline{1, n}$, cu ajutorul cărora construim variabila aleatoare

$$S_n = \sum_{k=1}^n X_k, \quad \text{cu} \quad E(S_n) = n\mu \quad \text{și} \quad Var(S_n) = n\sigma^2.$$

Dacă se aplică inegalitatea lui Cebîșev se obține

$$P\left(\left|\frac{S_n}{n} - \mu\right| < \varepsilon\right) \geq 1 - \frac{\sigma^2}{n\varepsilon^2}, \quad P\left(\left|\frac{S_n}{n} - \mu\right| \geq \varepsilon\right) \leq \frac{\sigma^2}{n\varepsilon^2}, \quad \forall \varepsilon > 0.$$

Cele două inegalități dau evaluări ale probabilităților ca media aritmetică a variabilelor aleatoare să se abată de la valoarea medie μ mai puțin, respectiv mai mult, decât $\varepsilon > 0$ fixat.

2.7.2 Tipuri de convergență

Să considerăm șirul de variabile aleatoare $(X_n)_{n \geq 1}$ cu șirul funcțiilor de repartiție corespunzătoare $(F_n)_{n \geq 1}$.

Definiția 2.7.3. Spunem că șirul de variabile aleatoare $(X_n)_{n \geq 1}$ converge în probabilitate la variabila aleatoare X , și vom nota $X_n \xrightarrow{p} X$, dacă pentru orice $\varepsilon > 0$, avem că

$$\lim_{n \rightarrow \infty} P(|X_n - X| < \varepsilon) = 1.$$

Definiția 2.7.4. Spunem că șirul de variabile aleatoare $(X_n)_{n \geq 1}$ converge în repartiție la variabila aleatoare X , și vom nota $X_n \xrightarrow{r} X$, dacă avem că

$$\lim_{n \rightarrow \infty} F_n(x) = F(x),$$

unde F este funcția de repartiție a variabilei aleatoare X , iar limita considerată există pentru orice punct de continuitate $x \in \mathbb{R}$ al funcției F .

Observația 2.7.5. Dacă $X_n \xrightarrow{p} X$, atunci $X_n \xrightarrow{r} X$. Implicația inversă nu are loc, dar dacă X este o constantă a , atunci afirmația inversă are loc de asemenea, adică dacă $X_n \xrightarrow{r} a$, pentru $a \in \mathbb{R}$, atunci $X_n \xrightarrow{p} a$.

2.7.3 Legea numerelor mari

Să considerăm șirul de variabile aleatoare $(X_n)_{n \geq 1}$.

Definiția 2.7.6. Spunem că șirul de variabile aleatoare $(X_n)_{n \geq 1}$ urmează legea (slabă) a numerelor mari, dacă pentru orice $\varepsilon > 0$ avem că

$$\lim_{n \rightarrow \infty} P \left(\left| \frac{1}{n} \sum_{k=1}^n X_k - \frac{1}{n} \sum_{k=1}^n M(X_k) \right| < \varepsilon \right) = 1,$$

adică

$$\frac{1}{n} \sum_{k=1}^n X_k - \frac{1}{n} \sum_{k=1}^n M(X_k) \xrightarrow{p} 0.$$

Observația 2.7.7. Dacă variabilele aleatoare ale șirului $(X_n)_{n \geq 1}$ sunt independente și identic repartizate, iar $E(X_n) = \mu$, atunci

$$\lim_{n \rightarrow \infty} P \left(\left| \frac{S_n}{n} - \mu \right| < \varepsilon \right) = 1, \quad \text{adică} \quad \frac{1}{n} S_n \xrightarrow{p} \mu,$$

unde $S_n = \sum_{k=1}^n X_k$.

Prin urmare, media aritmetică a primelor n variabile aleatoare din șirul considerat își pierde caracterul aleator, pentru $n \rightarrow \infty$.

Evident, deoarece convergența în probabilitate implică convergența în repartiție, avem și $\frac{1}{n} S_n \xrightarrow{r} \mu$, dar vom vedea imediat că există chiar un rezultat mai bun în ceea ce privește convergența în repartiție.

2.7.4 Teoreme limită

Am văzut că legea numerelor mari este în strânsă legătură cu convergența în probabilitate, chiar și cu așa numita convergența tare, dar care nu a fost prezentată în secțiunea precedentă.

Teoremele limită abordează problematica comportării asimptotice a șirurilor de variabile aleatoare din perspectiva convergenței în repartiție.

Teorema 2.7.8 (Lindeberg–Levy). Fie șirul de variabile aleatoare independente și identic repartizate $(X_n)_{n \geq 1}$, pentru care se introduc notațiile

$$E(X_n) = \mu, \quad \text{Var}(X_n) = \sigma^2, \quad S_n = \sum_{k=1}^n X_k.$$

Dacă se construiește șirul de variabile aleatoare $(Z_n)_{n \geq 1}$, cu termenul general

$$Z_n = \frac{\frac{1}{n}S_n - \mu}{\frac{\sigma}{\sqrt{n}}},$$

atunci $Z_n \xrightarrow{r} Z$, unde Z este o variabilă aleatoare ce urmează legea normală $\mathcal{N}(0, 1)$, adică

$$\lim_{n \rightarrow \infty} F_n(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{t^2}{2}} dt, \quad \forall x \in \mathbb{R},$$

F_n fiind funcția de repartiție a variabilei aleatoare Z_n .

Observația 2.7.9. Teorema limită de acest tip, în care legea de probabilitate limită este legea normală se numește *teoremă limită centrală*.

Teoremele limită au un rol aparte în statistică, unde sunt privite în următoarea manieră. Dacă n este mare ($n \rightarrow \infty$), atunci variabila aleatoare Z_n poate fi considerată ca urmând o lege de probabilitate cunoscută, în cazul teoremei precedente legea normală $\mathcal{N}(0, 1)$.

Mai observăm că variabila aleatoare $\frac{1}{n}S_n$, când $n \rightarrow \infty$, urmează legea normală $\mathcal{N}\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$.

Teorema 2.7.10 (Moivre–Laplace). Dacă variabila aleatoare Z_n urmează legea binomială, adică are distribuția

$$S_n \left(\begin{matrix} x \\ f(x | n, p) \end{matrix} \right)_{x=0, n}, \quad \text{unde } f(x | n, p) = \binom{n}{x} p^x q^{n-x}, \quad q = 1 - p,$$

atunci

$$\frac{S_n - np}{\sqrt{npq}} \xrightarrow{r} Z,$$

unde Z este o variabilă aleatoare ce urmează legea normală $\mathcal{N}(0, 1)$.

Rezultatul teoremei se exprimă de regulă sub forma

$$(2.7.1) \quad P\left(a \leq \frac{x - np}{\sqrt{npq}} < b\right) \approx \frac{1}{\sqrt{2\pi}} \int_a^b e^{-\frac{t^2}{2}} dt.$$

Observația 2.7.11. Teorema limită Moivre–Laplace este o consecință imediată a Teoremei Lindeberg–Levy, dacă se are în vedere că variabila aleatoare S_n este suma

a n variabile aleatoare independente X_k , $k = \overline{1, n}$, fiecare urmând aceeași lege a lui Bernoulli, adică având distribuția

$$X_k \left(\begin{matrix} x \\ f(x | p) \end{matrix} \right)_{x=1,0}, \quad \text{unde } f(x | p) = p^x q^{1-x}.$$

Mai remarcăm faptul că variabila aleatoare S_n , care urmează legea binomială $\mathcal{B}(n, p)$, poate fi privită, pentru $n \rightarrow \infty$, ca urmând legea normală $\mathcal{N}(np, \sqrt{npq})$.

Observația 2.7.12 (Regula celor trei σ). Dacă în formula Moivre–Laplace se consideră $b = -a = 3$, atunci

$$P \left(-3 \leq \frac{k - np}{\sqrt{npq}} < 3 \right) \approx \Phi(3) - \Phi(-3) = 0.997.$$

Prin urmare, rezultă că $P(|k - np| < 3\sqrt{npq}) \geq 0.99$. Deoarece, abaterea medie pătratică a variabilei aleatoare X ce urmează legea binomială este $\sigma = \sqrt{npq}$, avem regula următoare cunoscută sub denumirea de *regula celor trei σ* : probabilitatea ca frecvența absolută k să se abată de la valoarea medie $E(X) = np$ mai puțin de trei ori abaterea medie pătratică este mai mare decât 0.99.

Observația 2.7.13 (Corecție de continuitate). Având în vedere procesul de aproximare a unei legi de tip discret (legea binomială) cu una de tip continuu (legea normală), pentru obținerea unor rezultate mai bune în aplicarea formulei (2.7.1), se impune aplicarea unei corecții de continuitate, dată prin

$$\begin{aligned} f(x | n, p) &= P(S_n = x) \approx P\left(\frac{x - \frac{1}{2} - np}{\sqrt{npq}} \leq Z < \frac{x + \frac{1}{2} - np}{\sqrt{npq}}\right) \\ &= \Phi\left(\frac{x + \frac{1}{2} - np}{\sqrt{npq}}\right) - \Phi\left(\frac{x - \frac{1}{2} - np}{\sqrt{npq}}\right), \end{aligned}$$

unde

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{t^2}{2}} dt,$$

reprezintă funcția lui Laplace.

Folosind această corecție de continuitate avem următoarele formule de aproxi-

mare:

$$\begin{aligned}
 F_n(x) &= P(S_n \leq x) \approx P\left(Z \leq \frac{x + \frac{1}{2} - np}{\sqrt{npq}}\right) = \Phi\left(\frac{x + \frac{1}{2} - np}{\sqrt{npq}}\right), \\
 1 - F_n(x) &= P(S_n > x) \approx P\left(Z > \frac{x - \frac{1}{2} - np}{\sqrt{npq}}\right) = 1 - \Phi\left(\frac{x - \frac{1}{2} - np}{\sqrt{npq}}\right), \\
 P(k_1 \leq S_n \leq k_2) &\approx P\left(\frac{k_2 + \frac{1}{2} - np}{\sqrt{npq}} \leq Z \leq \frac{k_1 - \frac{1}{2} - np}{\sqrt{npq}}\right) \\
 &= \Phi\left(\frac{k_2 + \frac{1}{2} - np}{\sqrt{npq}}\right) - \Phi\left(\frac{k_1 - \frac{1}{2} - np}{\sqrt{npq}}\right).
 \end{aligned}$$

Programul 2.7.14. Pentru a ilustra rezultatele Teoremei Moivre–Laplace, programul Matlab care urmează, reprezintă grafic prin bare funcția de probabilitate a legii binomiale $\mathcal{B}(n, p)$ împreună cu densitatea de probabilitate a legii normale corespunzătoare, adică $\mathcal{N}(np, \sqrt{npq})$.

```

clf
n = input('n='); p = input('p=');
m = n*p; s = sqrt(n*p*(1-p));
x = 0:n; xx = -1:0.01:n+1;
P = binopdf(x,n,p); y = normpdf(xx,m,s);
bar(x,P,0.3), hold, plot(xx,y)

```

Executarea programului pentru $n=15$ și $p=0.4$, produce graficul din Figura 2.18.

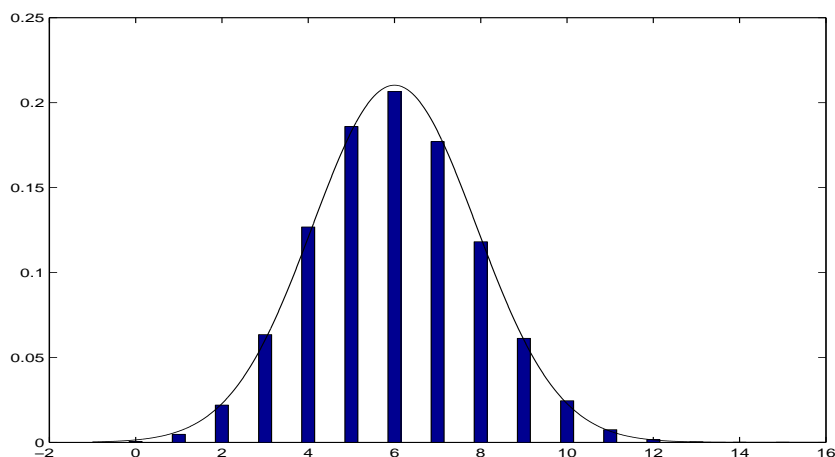


Figura 2.18: Legea $\mathcal{B}(15, 0.4)$ și legea $\mathcal{N}(6, \sqrt{3.6})$

Capitolul 3

Statistică descriptivă

Statistica se ocupă cu descrierea și analiza numerică a fenomenelor de masă, dezvăluind particularitățile lor de volum, structură, dinamică, conexiune, precum și regularitățile sau legile care le guvernează.

Etapile principale ce se disting în cercetarea statistică a unui fenomen aleator se pot considera ca fiind:

a). *Definirea (concepția)* obiectului studiat, care conține definirea unităților statistice, conceperea chestionarului (întrebărilor), planificarea culegerii datelor.

b). *Observarea fenomenului (culegerea datelor)* conform criteriilor stabilite la etapa precedentă.

c). *Descrierea statistică*, ce cuprinde reprezentarea grafică a datelor statistice, sistematizarea acestora, precum și calcularea indicatorilor numerici pentru punerea în evidență a unor proprietăți și pentru sugerarea unor ipoteze referitoare la legile care guvernează fenomenul cercetat.

d). *Modelarea probabilistică* a fenomenului cercetat, care are ca obiectiv principal cercetarea fenomenului folosind ca instrument de lucru teoria probabilităților relativ la datele statistice obținute conform etapelor precedente.

Dacă primele două etape pot fi considerate ca preliminare, etapa a treia conduce la ceea ce numim *statistică descriptivă*, iar a patra etapă conduce la *statistica matematică* sau *statistica inferențială*.

Desigur, nu se poate face o delimitare exactă între statistica descriptivă și statistica matematică, mai mult o descriere statistică fundamentată matematic va conduce la realizarea unei modelări a fenomenului, care pune în evidență proprietățile esențiale ale acestuia și legăturile dintre acestea. Ca o componentă de bază a statisticii descriptive, care s-a dezvoltat în mod impresionant odată cu explozia utilizării calculatoarelor în statistică, o reprezintă *analiza datelor*. Printre capitolele mai importante ale analizei datelor amintim *metodele de clasificare* și *metodele factoriale*.

Pe de altă parte, statistica matematică cuprinde ca și capitole principale *teoria selecției, teoria estimăției și verificarea ipotezelor statistice*, strâns legat de *teoria deciziilor*.

Din cele prezentate putem spune că *statistica descriptivă* are ca obiectiv principal culegerea și clasificarea datelor statistice în vederea descrierii numerice și grafice a fenomenului cercetat. *Statistica matematică* elaborează modelul matematic, folosind datele statistice, estimează parametrii modelului, stabilește metode pentru validarea ipotezelor statistice, toate acestea având la baza teoria probabilităților.

Amintim că studiul statistic al fenomenelor economice face obiectul *statisticii economice*, iar studiul statistic al fenomenelor ce pot constitui riscuri asigurate reprezintă obiectul *statisticii actuariale*.

3.1 Concepte de bază ale statisticii

Definiția 3.1.1. Numim colectivitate (populație) o mulțime C de elemente cercetată din punct de vedere a uneia sau mai multor proprietăți, elementele componente numindu-se indivizi sau unități statistice, iar numărul elementelor lui C se numește volumul colectivității.

Definiția 3.1.2. Numim caracteristică sau variabilă a colectivității C proprietatea supusă investigării statistice relativă la C . Când o caracteristică poate fi măsurată o numim caracteristică cantitativă sau numerică, iar dacă aceasta se exprimă printr-o însușire o numim caracteristică calitativă sau atribut.

În unele tratate de statistică variabilele (caracteristicile) calitative sunt numite *variabile nominale*. O clasă intermediară de variabile (caracteristici) este formată de așa numitele *variabile (caracteristici) ordinale*.

Observația 3.1.3. Dacă se consideră colectivitatea C a studenților dintr-un an de studiu, atunci înălțimea, greutatea, media obținută la sfârșitul anului sunt caracteristici cantitative (numerice), pe când sexul, culoarea ochilor, grupa sanguină sunt caracteristici calitative (nominale). Un exemplu de caracteristică ordinală ar fi gradul de pregătire (foarte bun, bun, satisfăcător, nesatisfăcător). Se observă că relativ la o variabilă ordinală se poate face ordonarea indivizilor, ca și pentru o variabilă numerică, ceea ce nu se întâmplă pentru o variabilă nominală, dar nu se pot face operații aritmetice, cum se întâmplă și la variabilele nominale, pe când pentru o variabilă numerică astfel de operații sunt posibile.

Observația 3.1.4. Din punct de vedere al teoriei probabilităților o caracteristică a unei populații C este o variabilă aleatoare X . Scopul principal al cercetării statistice

este de a stabili legea de probabilitate pe care o urmează caracteristica X , utilizând observațiile (datele statistice) relative la colectivitatea cercetată.

Dacă avem în vedere mai multe caracteristici X_1, \dots, X_n , acestea reprezintă componentele unui vector aleator \mathbf{X} , iar cercetarea fenomenului aleator modelat prin \mathbf{X} , se poate extinde și la determinarea legăturilor ce există între componentele vectorului aleator.

Definiția 3.1.5. *O caracteristică X ce ia o mulțime cel mult numărabilă de valori o numim caracteristică de tip discret, iar dacă ia valori dintr-un interval finit sau infinit o numim caracteristică de tip continuu.*

Observația 3.1.6. Dacă se consideră colectivitatea \mathcal{C} a bolnavilor externați dintr-un spital pe parcursul unei săptămâni, atunci caracteristica X ce reprezintă numărul zilelor de internare este o caracteristică de tip discret, pe când greutatea Y a bolnavilor externați reprezintă o caracteristică de tip continuu. Remarcăm faptul că pentru caracteristica Y caracterul continuu este estompat prin imposibilitatea măsurării exacte a acestei caracteristici, de aceea, practic, se consideră ca fiind de asemenea de tip discret.

3.2 Culegerea, prezentarea și prelucrarea datelor statistice

Modurile de culegere (observare) a datelor statistice cele mai des utilizate sunt:

a). *Observarea totală (exhaustivă, recensământ)*, când toți indivizii colectivității \mathcal{C} sunt înregistrați.

b). *Observare parțială (sondaj, selecție)*, când după criterii bine stabilite sunt înregistrați o parte din indivizii colectivității \mathcal{C} , numită *eșantion* sau *selecție*.

c). *Observare curentă*, când înregistrarea indivizilor se efectuează odată cu apariția (producerea) lor.

d). *Observarea periodică*, când înregistrarea fenomenului se efectuează la intervale de timp stabilite.

3.2.1 Generarea numerelor aleatoare în Matlab

Matematica, și prin teoria probabilităților și statistica matematică, printre multiplele și larg răspânditele aplicații, posedă, poate, una legată de *modelarea matematică*, care se distinge în mod special prin vastul domeniu de aplicabilitate. Modelarea matematică, după cum spune și numele, are drept scop construirea de *modele matematice* asociate unor fenomene sau sisteme pe care le întâlnim în diferite domenii: fizică, chimie, biologie, demografie, medicină, geografie, geologie, economie etc. Desigur,

astfel de modele matematice depind de niște parametri, fie aleatori, fie nealeatori. Odată modelul matematic construit se are în vedere realizarea practică a acestuia. Dar, apare o mare problemă, este acest model matematic (teoretic) construit viabil? Avem două alternative, sau realizăm practic modelul respectiv și urmărim cum funcționează acesta (ceea ce, în general, nu e de dorit), sau simulăm funcționarea modelului respectiv prin diferite particularizări ale parametrilor modelului și alegem varianta convenabilă. Dacă relativ la parametrii nealeatori care intervin în modelul teoretic nu se ridică probleme speciale, în privința parametrilor aleatori apare o nouă problemă, anume obținerea valorilor acestor parametri. Deoarece parametrii aleatori sunt, din punct de vedere al teoriei probabilităților, variabile aleatoare, apare așadar problema generării de valori numerice ale acestor variabile aleatoare și care vor purta denumirea de *numere aleatoare*. Desigur că într-un fel vor fi generate aceste numere aleatoare dacă parametrul urmează, să zicem, legea uniformă pe un anumit interval și altfel când urmează legea normală.

Deoarece pentru prezentarea noțiunilor și rezultatelor teoretice ale statisticii matematice sunt necesare exemplificări ale acestora, iar culegerea unor date statistice reale nu este întodeauna la îndemână, de multe ori vom considera numerele aleatoare ca fiind date statistice cărora li se vor aplica metodele cercetării statisticii.

Putem să amintim de asemenea că testarea programelor pe calculator, determinarea statistică a caracteristicilor numerice ale variabilelor aleatoare, care ne conduce la metoda Monte Carlo, fac de asemenea necesară generarea numerelor aleatoare.

Metodele de generare a numerelor aleatoare se împart în trei categorii:

- tabelele cu numere aleatoare (uniforme) obținute, de exemplu, din aruncarea monedei, aruncarea zarului, jocul de ruletă, tabele matematice, tabele de recensământ etc.
- procedeele fizice, care au la baza fenomene fizice, cum ar fi emiterea particulelor de o sursă radioactivă, intensitatea curentului la diferite momente, zgomotul electronic etc.
- procedeele aritmetice (analitice), care utilizează o formulă de calcul de forma

$$x_{n+1} = f(x_n, x_{n-1}, \dots, x_{n-m}), \quad n \geq m \geq 0.$$

Deoarece, de regulă, se folosește un volum mare de numere aleatoare, s-a impus generarea acestora cu ajutorul procedeelor aritmetice. Inconvenientul acestor procedee analitice este legat de faptul că nu posedă caracterul strict aleator, din moment ce există o legătură funcțională între acestea. Dar există posibilitatea de a alege astfel de procedee analitice care produc șiruri de *numere pseudoaleatoare*, foarte apropiate din punct de vedere statistic de numerele aleatoare propriu-zise.

În lucrările [3] și [5] sunt prezentate proceduri pentru generarea de numere aleatoare pentru diferite legi de probabilitate.

Sistemul Matlab are implementate astfel de proceduri pentru generarea numerelor aleatoare, fie prin *Statistics toolbox*, fie prin sistemul de bază din Matlab.

Funcția `random`

Instrucțiunile prin care sunt generate numere aleatoare ce urmează una din legile de probabilitate recunoscute de sistemul Matlab prin *Statistics toolbox* sunt de forma:

```
x=random('legea',par1,par2,...,m,n)
x=numef(par1,par2,...,m,n)
```

unde `legea` este un șir de caractere predefinit pentru fiecare din legile de probabilitate disponibile în *Statistics toolbox*, `numef` este un șir de caractere din care ultimele trei sunt `rnd`, iar cele ce le preced sunt cele care dau numele predefinit al legii de probabilitate (ca și cele din parametrul `legea`).

Executarea uneia din cele două forme ale instrucțiunii `random`, are ca efect generarea matricei `x`, cu `m` linii și `n` coloane, de numere aleatoare ce urmează legea de probabilitate corespunzătoare cu parametrii precizați prin matricele `p1`, `p2`, ... Aceste matrice trebuie să aibă aceleași dimensiuni cu matricea `x`, cu excepția când unele sunt scalari, caz în care acestea sunt extinse la matricele constante de aceleași dimensiuni cu celelalte și care iau valorile corespunzătoare scalarilor. Dacă toți parametrii legii de probabilitate sunt scalari, când lipsesc parametrii `m` și `n` se generează un număr aleator, când lipsește `n`, atunci trebuie ca `m` să fie un vector cu două componente, care conține dimensiunile matricei `x`. Dacă cel puțin unul din parametrii legii de probabilitate nu este scalar și dacă parametrii `m` și `n` sunt prezenți, aceștia trebuie să coincidă cu dimensiunile matricelor.

Funcțiile `mvnrnd` și `mvtrnd`

Funcțiile `mvnrnd` și `mvtrnd` sunt singurele funcții din *Statistics toolbox*, prin care se generează legi de probabilitate multidimensionale, respectiv legile normală și *t* (Student).

Apelurile acestor funcții se pot face respectiv prin:

```
x=mvnrnd(mu,v,m)
x=mvtrnd(r,n,m)
```

În primul caz, sunt generați `m` vectori aleatori ce urmează legea normală multidimensională, având vectorul valorilor medii `mu` și matricea covarianțelor `v` (matrice pătratică pozitiv definită). Cei `m` vectori aleatori generați sunt obținuți în matricea `x`, care în consecință va avea `m` linii, iar numărul coloanelor este dat de dimensiunea legii de probabilitate și care trebuie să coincidă cu lungimea vectorului `mu` și cu ordinul matricei pătratice `v`.

În al doilea caz, vor fi generați m vectori aleatori ce urmează legea Student multidimensională.

Parametrul r reprezintă matricea coeficienților de corelație, care este matrice pozitiv definită având pe diagonala principală toate elementele 1. Dacă această ultimă condiție nu este îndeplinită, adică dacă este matricea covarianțelor, atunci sistemul Matlab calculează prima dată, din matricea r , matricea coeficienților de corelație.

Parametrul n reprezintă numărul gradelor de libertate, putând fi un scalar sau un vector de lungime m .

Cei m vectori aleatori sunt obținuți în cele m linii ale matricei x și care are atâtea coloane cât este dimensiunea legii de probabilitate.

Trebuie să remarcăm faptul că o linie a matricei x se obține dintr-un vector aleator ce urmează legea normală multidimensională, având valoarea medie 0 și matricea covarianțelor dată prin parametrul r , ale cărui componente se împart la un număr aleator ce urmează legea χ^2 cu n grade de libertate.

Funcțiile **rand** și **randn**

Sistemul Matlab de bază conține două funcții pentru generarea de numere aleatoare ce urmează legea uniformă $\mathcal{U}(0, 1)$, prin **rand**, respectiv legea normală standard $\mathcal{N}(0, 1)$, prin **randn**.

Apelul acestor funcții se poate face prin instrucțiuni de forma:

```
x=rand
x=rand(m)
x=rand(m,n)
x=randn
x=randn(m)
x=randn(m,n)
```

În urma executării unor astfel de instrucțiuni este generat un număr aleator ce urmează legea $\mathcal{U}(0, 1)$ respectiv $\mathcal{N}(0, 1)$, când nu este folosit nici un argument, o matrice pătratică de ordin m de numere aleatoare, când apare un argument și o matrice de numere aleatoare de tipul (m, n) , când sunt folosite ambele argumente.

Unele situații impun regăsirea unui anumit șir de numere aleatoare generat prin una din cele două funcții **rand** și respectiv **randn**. Pentru aceasta avem la dispoziție ceea ce se numește *starea generatorului*.

Obținerea stării generatorului **rand** și **randn** la un moment dat se poate face prin:

```
s=rand('state')
s=randn('state')
```

Pentru reinițializarea sau schimbarea stării generatorului, se pot utiliza următoarele instrucțiuni:


```

rand('state',s)
randn('state',s)
rand('state',0)
randn('state',0)
rand('state',sum(100*clock))
randn('state',sum(100*clock))

```

Parametrul s precizează starea la care urmează să fie resetat generatorul, dacă acesta este 0, generatorul va fi resetat la starea inițială. Dacă apare parametrul $\text{sum}(100*\text{clock})$, generatorul este resetat la fiecare apel în funcție de timpul calculatorului.

Instrucțiuni moștenite din Matlab 4 și Matlab 5 sunt de asemenea active:

```

rand('seed')
randn('seed')
rand('seed',0)
randn('seed',0)
rand('seed',r)
randn('seed',r)

```

Funcția **randperm**

Generarea unei permutări a primelor n numere naturale se realizează prin instrucțiunea

```
perm=randperm(n)
```

și va avea ca efect obținerea vectorului perm ce conține o astfel de permutare.

3.2.2 Tabele statistice

Definiția 3.2.1. Numim tabel statistic (simplu, nesistematizat) un tablou în care înregistrările sunt trecute în ordinea apariției lor.

Observația 3.2.2. Dacă observațiile statistice s-au făcut relativ la p caracteristici, fie acestea X_1, X_2, \dots, X_p , asupra a n indivizi, atunci tabelul statistic simplu este matricea $\mathbf{X} = (x_{ij})_{i=\overline{1,n}, j=\overline{1,p}}$, unde x_{ij} este valoarea caracteristicii X_j pentru individul i .

Definiția 3.2.3. Numim tabel statistic (sistematizat) relativ la caracteristica X de tip discret, tabloul în care sunt trecute valorile distincte ale caracteristicii și frecvențele cu care au apărut aceste valori.

Observația 3.2.4. Având caracteristica X de tip discret, care ia valorile distincte $x_i, i = \overline{1,n}$, pentru care s-au obținut datele primare x'_1, x'_2, \dots, x'_N , tabelul statistic sistematizat are forma

x	f
x_1	f_1
x_2	f_2
\vdots	\vdots
x_n	f_n

unde f_i este *frecvența absolută* a apariției valorii x_i în datele primare $x'_k, k = \overline{1, N}$. Prin urmare are loc relația $\sum_{i=1}^n f_i = N$.

Definiția 3.2.5. Fie caracteristica de tip continuu X , care ia valori din intervalul (a, b) , descompus în intervale disjuncte, numite clase, prin punctele care satisfac relațiile

$$a = a_0 < a_1 < \dots < a_n = b.$$

Numim tabel statistic (sistematizat) relativ la caracteristica X , tabloul ce conține clasele caracteristicii și frecvențele cu care au apărut aceste clase.

Observația 3.2.6. Dacă datele primare ale caracteristicii continue X , care ia valori în intervalul (a, b) , sunt x'_1, x'_2, \dots, x'_N , atunci tabelul statistic sistematizat are forma

x	f		x	f
$[a_0, a_1)$	f_1	sau	x_1	f_1
$[a_1, a_2)$	f_2		x_2	f_2
\vdots	\vdots		\vdots	\vdots
$[a_{n-1}, a_n)$	f_n		x_n	f_n

unde f_i este *frecvența absolută* a apariției clasei $[a_{i-1}, a_i)$ printre datele primare $x'_k, k = \overline{1, N}$, iar $x_i = \frac{a_{i-1} + a_i}{2}$. Dacă se consideră varianta din dreapta pentru tabelul sistematizat, se obține aceeași formă de tabel ca și în cazul caracteristicii discrete. În acest fel s-a ajuns la o unificare a prezentării prin tabele sistematizate a tipurilor discrete și continue de caracteristici.

Observația 3.2.7. Prin stabilirea claselor unei caracteristici continue (care pot fi de lungimi egale sau nu) s-a efectuat o *grupare* a datelor statistice primare, și care se efectuează și în cazul caracteristicilor de tip discret, dacă au un număr mare de valori distincte.

Definiția 3.2.8. Numim amplitudinea clasei definită de intervalul $[a_{i-1}, a_i)$ lungimea acestui interval, $d_i = a_i - a_{i-1}$.

Observația 3.2.9. Când amplitudinile claselor sunt egale, două reguli de stabilire a numărului claselor sunt utilizate mai des:

$$n = 1 + \frac{10}{3} \lg N \quad (\text{regula lui Sturges}) \quad \text{și} \quad d = \frac{1}{100} \cdot 8 (x_{\max} - x_{\min}),$$

unde

$$x_{\max} = \max\{x'_1, x'_2, \dots, x'_N\} \quad \text{și} \quad x_{\min} = \min\{x'_1, x'_2, \dots, x'_N\}.$$

În acest fel, din prima regulă, se obține că

$$d = \frac{b - a}{n} \quad \text{și} \quad a_i = a + id, \quad i = \overline{0, n}.$$

Când (a, b) este infinit, atunci

$$d = \frac{x_{\max} - x_{\min}}{n} \quad \text{și} \quad a_i = x_{\min} + id, \quad i = \overline{0, n}.$$

Exemplul 3.2.10. Se consideră un lot de 70 becuri din punct de vedere al caracteristicii X , ce reprezintă durata de viață în mii ore. Datele statistice obținute sunt:

1.318	3.128	2.758	1.583	2.517	2.304	1.155	3.156	2.807	2.879
3.426	1.690	3.537	2.214	2.219	2.072	2.726	1.403	2.493	1.560
3.972	2.637	0.842	2.256	1.708	1.628	2.345	1.855	1.546	3.852
2.128	2.465	2.316	2.262	1.962	1.802	2.230	3.460	1.493	3.093
1.548	2.298	1.875	3.394	1.931	1.179	1.946	1.355	3.006	2.455
1.937	1.977	2.206	1.681	1.960	3.281	2.838	2.525	1.553	2.676
2.500	2.641	1.631	1.864	2.015	2.502	2.444	2.636	2.337	1.966

Vrem să scriem tabelul sistematizat al datelor statistice, considerând clase de amplitudini egale.

Folosind regula lui Sturges avem că numărul n al claselor poate fi calculat cu ajutorul numărului N al datelor statistice prin

$$n = 1 + \frac{10}{3} \lg N = 1 + \frac{10}{3} \times 1.8451,$$

de unde rezultă că $n = 7$. În acest fel avem că amplitudinea d a claselor este

$$d = \frac{x_{\max} - x_{\min}}{n} = \frac{3.972 - 0.842}{7} = 0.45.$$

Când se folosește formula

$$d = \frac{1}{100} \cdot 8 (x_{\max} - x_{\min})$$

pentru calculul amplitudinii claselor, se obține $d = 0.25$. În acest caz numărul n al claselor va fi

$$n = \frac{x_{\max} - x_{\min}}{d} = \frac{3.972 - 0.842}{0.25},$$

de unde $n = 12$.

Se observă că folosind cele două metode s-au obținut valori distincte pentru numărul claselor. Aceste formule au mai mult rolul de a da o primă informație relativă la numărul claselor.

În continuare vom considera numărul claselor $n = 11$, $a_i = 0.7 + i \times d$, $i = \overline{0, 11}$, și $d = 0.3$.

Folosind această grupare se obține tabelul sistematizat

x	f		x	f
[0.7; 1.0)	1		0.85	1
[1.0; 1.3)	2		1.15	2
[1.3; 1.6)	9		1.45	9
[1.6; 1.9)	9		1.75	9
[1.9; 2.2)	10		2.05	10
[2.2; 2.5)	15	sau	2.35	15
[2.5; 2.8)	10		2.65	10
[2.8; 3.1)	5		2.95	5
[3.1; 3.4)	4		3.25	4
[3.4; 3.7)	3		3.55	3
[3.7; 4.0)	2		3.85	2

Programul 3.2.11. În urma executării următorului program Matlab

```
load x.dat -ascii
N=length(x);
ns=fix(1+10/3*log10(N));
ds=(max(x)-min(x))/ns;
fprintf('    ns= %2d,    ds= %4.2f\n',ns,ds)
da=8*(max(x)-min(x))/100;
na=fix((max(x)-min(x))/da);
a=0.7:0.3:4; n=length(a)-1;
fprintf('    na= %2d,    da= %4.2f\n',na,da)
for i=1:N
    for j=1:n
        if (a(j)<=x(i))&(x(i)<a(j+1))
            cl(i)=j;
        end
    end
end
t=tabulate(cl);
fprintf('  Tabelul sistematizat\n')
for j=1:n
    t(j,1)=(a(j)+a(j+1))/2;
    fprintf('%10.2f | %2d\n',t(j,1),t(j,2))
end
```

se obțin rezultatele:

```
ns= 7,    ds= 0.45
na= 12,   da= 0.25
Tabelul sistematizat
```

0.85		1
1.15		2
1.45		9
1.75		9
2.05		10
2.35		15
2.65		10
2.95		5
3.25		4
3.55		3
3.85		2

Se observă că prima comandă citește datele statistice din fișierul `x.dat`, care este un fișier `ascii`.

Definiția 3.2.12. Numim frecvență relativă a clasei x_i raportul $p_i = \frac{f_i}{N}$.

Definiția 3.2.13. Numim frecvențe cumulate ascendente, respectiv frecvențe cumulate descendente *frecvențele date de relațiile*

$$F_k = \sum_{i=1}^k f_i, \quad F'_k = \sum_{i=k+1}^n f_i, \quad k = \overline{0, n},$$

unde $F_0 = 0$ și $F'_n = 0$.

Observația 3.2.14. Pentru frecvențele relative are loc relația

$$\sum_{i=1}^n p_i = 1,$$

iar pentru cele cumulate au loc relațiile $F_k + F'_k = N$, $F_n = N$ și $F'_0 = N$.

Definiția 3.2.15. Numim distribuție statistică a caracteristicii X tabloul de forma

$$X \begin{pmatrix} x_i \\ f_i \end{pmatrix}_{i=\overline{1, n}} \quad \text{sau} \quad X \begin{pmatrix} x_i \\ p_i \end{pmatrix}_{i=\overline{1, n}},$$

unde x_i , $i = \overline{1, n}$, sunt clasele considerate, iar f_i și p_i , $i = \overline{1, n}$, sunt respectiv frecvențele absolute și frecvențele relative.

Exemplul 3.2.16. Considerăm datele statistice de la Exemplul 3.2.10. Atunci distribuția statistică a caracteristicii X poate fi scrisă, fie cu ajutorul frecvențelor absolute, fie cu ajutorul frecvențelor relative. Astfel avem că

$$X \begin{pmatrix} 0.85 & 1.15 & 1.45 & 1.75 & 2.05 & 2.35 & 2.65 & 2.95 & 3.25 & 3.55 & 3.85 \\ 1 & 2 & 9 & 9 & 10 & 15 & 10 & 5 & 4 & 3 & 2 \end{pmatrix}$$

$X \setminus Y$	y_1	y_2	\dots	y_n	
x_1	f_{11}	f_{12}	\dots	f_{1n}	$f_{1.}$
x_2	f_{21}	f_{22}	\dots	f_{2n}	$f_{2.}$
\vdots	\vdots	\vdots		\vdots	\vdots
x_m	f_{m1}	f_{m2}	\dots	f_{mn}	$f_{m.}$
	$f_{.1}$	$f_{.2}$	\dots	$f_{.n}$	$f_{..} = N$

Tabelul 3.1: Tabel de contingență

sau

$$X \left(\begin{array}{cccccccccccc} 0.85 & 1.15 & 1.45 & 1.75 & 2.05 & 2.35 & 2.65 & 2.95 & 3.25 & 3.55 & 3.85 \\ \frac{1}{70} & \frac{2}{70} & \frac{9}{70} & \frac{9}{70} & \frac{10}{70} & \frac{15}{70} & \frac{10}{70} & \frac{5}{70} & \frac{4}{70} & \frac{3}{70} & \frac{2}{70} \end{array} \right)$$

Definiția 3.2.17. Fie colectivitatea C relativ la care sunt cercetate două caracteristici X și Y . Numim tabel de contingență, un tablou care conține clasele caracteristicilor X și respectiv Y , împreună cu frecvențele absolute ale acestor clase.

Observația 3.2.18. Dacă pentru cele două caracteristici X și Y avem respectiv clasele date prin $x_i, i = \overline{1, m}$, și $y_j, j = \overline{1, n}$, iar datele primare sunt date prin perechile $(x'_1, y'_1), (x'_2, y'_2), \dots, (x'_N, y'_N)$, atunci tabelul de contingență este Tabelul 3.1 unde f_{ij} este *frecvența absolută* a apariției clasei (x_i, y_j) în datele primare (x'_k, y'_k) , $k = \overline{1, N}$. După cum se vede tabelul a fost completat cu o linie și o coloană ale căror elemente sunt date prin formulele

$$f_{.j} = \sum_{i=1}^m f_{ij}, \quad j = \overline{1, n}, \quad f_{i.} = \sum_{j=1}^n f_{ij}, \quad i = \overline{1, m},$$

$$f_{..} = \sum_{j=1}^n f_{.j} = \sum_{i=1}^m f_{i.} = \sum_{i=1}^m \sum_{j=1}^n f_{ij} = N.$$

Observația 3.2.19. Când caracteristicile X și Y sunt caracteristici cantitative și între ele există o relație de dependență, tabelul de contingență se numește *tabel de corelație*.

Exemplul 3.2.20. Un astfel de tabel de corelație este prezentat pentru datele statistice ce reprezintă 425 de copii de 10 ani cercetați din punct de vedere al înălțimii X (în centimetri) și al greutății Y (în kilograme). Datele au fost trecute în tabelul de corelație care urmează

$X \backslash Y$	24	25	26	27	28	29	30	31	32	33	34	35	$f_{i.}$
126	1	1	3	1									6
127		3	4	1		1							9
128	1	4	5	7	1	2	1						21
129		1	2	6	9	6	4	1					29
130			1	7	36	17	6	2	1				70
131			1	6	27	56	39	18	4	1			152
132				2	10	16	26	7	2	1			64
133					1	7	10	11	6	2	1		38
134						1	2	4	7	3	1		18
135							1	2	3	1	1		8
136									1	2	2	1	6
137										1	2	1	4
$f_{.j}$	2	9	16	30	84	106	89	45	24	11	7	2	425

3.2.3 Funcțiile tabulate și crosstab

Obținerea de tabele sistematizate, în Matlab, este posibilă cu ajutorul acestor două funcții din *Statistics toolbox*.

Funcția tabulate

Apelul funcției `tabulate` se poate face prin:

```
tabulate(x)
t=tabulate(x)
```

unde x este un vector cu valori pozitive. Prin executarea unei astfel de instrucțiuni, se obține pe ecran respectiv în t un tabel statistic sistematizat având trei coloane. Prima coloană conține valorile distincte ale lui x , a doua conține frecvențele absolute ale acestor valori distincte, iar ultima reprezintă frecvențele relative în procente.

Programul 3.2.21. Vom scrie un program Matlab, care generează N numere aleatoare ce urmează legea uniformă $\mathcal{U}(a, b)$, după care construiește tabelul sistematizat, considerând numărul claselor n dat prin regula lui Sturges, iar clasele fiind identificate prin mijloacele lor.

```
clear
a=input('a:'); b=input('b (a<b):');
N=input('N='); n=fix(1+10/3*log10(N));
xx=unifrnd(a,b,1,N); c=a:(b-a)/n:b;
for i=1:N
    for j=1:n
        if c(j)<=xx(i) & xx(i)<c(j+1)
            x(i)=j;
        end
    end
end
```

```

t=tabulate(x);
for j=1:n
    t(j,1) =(c(j)+c(j+1))/2;
end
fprintf(' Clasa      f      frel\n')
fprintf(' _____\n')
fprintf('%7.3f %5d %8.2f%%\n',t')

```

Executarea programului, cu datele de intrare $a=0$, $b=5$ și $N=200$, produce tabel sistematizat:

Clasa	f	frel
0.313	22	11.00%
0.938	27	13.50%
1.563	27	13.50%
2.188	32	16.00%
2.813	27	13.50%
3.438	23	11.50%
4.063	25	12.50%
4.688	17	8.50%

Funcția crosstab

Apelul funcției crosstab se poate face prin:

```

crosstab(x,y,z,...)
t=crosstab(x,y,z,...)

```

unde x, y, z, \dots sunt vectori cu valori întregi pozitive având aceleași lungimi. Prin executarea unei astfel de instrucțiuni, se obține pe ecran respectiv în t un tabel statistic sistematizat t , cu atâtea intrări câți parametri (vectori) există. Dacă sunt doi parametri, x și y , se obține tabelul de contingență.

Programul 3.2.22. Programul Matlab care urmează, generează N vectori aleatori ce urmează legea normală bidimensională, după care construiește tabelul de corelație cu m clase în raport cu prima variabilă, respectiv n clase în raport cu a doua variabilă.

```

clear, mu(1)=input('m1='); mu(2)=input('m2=');
v(1,1)=input('sigma1^2=');
v(2,2)=input('sigma2^2=');
v(1,2)=input('Cov(X,Y)='); v(2,1) =v(1,2);
N=input('N='); m=input('m='); n=input('n=');
if det(v) <= 0
    error('Matricea v nu e pozitiv definita!')
end
X=mvnrnd(mu,v,N); xx=X(:,1); yy=X(:,2);
xmin=min(xx); xmax=max(xx);
ymin=min(yy); ymax=max(yy);
cx =xmin:(xmax-xmin)/m:xmax;
cy =ymin:(ymax-ymin)/n:ymax;

```



```

for k=1:N
    for i=1:m-1
        if (cx(i)<=xx(k)) & (xx(k)<cx(i+1))
            x(k)=i;
        end
    end
    if (cx(m)<=xx(k)) & (xx(k)<=cx(m+1))
        x(k)=m;
    end
    for j=1:n-1
        if (cy(j)<=yy(k)) & (yy(k)<cy(j+1))
            y(k)=j;
        end
    end
    if (cy(n)<=yy(k)) & (yy(k)<=cy(n+1))
        y(k)=n;
    end
end
t=crosstab(x,y);
disp(t)

```

Pentru $N=100$, $\mu=(5,10)$, $v = \begin{pmatrix} 2 & 1 \\ 1 & 3 \end{pmatrix}$, $m=6$ și $n=4$, se obține tabelul de corelație:

1	2	0	0
2	7	0	0
5	14	6	1
2	17	8	1
1	12	12	3
1	1	1	3

3.2.4 Funcțiile caseread, casewrite, tblread și tblwrite

Pentru completarea tabelului de contingență construit prin funcția `crosstab`, se poate face apel la aceste funcții.

Funcția casewrite

Apelul funcției `casewrite` are sintaxa

```
casewrite(obs,'file')
```

și are ca efect scrierea datelor din matricea `obs` de tip caracter în fișierul cu numele `file`.

Funcția caseread

Apelul funcției `caseread` are sintaxa

```
obs=caseread('file')
```

și are ca efect citirea datelor din fișierul cu numele `file` în matricea `obs` de tip caracter.

Funcția `tblwrite`

Sintaxa funcției `tblwrite` este:

```
tblwrite(t,va,obs,'file')
```

și are ca efect scrierea în fișierul cu numele `file` a matricelor `va`, `obs` de tip caracter și a datelor din matricea numerică `t`.

Matricele `va` și `obs` trebuie să aibă atâtea linii câte coloane, respectiv câte linii are matricea `t`. Liniile celor două matrice de tip caracter vor reprezenta etichete pentru coloanele, respectiv liniile matricei `t`, și drept urmare vor fi așezate în dreptul coloanelor și liniilor corespunzătoare.

Funcția `tblread`

Sintaxa funcției `tblread` este:

```
[t,va,obs]=tblread('file')
```

și are ca efect citirea din fișierul cu numele `file` a matricelor `va`, `obs` de tip caracter și a datelor din matricea numerică `t`.

Programul 3.2.23. Programul care urmează completează Programul 3.2.22, pentru obținerea în tabelul de contingență și a mijloacelor claselor pentru cele două variabile.

```
. . . . .
t=crosstab(x,y);
fy=fopen('ly.txt','w+');
for j=1:n
    laby(j)=(cy(j)+cy(j+1))/2;
    ly=num2str(laby(j));
    fprintf(fy,'%s\n',ly);
end
fx=fopen('lx.txt','w+');
for i=1:m
    labx(i)=(cx(i)+cx(i+1))/2;
    lx=num2str(labx(i));
    fprintf(fx,'%s\n',lx);
end
va=caseread('ly.txt'); obs=caseread('lx.txt');
tblwrite(t,va,obs,'clase.txt');
type clase.txt, fclose('all');
```

În urma executării programului cu aceleași date de intrare folosite în Programul 3.2.22, prin instrucțiune `type` este afișat pe ecran tabelul de contingență:

	7.5244	9.7152	11.906	14.0967
2.5895	4	6	2	0
3.8424	4	16	7	0
5.0952	1	13	12	0
6.348	2	12	9	1
7.6008	1	1	4	3
8.8537	0	0	1	1

Se observă că datele din tabelul obținut nu sunt aceleași cu cele obținute prin executarea Programului 3.2.22, deși se dau aceleași date de intrare. Lucru de așteptat, având în vedere că se generează un alt set de numere aleatoare.

Se observă de asemenea, în noul program, că apar instrucțiunile `fopen`, `fclose`, `caseread`, `tblwrite`.

Instrucțiunea `fopen` este folosită pentru deschiderea fișierelor `lx.txt` și `ly.txt`, care vor conține etichetele claselor celor două variabile. Acestea sunt date de mijloacele claselor și sunt transformate în șiruri de caractere prin comanda `num2str`. Parametrul `'w+'`, are ca efect crearea unui fișier nou sau actualizarea unui deja existent, după ce a fost anulat conținutul său, în vederea scrierii sau citirii în respectivul fișier. Evident instrucțiunea `fclose` este folosită în vederea închiderii fișierelor.

Prin instrucțiunile `caseread` sunt citite etichetele claselor din fișierele `lx.txt` și `ly.txt`, care apoi sunt folosite pentru scrierea acestora în fișierul `clase.txt`, cu ajutorul comenzii `tblwrite`, împreună cu elementele tabelului de contingență.

3.2.5 Reprezentări grafice

Definiția 3.2.24. Numim diagramă prin batoane (bare) a unei distribuții statistice X de tip discret, reprezentarea grafică într-un sistem de axe rectangulare a segmentelor (batoanelor) date prin $\{(x_i, y) \mid 0 \leq y \leq \alpha f_i\}$, $i = \overline{1, n}$, $\alpha > 0$ fiind un factor de proporționalitate, iar f_i este frecvența absolută a valorii x_i .

Definiția 3.2.25. Numim diagramă cumulativă (ascendentă) a unei distribuții statistice X de tip discret, linia poligonală care unește punctele de coordonate $(x_1, \alpha F_0)$, $(x_1, \alpha F_1)$, $(x_2, \alpha F_1)$, $(x_2, \alpha F_2)$, $(x_3, \alpha F_2)$, \dots , $(x_n, \alpha F_n)$, unde F_i este frecvența cumulată (ascendentă) atașată valorii x_i , iar $\alpha > 0$ este un factor de proporționalitate.

Definiția 3.2.26. Numim histograma unei distribuții statistice X de tip continuu, diagrama obținută prin construirea de dreptunghiuri având drept baze clasele distribuției statistice și înălțimile astfel considerate încât ariile dreptunghiurilor să fie proporționale cu frecvențele claselor.

Observația 3.2.27. Când clasele distribuției statistice sunt de amplitudini egale, atunci înălțimile dreptunghiurilor histogramei sunt proporționale cu frecvențele claselor. Dacă factorul de proporționalitate este $\frac{1}{N}$, atunci se obține histograma frecvențelor relative.

Se știe că aria delimitată de curba ce reprezintă graficul unei densități de probabilitate și axa absciselor este 1. Din acest motiv, este de preferat ca factorul de proporționalitate în construirea histogramei frecvențelor relative să fie astfel ales încât

aria totală a dreptunghiurilor histogramei să fie de asemenea 1. Se verifică ușor că dacă N este volumul datelor, iar clasa $[a_{i-1}, a_i)$ are amplitudinea $h_i = a_i - a_{i-1}$, pentru $i = \overline{1, n}$, atunci funcția în scară ce delimitează aceste dreptunghiuri este dată prin:

$$\hat{f}(x) = \frac{f_i}{Nh_i}, \quad x \in [a_{i-1}, a_i), \quad i = \overline{1, n}.$$

Exemplul 3.2.28. Dacă se consideră datele statistice din Exemplul 3.2.10, și având în vedere sistematizarea acestora, atunci histograma frecvențelor absolute este prezentată în Figura 3.1.

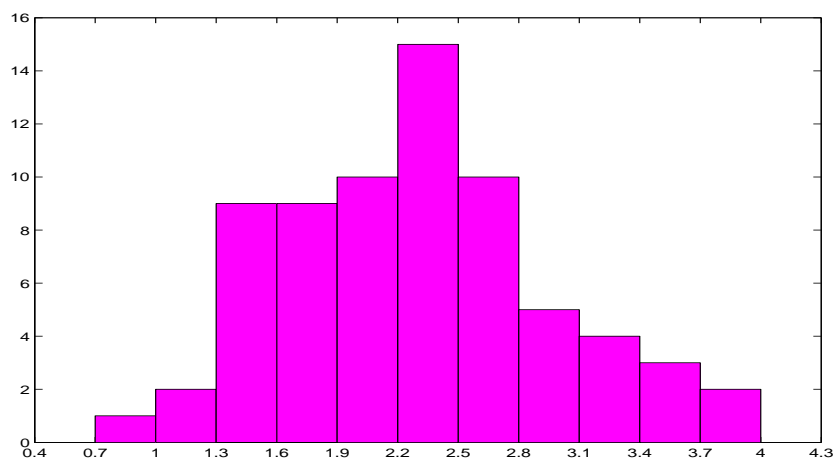


Figura 3.1: Histograma frecvențelor absolute

Metoda ferestrei mobile

Histograma frecvențelor relative a distribuției statistice reprezintă o aproximare destul de rudimentară a graficului densității de probabilitate a caracteristicii X . O ameliorare a acestei aproximări se obține dacă se aplică *metoda ferestrei mobile*.

Metoda ferestrei mobile constă în construirea clasei $(x - \frac{h}{2}, x + \frac{h}{2})$, $h > 0$, pentru fiecare $x \in (a, b)$, și determinarea frecvenței f_x a acestei clase, cu ajutorul căreia se obține curba dată prin $\hat{f}(x) = \frac{f_x}{Nh}$, $x \in (a, b)$.

Dacă se consideră funcția indicatoare K a intervalului $(-\frac{1}{2}, \frac{1}{2})$, adică

$$(3.2.1) \quad K(x) = \begin{cases} 1, & \text{dacă } x \in (-\frac{1}{2}, \frac{1}{2}), \\ 0, & \text{dacă } x \notin (-\frac{1}{2}, \frac{1}{2}), \end{cases}$$

atunci

$$\hat{f}(x) = \frac{1}{Nh} \sum_{i=1}^N K\left(\frac{x - x'_i}{h}\right).$$

Pentru a obține o formă mai regulată pentru \hat{f} se pot considera și alte tipuri de funcție nucleu K , care de regulă este o densitate de probabilitate simetrică.

Mai des se utilizează nucleul *Gaussian*

$$(3.2.2) \quad K(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}, \quad x \in \mathbb{R},$$

sau nucleul *parabolic* (Epanechnikov)

$$(3.2.3) \quad K(x) = \begin{cases} \frac{3}{4\sqrt{5}} \left(1 - \frac{x^2}{5}\right), & \text{dacă } |x| < \sqrt{5}, \\ 0, & \text{dacă } |x| \geq \sqrt{5}. \end{cases}$$

Alte nuclee ce s-ar putea folosi sunt:

$$(3.2.4) \quad K(x) = \begin{cases} \frac{15}{16} (1 - x^2)^2, & \text{dacă } |x| < 1, \\ 0, & \text{dacă } |x| \geq 1, \end{cases}$$

$$(3.2.5) \quad K(x) = \begin{cases} 1 + \cos(2\pi x) & \text{dacă } 0 < x < 1, \\ 0, & \text{dacă } |x| \geq 1, \end{cases}$$

$$(3.2.6) \quad K(x) = \begin{cases} 1 - |x| & \text{dacă } |x| < 1, \\ 0, & \text{dacă } |x| \geq 1. \end{cases}$$

Programul 3.2.29. Programul ce urmează generează N numere aleatoare ce urmează legea normală $\mathcal{N}(0, 1)$, după care, folosind același pas h pozitiv, reprezintă grafic funcția \hat{f} împreună cu densitatea de probabilitate a legii normale, pentru fiecare din cele șase nuclee, mai sus definite.

Calculul valorilor celor șase nuclee se face cu următoarea funcție Matlab:

```
function y = K(k,x)
switch k
case 1
    y=(abs(x)<1/2);
case 2
    y=1/(sqrt(2*pi))*exp(-x.^2/2);
case 3
    y=(abs(x)<sqrt(5));
```

```

        y=3/(4*sqrt(5))*(1-x.^2/5).*y;
    case 4
        y=(abs(x)<1);
        y=15/16*(1-x.^2).^2.*y;
    case 5
        y=((x>0)&(x<1));
        y=(1+cos(2*pi*x)).*y;
    case 6
        y=(abs(x)<1); y=abs(x).*y;
    otherwise
        error('Eroare')
    end
end

```

iar programul Matlab este

```

clf,clear
N=input('N='); h=input('h=');
X=randn(1,N); x=-3:0.01:3;
n=length(x); p=pdf('norm',x,0,1);
for ka=1:6
    s=zeros(1,n);
    for i=1:N
        s=s+K(ka,(x-X(i))/h);
    end
    f=s/(N*h);
    subplot(3,2,ka), plot(x,f,'-',x,p,'--')
end

```

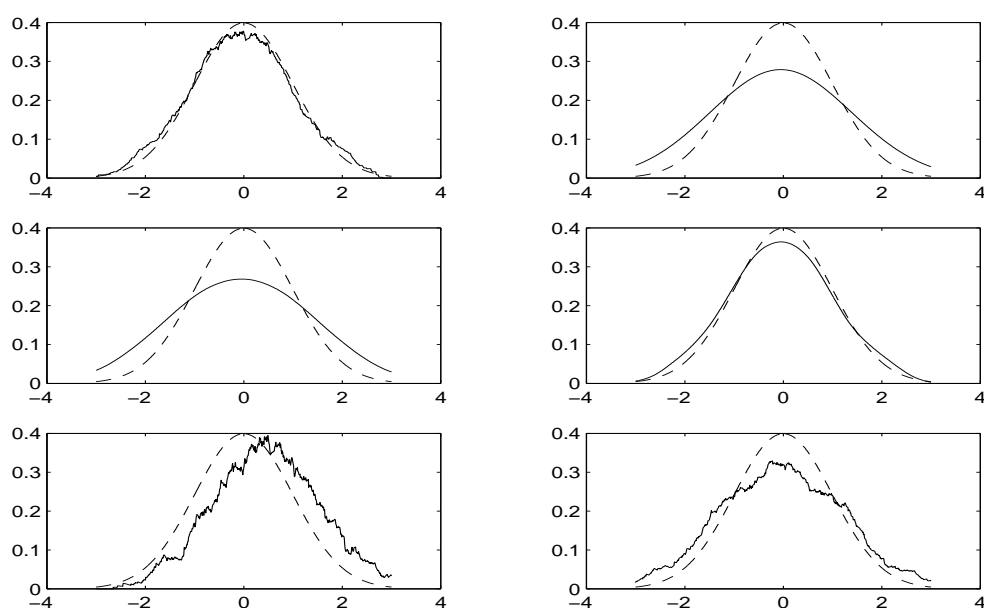
În urma executării programului, cu valorile $N=500$ și $h=1$, se obțin graficele din Figura 3.2.

Definiția 3.2.30. Numim poligonul frecvențelor al unei distribuții statistice X de tip continuu, poligonul obținut prin unirea punctelor de coordonate $(x_i, \alpha_i f_i)$, $i = \overline{1, n}$, unde α_i este factor de proporționalitate, iar f_i este frecvența clasei x_i .

Exemplul 3.2.31. Dacă se consideră datele statistice din Exemplul 3.2.10, și având în vedere sistematizarea acestora, atunci poligonul frecvențelor absolute este prezentată în Figura 3.3.

Definiția 3.2.32. Numim diagrame integrale (cumulative) ale frecvențelor cumulate ascendente, respectiv descendente, relative la distribuția statistică X de tip continuu, liniile poligonale obținute prin unirea punctelor de coordonate (a_k, F_k) , $k = \overline{0, n}$, și respectiv (a_k, F'_k) , $k = \overline{0, n}$.

Exemplul 3.2.33. Dacă se consideră datele statistice din Exemplul 3.2.10, și având în vedere sistematizarea acestora, atunci diagramele integrale ale frecvențelor cumulate (ascendente și descendente) sunt prezentate în Figura 3.4.

Figura 3.2: Grafice pentru $\hat{f}(x)$

Pentru trasarea diagramelor integrale, în tabelul următor calculăm frecvențele absolute ascendente și respectiv descendente, conform formulelor

$$F_k = \sum_{i=1}^k f_i, \quad F'_k = \sum_{i=k+1}^n f_i, \quad k = \overline{0, n}.$$

x	$a_0 - a_1$	$a_1 - a_2$	$a_2 - a_3$	$a_3 - a_4$	$a_4 - a_5$	$a_5 - a_6$	$a_6 - a_7$	$a_7 - a_8$	$a_8 - a_9$	$a_9 - a_{10}$	$a_{10} - a_{11}$
f	1	2	9	9	10	15	10	5	4	3	2
F	1	3	12	21	31	46	56	61	65	68	70
F'	69	67	58	49	39	24	14	9	5	2	0

Definiția 3.2.34. Numim nor statistic atașat caracteristicilor X și Y , punctele din plan obținute prin reprezentarea grafică a datelor primare (x'_k, y'_k) , $k = \overline{1, N}$.

Observația 3.2.35. Norul statistic este utilizat pentru observarea formei legăturii funcționale care există între cele două caracteristici. Amintim câteva legături mai des întâlnite:

$$y = ax + b, \quad x \in \mathbb{R}, \quad a, b \in \mathbb{R},$$

$$y = ax^b, \quad x > 0, \quad a, b \in \mathbb{R},$$

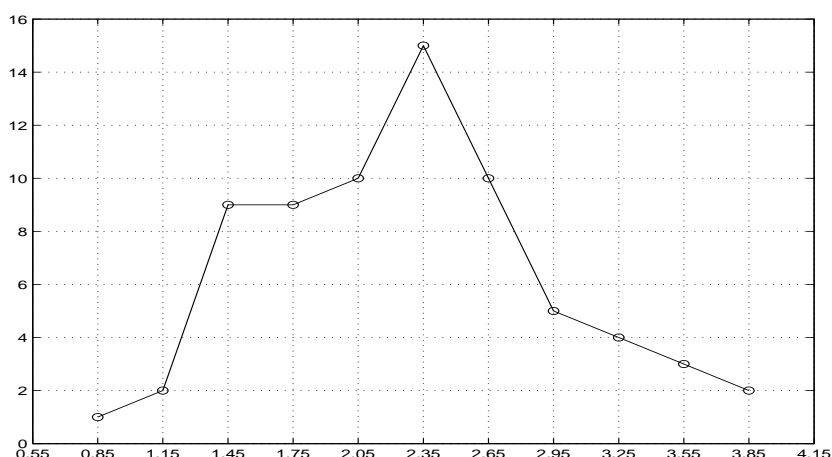


Figura 3.3: Poligonul frecvențelor absolute

$$y = ab^x, \quad x > 0, \quad a > 0, b > 0,$$

$$y = ax^2 + bx + c, \quad x \in \mathbb{R}, \quad a, b, c \in \mathbb{R},$$

$$y = a + \frac{b}{x}, \quad x \neq 0, \quad a, b \in \mathbb{R}.$$

Dacă legătura liniară $y = ax + b$ este ușor de observat din norul statistic, pentru celelalte se încearcă liniarizarea. De exemplu, dacă $y = ax^b$, prin logaritmare se ajunge la $\log y = \log a + b \log x$. Se notează $\log y = v$, $\log a = A$, $\log x = u$ și se obține legătura liniară $v = A + bu$, pentru $\log X$ și $\log Y$. Pentru aceasta se va reprezenta grafic norul statistic dat de punctele $(\log x'_k, \log y'_k)$, $k = \overline{1, N}$. Dacă punctele din acest nor statistic sunt situate în jurul unei drepte, atunci se poate considera că legătura dintre X și Y este de forma $y = ab^x$.

3.2.6 Funcții Matlab pentru reprezentarea grafică a datelor statistice

Sistemul Matlab, atât prin sistemul de bază, cât și prin *Statistics toolbox*, dispune de funcții pentru reprezentarea grafică a datelor statistice.

Cazul discret

Dacă în cercetare se consideră o caracteristică (variabilă) X de tip discret, funcțiile Matlab, deja prezentate, `bar` și `stairs` se pot utiliza cu succes.

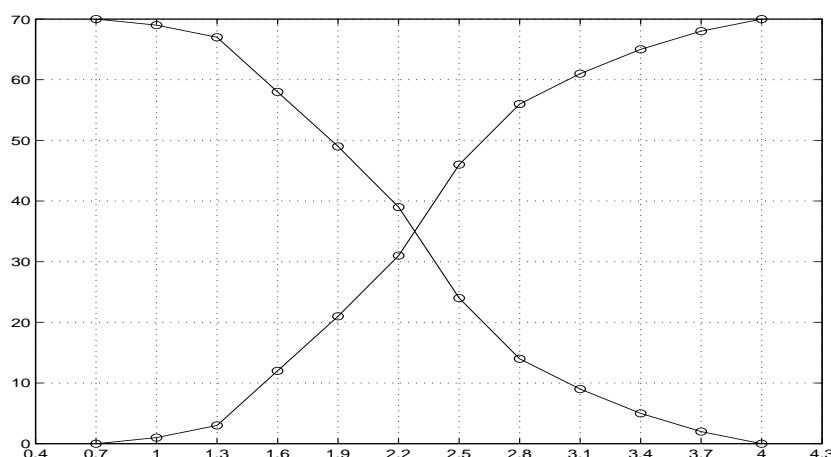


Figura 3.4: Diagramele integrale ale frecvențelor absolute

Pentru exemplificarea modului de utilizare a acestor funcții Matlab, vom scrie două funcții Matlab pentru construirea diagramei prin batoane și respectiv pentru construirea diagramei cumulative ascendente.

Funcția 3.2.36. Funcția Matlab pe care o scriem, va genera N numere aleatoare ce urmează o lege de probabilitate de tip discret, va construi diagrama prin batoane pentru datele astfel obținute, iar pe aceeași figură va reprezenta, pentru comparație, tot prin batoane, funcția de probabilitate a legii de probabilitate considerate. Pentru o vizualizare corectă a comparației, se impune, fie amplificarea funcției de probabilitate cu factorul N , fie împărțirea frecvențelor distribuției statistice cu N , adică prin considerarea frecvențelor relative. În cele ce urmează vom considera prima variantă.

```
function diag1(N,lege)
% Functia diag1 produce diagrama prin batoane
% N - reprezinta volumul datelor ce urmeaza a fi generate
% lege - reprezinta denumirea Matlab a legii de probabilitate
clf, colormap spring
switch lege
case 'unid'
    n = input('n (parametrul legii uniforme discrete):');
    x = random(lege,n,1,N);
    t = tabulate(x);
    P = N*pdf(lege,t(:,1),n);
    bar(t(:,1),[t(:,2),P])
    legend('Frecvente observate', 'Frecvente teoretice',1)
    return
case 'bino'
    n = input('n='); p = input('p=');
    x = random(lege,n,p,1,N);
```

```

t = tabulate(x+1);
P =N*pdf(lege,t(:,1)-1,n,p);
case 'hyge'
M = input('M='); K = input('K=');
n = input('n=');
x = random(lege,M,K,n,1,N);
t = tabulate(x+1);
P =N*pdf(lege,t(:,1)-1,M,K,n);
case 'poiss'
lambda = input('lambda=');
x = random(lege,lambda,1,N);
t = tabulate(x+1);
P =N*pdf(lege,t(:,1)-1,lambda);
case 'nbin'
r = input('r='); p = input('p=');
x = random(lege,r,p,1,N);
t = tabulate(x+1);
P =N*pdf(lege,t(:,1)-1,r,p);
case 'geo'
p = input('p=');
x = random(lege,p,1,N);
t = tabulate(x+1);
P =N*pdf(lege,t(:,1)-1,p);
otherwise
error('Lege discreta necunoscuta')
end
bar(t(:,1),[t(:,2),P])
legend('Frecvente observate', 'Frecvente teoretice',1)

```

Apelul funcției `diag1` se face prin

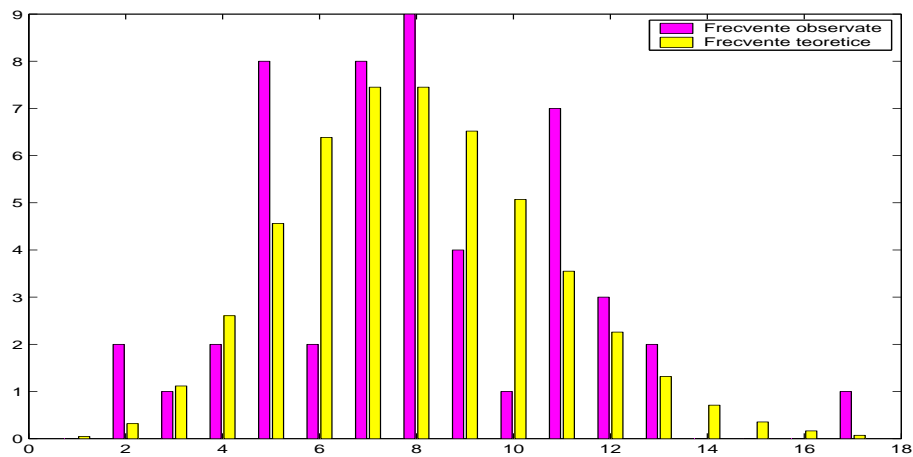
```
>>diag1(N,'lege')
```

unde valorile pentru `N` și `lege`, fie că sunt precizate în acest apel, fie sunt precizate înainte. De exemplu, comanda

```
>>diag1(50,'poiss')
```

are ca efect apelul funcției pentru legea lui Poisson, iar pe ecran se va cere introducerea parametrului `lambda`, după care pe ecran va fi reprezentat graficul din Figura 3.5, în cazul în care `lambda=7`. Să remarcăm totuși că la un nou apel, cu aceeași parametri, graficul diferă, deoarece sunt generate alte numere aleatoare.

Funcția 3.2.37. Următoarea funcție Matlab, va genera `N` numere aleatoare ce urmează o lege de probabilitate de tip discret, va construi diagrama cumulativă ascendentă pentru datele astfel obținute, iar pe aceeași figură va reprezenta, pentru comparație, funcția de repartiție a legii de probabilitate considerate. Pentru o vizualizare corectă a comparației, se impune, fie amplificarea funcției de repartiție cu factorul `N`, fie împărțirea frecvențelor cumulate ale distribuției statistice cu `N`, adică prin considerarea frecvențelor relative cumulate. În cele ce urmează vom considera a doua variantă.

Figura 3.5: Legea $Po(7)$

```
function diag2(N,lege)
% Functia diag2 produce diagrama cumulativa
% N - reprezinta volumul datelor generate
% lege - reprezinta denumirea legii de probabilitate
clf
switch lege
case 'unid'
    n=input('n (parametrul legii uniforme discrete):');
    x=random(lege,n,1,N);
    t=tabulate(x);
    xx=[t(1,1)-1;t(:,1);t(end,1)+1];
    cf=[0;cumsum(t(:,2))/N]/N;
    P=cdf(lege,xx,n);
    stairs(xx,cf,'k-'),hold on,
    stairs(xx,P,'k-.'),return
case 'bino'
    n=input('n='); p=input('p=');
    x=random(lege,n,p,1,N);
    t=tabulate(x+1);
    xx=[t(1,1)-2;t(:,1)-1;t(end,1)];
    cf=[0;cumsum(t(:,2))/N]/N;
    P=[0;cdf(lege,t(:,1)-1,n,p);1];
case 'hyge'
    M=input('M='); K=input('K=');
    n=input('n=');
    x=random(lege,M,K,n,1,N);
    t=tabulate(x+1);
    xx=[t(1,1)-2;t(:,1)-1;t(end,1)];
```

```

    cf=[0;cumsum(t(:,2))/N];
    P=[0;cdf(lege,t(:,1)-1,M,K,n);1];
case 'poiss'
    lambda=input('lambda=');
    x=random(lege,lambda,1,N);
    t=tabulate(x+1);
    xx=[t(1,1)-2;t(:,1)-1;t(end,1)];
    cf=[0;cumsum(t(:,2))/N];
    P=[0;cdf(lege,t(:,1)-1,lambda);1];
case 'nbin'
    r=input('r='); p=input('p=');
    x=random(lege,r,p,1,N);
    t=tabulate(x+1);
    xx=[t(1,1)-2;t(:,1)-1;t(end,1)];
    cf=[0;cumsum(t(:,2))/N];
    P=[0;cdf(lege,t(:,1)-1,r,p);1];
case 'geo'
    p=input('p=');
    x=random(lege,p,1,N);
    t=tabulate(x+1);
    xx=[t(1,1)-2;t(:,1)-1;t(end,1)];
    cf=[0;cumsum(t(:,2))/N];
    P=[0;cdf(lege,t(:,1)-1,p);1];
otherwise
    error('Lege discreta necunoscuta')
end
stairs(xx,cf,'k-'),hold on,
stairs(xx,P,'k-.')

```

Apelul funcției `diag2` se face analog apelului funcției `diag1` prezentată în Funcția 3.2.36:

```
>>diag2(N,'lege')
```

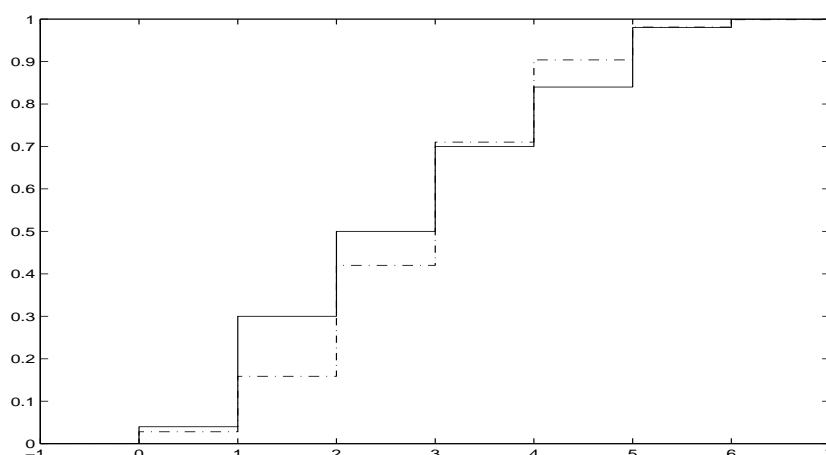
De exemplu, comanda

```
>>diag2(50,'bino')
```

are ca efect apelul funcției pentru legea binomială, iar pe ecran se vor cere introducerea parametrilor n și p , după care pe ecran va fi reprezentat graficul din Figura 3.6, în cazul în care $n=7$ și $p=0.4$. Diagrama cumulativă este reprezentată prin linie continuă, iar funcția de repartiție teoretică prin linie-punct. Și aici remarcăm că la un nou apel, cu aceiași parametri, graficul diferă, deoarece sunt generate alte numere aleatoare.

Cazul continuu

Reprezentări grafice pentru distribuția statistică a unei caracteristici (variabile) de tip continuu se obțin, fie prin utilizarea funcției `plot` și `bar`, fie prin funcțiile `hist`,

Figura 3.6: Legea $\mathcal{B}(7, 0.4)$

`histc`. Remarcăm faptul că `hist` și `histc` sunt funcții specifice sistemului de bază Matlab.

Funcția `hist`

Modurile de apel ale funcției `hist` sunt:

```
hist(y)
hist(y,nb)
hist(y,x)
n=hist(y)
n=hist(y,nb)
n=hist(y,x)
[n,x]=hist(y)
[n,x]=hist(y,nb)
```

Primele trei forme au ca efect reprezentarea grafică a histogramei corespunzătoare datelor conținute de vectorul `y`. Parametrul `nb` reprezintă numărul claselor de amplitudini egale, iar dacă acesta lipsește, se consideră `nb=10`. Amplitudinea acestor clase este obținută prin împărțirea lungimii intervalului definit prin valoarea minimă și valoarea maximă ale componentelor vectorului `y` la `nb`. Când se utilizează parametrul `x`, acesta trebuie să conțină mijloacele claselor.

Celelalte forme nu produc reprezentări grafice pentru histograme, ci doar calculează frecvențele absolute și respectiv mijloacele claselor în vectorii `n` și `x`, care au lungimea dată prin `nb`.

Toate formele acceptă parametrul y ca fiind matrice, caz în care operația se execută pentru fiecare coloană a matricei.

Funcția `histc`

Apelul funcției `histc` se poate face prin:

```
n=histc(y,a)
[n,nb]=hist(y,a)
```

care au ca efect, pentru datele vectorului y , obținerea în n a frecvențelor claselor precizate prin componentele vectorului a , care se impune a fi ordonate crescător. Valorile lui y care sunt înafara componentelor lui a nu sunt luate în considerare. De aceea e de dorit ca să fie folosite valorile $-\text{inf}$ și inf pentru a fi incluse toate valorile lui y .

Dacă y este o matrice, atunci funcția operează pentru fiecare coloană în parte.

Parametrul nb va conține numerele de ordine ale claselor în care au intrat fiecare din elementele lui y . Pentru elementele care nu au fost clasificate se obține valoarea 0.

Să remarcăm că Figura 3.1, Figura 3.3 și Figura 3.4 au fost realizate, respectiv cu programele Matlab:

```
clear, clf, load x.dat -ascii
N=length(x);
a=0.7:0.3:4; n=length(a)-1;
for j=1:n
    mij(j)=(a(j)+a(j+1))/2;
end
fr=histc(x,a); fr=fr(1:end-1);
bar(mij,fr,1)
axis([0.4 4.3 0 16])
set(gca,'xtick',[0.4:0.3:4.3]), colormap spring
```

```
clear, clf, load x.dat -ascii
N=length(x);
a=0.7:0.3:4; n=length(a)-1;
for j=1:n
    mij(j)=(a(j)+a(j+1))/2;
end
fr=histc(x,a); fr=fr(1:end-1);
plot(mij,fr,'k-o')
axis([0.55 4.15 0 16])
set(gca,'xtick',[0.55:0.3:4.15]), grid on
```

```
clear, clf, load x.dat -ascii
N=length(x);
a=0.7:0.3:4; n=length(a)-1;
fr=histc(x,a); fr=fr(1:end-1);
FA=[0;cumsum(fr)], FD=N-FA,
```

```
plot(a,FA,'k-o',a,FD,'k-o')
axis([0.4 4.3 0 70])
set(gca,'xtick',[0.4:0.3:4.3]), grid on
```

care utilizează funcțiile `histc` și `bar`.

Funcția 3.2.38. Următoarea funcție Matlab, va genera N numere aleatoare ce urmează o lege de probabilitate de tip continuu și va construi histograma pentru datele astfel obținute. Într-o figură alăturată se va construi histograma folosind funcția `bar`, împreună cu densitatea de probabilitate a legii de probabilitate considerate. Pentru o vizualizare corectă a comparației, se impune, împărțirea frecvențelor absolute cu Nh , unde h este amplitudinea claselor.

```
function hist0(N,lege)
% Functia hist0 produce histograma
% N - reprezinta volumul datelor generate
% lege - reprezinta denumirea legii de probabilitate
clf, colormap spring
n=fix(1+10/3*log10(N));
switch lege
case 'unif'
    a=input('a:'); b=input('b(a<b):');
    x=random(lege,a,b,1,N);
    [f,c]=hist(x,n); h=c(2)-c(1);
    t=c(1)-h:0.01:c(end)+h;
    y=pdf(lege,t,a,b);
case 'norm'
    mu=input('mu='); s=input('sigma=');
    x=random(lege,mu,s,1,N);
    [f,c]=hist(x,n); h=c(2)-c(1);
    t=c(1)-h:0.01:c(end)+h;
    y=pdf(lege,t,mu,s);
case 'logn'
    mu=input('mu='); s=input('sigma=');
    x=random(lege,mu,s,1,N);
    [f,c]=hist(x,n); h=c(2)-c(1);
    t=c(1)-h:0.01:c(end)+h;
    y=pdf(lege,t,mu,s);
case 'gam'
    a=input('a='); b=input('b=');
    x=random(lege,a,b,1,N);
    [f,c]=hist(x,n); h=c(2)-c(1);
    t=c(1)-h:0.01:c(end)+h;
    y=pdf(lege,t,a,b);
case 'exp'
    mu=input('mu=');
    x=random(lege,mu,1,N);
    [f,c]=hist(x,n); h=c(2)-c(1);
    t=c(1)-h:0.01:c(end)+h;
    y=pdf(lege,t,mu);
```

```

case 'beta'
    a=input('a='); b=input('b=');
    x=random(lege,a,b,1,N);
    [f,c]=hist(x,n); h=c(2)-c(1);
    t=c(1)-h:0.01:c(end)+h;
    y=pdf(lege,t,a,b);
case 'weib'
    a=input('a='); b=input('b=');
    x=random(lege,a,b,1,N);
    [f,c]=hist(x,n); h=c(2)-c(1);
    t=c(1)-h:0.01:c(end)+h;
    y=pdf(lege,t,a,b);
case 'rayl'
    b=input('b=');
    x=random(lege,b,1,N);
    [f,c]=hist(x,n); h=c(2)-c(1);
    t=c(1)-h:0.01:c(end)+h;
    y=pdf(lege,t,a);
otherwise
    error('Lege continua necunoscuta')
end
subplot(1,2,1), hist(x,n)
subplot(1,2,2), bar(c,f/(N*h),1)
hold on, plot(t,y,'k-')

```

Apelul funcției `hist0` se face prin

```
>>hist0(N,'lege')
```

De exemplu, comanda

```
>>hist0(500,'norm')
```

are ca efect apelul funcției pentru legea normală, iar pe ecran se va cere introducerea parametrilor μ și σ , după care pe ecran va fi reprezentat graficul din Figura 3.7, în cazul în care $\mu=0$ și $\sigma=1$.

Funcția `scatter`

Sistemul Matlab este prevăzut cu funcția `scatter` pentru reprezentarea norului statistic în cazul bidimensional. Apelul funcției se face prin:

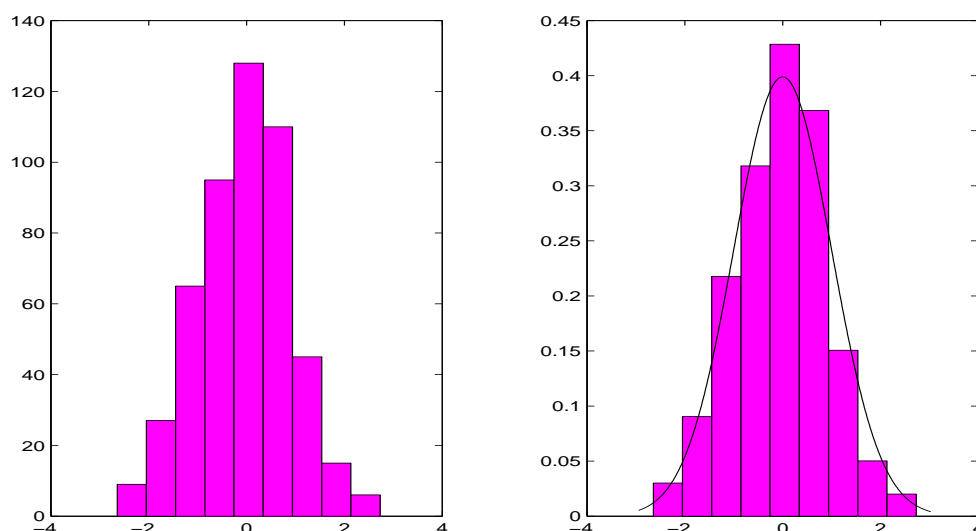
```

scatter(x,y)
scatter(x,y,s)
scatter(x,y,s,c)
scatter(x,y,s,c,m)
scatter(x,y,s,c,m,'filled')

```

Efectul executării unei astfel de instrucțiuni este producerea norului statistic precizat prin vectorii x și y , care au aceeași lungime.

Parametrul opțional s specifică mărimea marcajelor (cerculețe) punctelor norului. Acesta poate să fie un scalar, în care caz toate marcajele vor avea aceleași mărimi

Figura 3.7: Legea $\mathcal{N}(0, 1)$

sau poate să fie un vector de aceeași lungime cu x și y , caz în care se pot specifica mărimi diferite pentru puncte diferite.

Parametrul c specifică culoarea marcajelor și poate lua una din cele opt culori folosite și la funcția `plot`. Dacă c este vector, atunci trebuie să aibă aceeași lungime cu ceilalți parametri de tip vector, caz în care se pot specifica culorile fiecărui marcaj în parte.

Parametrul m permite înlocuirea marcajului implicit (cerculeț), prin marcajul precizat prin m .

Prezența opțiunii `'filled'` atrage după sine umbrirea interiorului marcajului.

Remarcăm faptul că se poate folosi cu succes funcția `plot`, dacă se folosește un singur tip de marcaj.

Funcția `scatter3`

Sistemul Matlab este prevăzut cu funcția `scatter3` pentru reprezentarea norului statistic în cazul tridimensional. Apelul funcției se face cu aceleași instrucțiuni ca și în cazul funcției `scatter`, doar că se mai înserează încă un parametru z , a treia dimensiune, după primii doi parametri x și y .

Funcția `gscatter`

Statistics toolbox dispune de funcția `gscatter` pentru reprezentarea norului statistic în cazul bidimensional pe grupe de puncte. Apelul funcției se face prin:

```
gscatter(x,y,g)
gscatter(x,y,s)
gscatter(x,y,'c','m',s)
gscatter(x,y,'c','m',s,'leg')
gscatter(x,y,'c','m',s,'leg','xl','yl')
```

Efectul executării unei astfel de instrucțiuni este producerea norului statistic precizat prin vectorii x și y , care au aceeași lungime, folosind gruparea specificată prin vectorul g de aceeași lungime. Punctele norului vor fi marcate la fel pentru valori identice corespunzătoare ale lui g .

Parametrul `'c'` specifică culorile marcajelor și are valoarea implicită dată prin `'bgrcmk'`.

Parametrul opțional `'m'` este un tablou de tip caractere recunoscute de funcția `plot`, valoarea implicită fiind caracterul `'.'` (punct).

Parametrul `s` specifică mărimea marcajelor punctelor norului. Acesta poate să fie un scalar, în care caz toate marcajele vor avea aceleași mărimi sau poate să fie un vector.

Dacă nu sunt specificate valori suficiente pentru toate grupele, sistemul Matlab efectuează o ciclare a valorilor specificate.

Parametrul `leg` precizează dacă se dorește afișarea unei legende (implicit acesta este `'on'`), sau nu se dorește legendă, când se va specifica valoarea `'off'`.

Parametrii `'xl'` și `'yl'` precizează etichetele pentru cele două axe de coordonate. Dacă acești parametri lipsesc, sistemul Matlab etichetează axele de coordonate cu denumirile variabilelor x și respectiv y .

Programul 3.2.39. Programul ce urmează generează n vectori aleatori ce urmează legea normală bidimensională și n numere aleatoare ce urmează legea uniformă discretă. Folosind aceste date, se va reprezenta norul vectorilor aleatori bidimensionali, efectuând gruparea conform numerelor aleatoare uniforme generate.

```
clear, mu(1)=input('m1='); mu(2)=input('m2=');
v(1,1)=input('sigma1^2=');
v(2,2)=input('sigma2^2=');
v(1,2)=input('Cov(X,Y)='); v(2,1) =v(1,2);
if det(v) <= 0
    error('Matricea v nu e pozitiv definita!')
end
N=input('N(parametrul legii uniforme):');
n=input('n='); X=mvnrnd(mu,v,n);
x=X(:,1); y=X(:,2); g=unidrnd(N,n,1);
gscatter(x,y,g,'k','o*.*')
```

Pentru $n=25$, $\mu=(5, 10)$, $v = \begin{pmatrix} 2 & -1 \\ -1 & 3 \end{pmatrix}$ și $N=4$, se obține Figura 3.8.

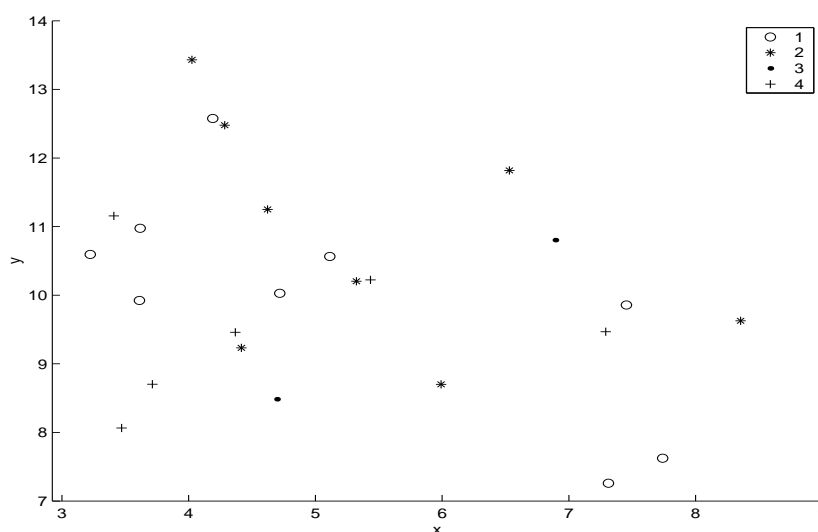


Figura 3.8: Nor statistic pe grupe

Funcția `plotmatrix`

Sistemul de bază Matlab conține funcția `plotmatrix`, care reprezintă pe aceeași figură mai mulți nori statistici. Formele de apel ale funcției sunt

```
plotmatrix(x)
plotmatrix(x,y)
plotmatrix(x,y,s)
```

unde x și y sunt matrice cu aceleași număr de linii, dar numărul coloanelor poate fi diferit, fie acesta m și respectiv n . Executarea unei instrucțiuni de acest tip, produce $m \times n$ nori statistici, pentru fiecare coloană a matricei x cu fiecare coloană a matricei y , cu marcasele specificate prin parametrul opțional s , care are sintaxa de la funcția `plot`.

Prima formă este echivalentă cu

```
plotmatrix(x,x)
```

cu excepția că pe diagonala principală a matricei norilor statistici vor fi reprezentate histogramele coloanelor matricei x .

Programul 3.2.40. Programul următor va genera o matrice x de tipul $(N, 3)$, care conține vectori aleatori ce urmează legea normală tridimensională, după care se va aplica funcția `plotmatrix`.

```

clear, mu(1)=input('m1=');
mu(2)=input('m2='); mu(3)=input('m3=');
v(1,1)=input('sigma1^2=');
v(2,2)=input('sigma2^2=');
v(3,3)=input('sigma3^2=');
v(1,2)=input('Cov(X,Y)='); v(2,1) =v(1,2);
v(1,3)=input('Cov(X,Z)='); v(3,1) =v(1,3);
v(2,3)=input('Cov(Y,Z)='); v(3,2) =v(1,3);
N=input('N='); X=mvnrnd(mu,v,N);
plotmatrix(X,'o')

```

Pentru $N=10$, $\mu=(5, 10, 15)$ și $v = \begin{pmatrix} 2 & -1 & 1 \\ -1 & 3 & 0 \\ 1 & 0 & 1 \end{pmatrix}$, se obține Figura 3.9.

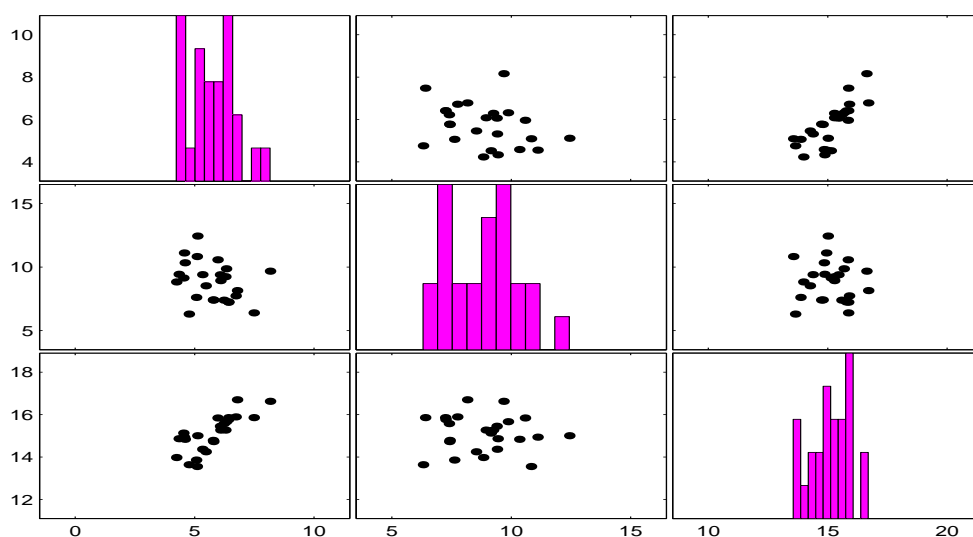


Figura 3.9: Funcția plotmatrix

Funcția gplotmatrix

Statistics toolbox conține funcția `gplotmatrix`, care are în principal același efect ca și funcția `plotmatrix`, având însă posibilitatea afișării punctelor norului statistic în funcție de grupele la care aparțin.

Formele de apel sunt cele de la funcția `gscatter`, doar că x și y , în acest caz, sunt matrice cu același număr de linii, iar numărul coloanelor putând fi diferite.

3.2.7 Funcția randtool

Având la dispoziție funcțiile `random` și `hist`, putem să prezentăm funcția demonstrativă `randtool`. Lansarea acestei funcții se face prin:

```
>>randtool
```

în urma căreia se produce o fereastră grafică interactivă (demonstrativă) privind generarea numerelor aleatoare și ilustrarea acestora cu ajutorul histogramelor.

Fixarea (stabilirea) legii de probabilitate în scop demonstrativ se face prin alegerea din meniul legilor de probabilitate situat în partea stângă sus a ferestrei.

Volumul numerelor aleatoare, ce urmează a fi generate, se precizează prin introducerea acestuia în fereastra din partea dreaptă sus.

Pentru stabilirea valorilor parametrilor legii de probabilitate considerate se poate proceda în două moduri. Fie prin introducerea în ferestrele corespunzătoare ale valorilor dorite, fie prin deplasarea barelor atașate acestora. În plus, limite pentru parametrii legii de probabilitate considerate pot fi precizate prin introducerea acestora în ferestrele considerate.

Activarea butonului `output` are ca efect salvarea numerelor aleatoare curente în variabila `ans` sau în variabila precizată de utilizator, iar butonul `resample` permite repetarea generării de numere aleatoare cu același volum și aceeași parametri.

3.2.8 Funcțiile pie și pie3

Reprezentări ale datelor prin sectoare circulare în plan, respectiv în spațiu, se obțin prin funcțiile `pie` și `pie3`, care se apelează prin formulele:

```
pie(x,ex,label)
pie3(x,ex,label)
```

unde `x` este un vector având componentele pozitive și normalizate, adică astfel încât suma acestora să fie 1. Dacă suma acestora va fi mai mică decât 1, atunci numai o parte a cercului va fi reprezentată grafic.

Parametrul `ex` este opțional și specifică sectoarele de cerc ce urmează să fie detașate pe figură, acesta, dacă este prezent, trebuie să fie un vector numeric de aceeași lungime cu vectorul `x`.

Parametrul opțional `label`, specifică etichetele sectoarele de cerc, acesta, dacă este prezent, trebuie să fie un vector de aceeași lungime cu vectorul `x`, dar care conțin numai șiruri de caractere.

Programul 3.2.41. Programul ce urmează generează n numere aleatoare ce urmează legea uniformă discretă $\mathcal{U}(N)$, construiește tabelul sistematizat, după care reprezintă datele sistematizate folosind funcțiile `pie` și `pie3`.

```
n=input('n='); N=input('N=');
x=unidrnd(N,1,n);
t=tabulate(x); d=length(t(:,3));
```

```
ex=zeros(1,d); ex(1)=1;
subplot(2,1,1), pie(t(:,3)'/100,ex)
subplot(2,1,2), pie3(t(:,3)'/100,ex)
colormap summer
```

Prin execuția acestui program, cu datele $n=50$ și $N=4$, se obțin reprezentările grafice din Figura 3.10.

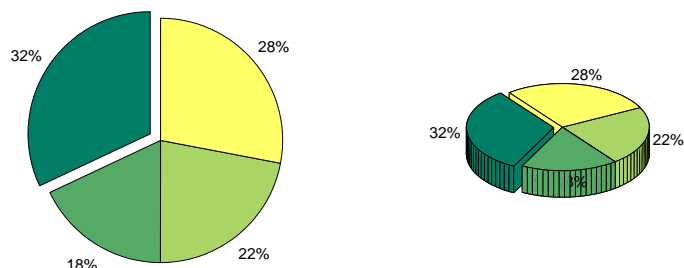


Figura 3.10: Diagrame circulare

3.3 Parametrii distribuțiilor statistice

Se consideră datele statistice primare x'_k , $k = \overline{1, N}$, relative la caracteristica X , din care se obține distribuția statistică

$$X \begin{pmatrix} x_i \\ f_i \end{pmatrix}_{i=\overline{1,n}}.$$

3.3.1 Parametri statistici ce măsoară tendința

Definiția 3.3.1. Media (aritmetică) a distribuției statistice a caracteristicii X este dată prin

$$\bar{x}_a = \frac{1}{N} \sum_{k=1}^N x'_k = \frac{1}{N} \sum_{k=1}^n f_k x_k = \sum_{k=1}^n p_k x_k.$$

Definiția 3.3.2. Media geometrică a distribuției statistice a caracteristicii pozitive X este dată prin

$$\bar{x}_g = \sqrt[N]{\prod_{k=1}^N x'_k} = \sqrt[n]{\prod_{k=1}^n x_k^{f_k}}.$$

Observația 3.3.3. În aplicații se lucrează mai ușor cu

$$\log \bar{x}_g = \frac{1}{N} \sum_{k=1}^N \log x'_k = \frac{1}{N} \sum_{k=1}^n f_k \log x_k = \sum_{k=1}^n p_k \log x_k.$$

Definiția 3.3.4. Media armonică a distribuției statistice a caracteristicii nenule X este dată prin

$$\bar{x}_h = \frac{N}{\sum_{k=1}^N \frac{1}{x'_k}} = \frac{N}{\sum_{k=1}^n f_k \frac{1}{x_k}} = \frac{1}{\sum_{k=1}^n p_k \frac{1}{x_k}}.$$

Observația 3.3.5. Între cele trei medii definite mai înainte există relațiile cunoscute $\bar{x}_h \leq \bar{x}_g \leq \bar{x}_a$.

3.3.2 Funcțiile mean, geomean și harmean

Sistemul de bază Matlab conține funcția mean, iar *Statistics toolbox* dispune de celelalte două funcții, geomean și harmmean.

Apelul acestor funcții se face prin:

```
ma=mean(x)
mg=geomean(x)
mh=harmmean(x)
```

și au ca efect calculul mediilor corespunzătoare, pentru componentele vectorului x , iar dacă x este matrice, atunci operația se execută pentru fiecare coloană a matricei. Prin urmare, în acest ultim caz, rezultatul este un vector având atâtea componente câte coloane are matricea x .

3.3.3 Funcția trimmean

Sistemul Matlab, prin *Statistics toolbox*, dispune de funcția trimmean, care se poate apela prin

```
m=trimmean(x,p)
```

având ca efect calculul mediei aritmetice a vectorului x , după ce au fost eliminate cele mai mici $\frac{p}{2}\%$ componente și cele mai mari $\frac{p}{2}\%$ componente, iar dacă x este matrice, această operație se face pentru fiecare coloană în parte.

Definiția 3.3.6. Numim mediana distribuției statistice a caracteristicii X , valoarea numerică \bar{m} care împarte datele statistice, ordonate crescător, în două părți egale.

Observația 3.3.7. Fie datele statistice primare ordonate în mod nedescrescător

$$x'_{(1)} \leq x'_{(2)} \leq \dots \leq x'_{(N)},$$

atunci mediana va fi dată prin

$$\overline{m} = \begin{cases} x'_{(k)}, & \text{dacă } N = 2k - 1, \\ \frac{x'_{(k)} + x'_{(k+1)}}{2}, & \text{dacă } N = 2k. \end{cases}$$

Observația 3.3.8. Când datele statistice sunt grupate, atunci se determină prima dată *intervalul median* $[a_{j-1}, a_j)$ astfel încât pentru frecvențele cumulate F_{j-1} și F_j să fie satisfăcute inegalitățile $F_{j-1} < \frac{N}{2}$ și $F_j > \frac{N}{2}$. Folosind apoi interpolarea liniară se ia ca mediană

$$\overline{m} = a_{j-1} + d_j \frac{\frac{N}{2} - F_{j-1}}{f_j},$$

unde d_j este amplitudinea intervalului median.

Când caracteristica cercetată este de tip discret, mediana este $\overline{m} = x_j$.

3.3.4 Funcția median

Sistemul de bază Matlab dispune de funcția `median`. Apelarea acesteia se face prin

```
m=median(x)
```

și are ca efect calculul mediane pentru componentele vectorului x , iar dacă x este matrice, atunci operația se execută pentru fiecare coloană a matricei. Prin urmare, în acest ultim caz, rezultatul este un vector având atâtea componente câte coloane are matricea x .

3.3.5 Parametri statistici ce măsoară dispersarea

Definiția 3.3.9. Numim *cuartile ale distribuției statistice a caracteristicii X* , valorile numerice \overline{Q}_1 (cuartila inferioară), $\overline{Q}_2 = \overline{m}$, \overline{Q}_3 (cuartila superioară), care împart datele statistice, ordonate crescător, în patru părți egale.

Observația 3.3.10. Când datele statistice sunt grupate, ca și în cazul mediane, folosind interpolarea liniară considerăm

$$\overline{Q}_1 = a_{i-1} + d_i \frac{\frac{N}{4} - F_{i-1}}{f_i},$$

$$\overline{Q}_3 = a_{k-1} + d_k \frac{\frac{3N}{4} - F_{k-1}}{f_k},$$

după ce s-a determinat *intervalul cuartilic inferior* $[a_{i-1}, a_i)$, astfel încât să aibă loc $F_{i-1} < \frac{N}{4}$ și $F_i > \frac{N}{4}$, respectiv *intervalul cuartilic superior* $[a_{k-1}, a_k)$, astfel încât $F_{k-1} < \frac{3N}{4}$ și $F_k > \frac{3N}{4}$.

Când caracteristica cercetată este de tip discret, avem că $\overline{Q}_1 = x_i$ și $\overline{Q}_3 = x_k$.

Observația 3.3.11. În mod analog se definesc parametri numerici cum ar fi *decilele* și *centilele*.

3.3.6 Funcția prctile

Statistics toolbox conține funcția `prctile`, care se poate apela prin

`c=prctile(x,p)`

și care calculează pentru componentele vectorului x , după ce acestea au fost ordonate crescător, centilele precizate prin parametrul p , care conține unul sau mai multe numere întregi de la 1 la 99. De exemplu, $p=[25, 50, 75]$ produce respectiv cuartila inferioară (\overline{Q}_1), mediana, cuartila superioară (\overline{Q}_3):

$$\overline{Q}_1 = \frac{x_{(i)} + x_{(j)}}{2}, \quad i = \text{floor}\left(\frac{n+1}{4}\right), \quad j = \text{ceil}\left(\frac{n}{4}\right),$$

$$\overline{Q}_3 = \frac{x_{(i)} + x_{(j)}}{2}, \quad i = \text{floor}\left(\frac{3(n+1)}{4}\right), \quad j = \text{ceil}\left(\frac{3n}{4}\right).$$

Dacă x este matrice, atunci se operează pentru fiecare coloană în parte.

Definiția 3.3.12. Numim *mod al distribuției statistice a caracteristicii X orice punct \overline{m}_0 de maxim local al distribuției statistice*.

Observația 3.3.13. Când distribuția statistică are un singur mod spunem că avem distribuție statistică *unimodală*. Dacă există două sau mai multe moduri spunem că avem distribuții statistice *bimodale*, respectiv *multimodale*.

Observația 3.3.14. Când datele statistice sunt grupate, pentru determinarea modului, se determină *intervalul modal*, adică intervalul cu frecvența maximă locală. Dacă intervalul modal este $[a_{k-1}, a_k)$, atunci se consideră

$$\overline{m}_0 = a_{k-1} + d_k \frac{\Delta f_k}{\Delta f_k - \Delta f_{k+1}},$$

unde $d_k = a_k - a_{k-1}$, $\Delta f_k = f_k - f_{k-1}$, $\Delta f_{k+1} = f_{k+1} - f_k$. Formula se obține ușor dacă se intersectează interpolantul liniar al punctelor (a_{k-1}, f_{k-1}) , (a_k, f_k) cu interpolantul liniar al punctelor (a_{k-1}, f_k) , (a_k, f_{k+1})

Când caracteristica cercetată este de tip discret, avem că $\overline{m}_0 = x_k$.

Definiția 3.3.15. Numim *moment de ordin k al distribuției statistice a caracteristicii X , valoarea numerică*

$$\overline{\nu}_k = \frac{1}{N} \sum_{i=1}^N x_i^k = \frac{1}{N} \sum_{i=1}^n f_i x_i^k = \sum_{i=1}^n p_i x_i^k.$$

Definiția 3.3.16. Numim amplitudinea (interval de variație) distribuției statistice a caracteristicii X , valoarea numerică

$$\overline{w} = \max\{x'_1, x'_2, \dots, x'_N\} - \min\{x'_1, x'_2, \dots, x'_N\} = x_{\max} - x_{\min}.$$

Definiția 3.3.17. Numim abatere cuartilică (interval intercuartilic) a distribuției statistice a lui X , diferența dintre cuartila superioară și cuartila inferioară, adică diferența $\overline{Q}_3 - \overline{Q}_1$.

Observația 3.3.18. Mai întâlnim, în aplicații, variația intercuartilică dată prin formula

$$\overline{Q} = \frac{(\overline{Q}_3 - \overline{m}) + (\overline{m} - \overline{Q}_1)}{2} = \frac{\overline{Q}_3 - \overline{Q}_1}{2},$$

respectiv abaterea cuartilă relativă

$$\overline{Q}_r = \frac{\overline{Q}_3 - \overline{Q}_1}{\overline{Q}_2} = \frac{\overline{Q}_3 - \overline{Q}_1}{\overline{m}}.$$

Observația 3.3.19. Dacă $\overline{Q}_3 - \overline{Q}_1 < \frac{\overline{w}}{2}$, atunci distribuția se consideră intens concentrată, iar în caz contrar, intens dispersată.

Definiția 3.3.20. Numim abatere medie (absolută) a distribuției statistice X , valoarea numerică

$$\overline{\delta} = \frac{1}{N} \sum_{k=1}^N |x'_k - \overline{x}| = \frac{1}{N} \sum_{k=1}^n f_k |x_k - \overline{x}| = \sum_{k=1}^n p_k |x_k - \overline{x}|,$$

unde $\overline{x} = \overline{x}_a$.

Definiția 3.3.21. Numim moment centrat de ordin k al distribuției statistice X , valoarea numerică

$$\overline{\mu}_k = \frac{1}{N} \sum_{i=1}^N (x'_i - \overline{x})^k = \frac{1}{N} \sum_{i=1}^n f_i (x_i - \overline{x})^k = \sum_{i=1}^n p_i (x_i - \overline{x})^k.$$

Definiția 3.3.22. Momentul centrat de ordinul doi al distribuției statistice X se numește dispersie și se notează $\overline{\sigma}^2 = \overline{\mu}_2$, iar $\overline{\sigma} = \sqrt{\overline{\mu}_2}$ se numește abatere medie pătratică sau abatere standard.

Observația 3.3.23. Alte formule de calcul pentru dispersie sunt

$$\overline{\sigma}^2 = \frac{1}{N} \left[\sum_{i=1}^n f_i x_i^2 - \frac{1}{N} \left(\sum_{i=1}^n f_i x_i \right)^2 \right],$$

$$\overline{\sigma}^2 = \frac{1}{N} \sum_{i=1}^n f_i (x_i - a)^2 - (\overline{x} - a)^2, \quad a \in \mathbb{R}, \quad (\text{formula lui König}).$$

Observația 3.3.24. Se observă că pentru calculul dispersiei ar fi necesar să fie parcurse datele statistice de două ori, odată pentru calculul mediei aritmetice \bar{x} , iar apoi pentru calculul lui $\bar{\sigma}^2$, care folosește pe \bar{x} . Prezentăm o metodă pentru calculul dispersiei printr-o singură parcurgere a datelor statistice.

Pentru aceasta introducem notațiile

$$t_j = \sum_{i=1}^j x'_i \quad \text{și} \quad s_j^2 = \sum_{i=1}^j \left(x'_i - \frac{t_j}{j} \right)^2, \quad j = \overline{1, N}.$$

Vom arăta, prima dată, că are loc relația de recurență

$$s_j^2 = s_{j-1}^2 + \frac{1}{j(j-1)} (j x'_j - t_j)^2, \quad j = \overline{2, N}.$$

Pentru stabilirea relației de recurență se scrie succesiv

$$\begin{aligned} s_j^2 - s_{j-1}^2 &= \sum_{i=1}^j \left(x'_i - \frac{t_j}{j} \right)^2 - \sum_{i=1}^{j-1} \left(x'_i - \frac{t_j - x'_j}{j-1} \right)^2 \\ &= \sum_{i=1}^j x_i'^2 - \frac{2t_j}{j} \sum_{i=1}^j x'_i + \frac{t_j^2}{j} - \sum_{i=1}^{j-1} x_i'^2 + \frac{2(t_j - x'_j)}{j-1} \sum_{i=1}^{j-1} x'_i - \frac{(t_j - x'_j)^2}{j-1} \\ &= x_j'^2 - \frac{2t_j^2}{j} + \frac{t_j^2}{j} + \frac{2(t_j - x'_j)^2}{j-1} - \frac{(t_j - x'_j)^2}{j-1} = x_j'^2 - \frac{t_j^2}{j} + \frac{(t_j - x'_j)^2}{j-1}. \end{aligned}$$

Dacă reținem extremitățile șirului de egalități avem:

$$s_j^2 - s_{j-1}^2 = \frac{1}{j(j-1)} (j^2 x_j'^2 - 2j t_j x'_j + t_j^2) = \frac{1}{j(j-1)} (j x'_j - t_j)^2,$$

ce trebuie arătat.

Observăm că $\bar{x} = \frac{1}{N} t_N$, iar $\bar{\sigma}^2 = \frac{1}{N} s_N^2$. Prin urmare formula de recurență de la punctul precedent permite calculul lui $\bar{\sigma}^2$ odată cu calculul lui \bar{x} , deci printr-o parcurgere a datelor statistice.

Algoritmul de calcul poate fi descris prin următoarele etape:

1. $s := 0, t := x'_1$.
2. Pentru $j = \overline{2, N}$: $t := t + x'_j, s := s + \frac{1}{j(j-1)} (j x'_j - t)^2$.
3. $\bar{x} = \frac{t}{N}; \quad \bar{\sigma}^2 = \frac{s}{N}$.

Definiția 3.3.25. Numim coeficient de variație al distribuției statistice X , raportul

$$\bar{v} = \frac{\bar{\sigma}}{\bar{x}},$$

care de regulă se exprimă în procente ($100 \bar{v} \%$).

Definiția 3.3.26. Numim coeficient de variație intercuartil al distribuției statistice X , raportul

$$\bar{q} = \frac{\bar{Q}}{\bar{m}} = \frac{\frac{\bar{Q}_3 - \bar{Q}_1}{2}}{\bar{Q}_2} = \frac{\bar{Q}_3 - \bar{Q}_1}{2\bar{Q}_2}.$$

Definiția 3.3.27. Numim coeficientul de asimetrie intercuartil (Yule) al distribuției statistice X , raportul

$$\bar{Y} = \frac{(\bar{Q}_3 - \bar{Q}_2) - (\bar{Q}_2 - \bar{Q}_1)}{(\bar{Q}_3 - \bar{Q}_2) + (\bar{Q}_2 - \bar{Q}_1)} = \frac{\bar{Q}_3 + \bar{Q}_1 - 2\bar{m}}{\bar{Q}_3 - \bar{Q}_1}.$$

Observația 3.3.28. Coeficientul lui Yule satisface relațiile $-1 \leq \bar{Y} \leq 1$, iar dacă $\bar{Q}_3 - \bar{Q}_2 > \bar{Q}_2 - \bar{Q}_1$ ($\bar{Y} > 0$), atunci $\bar{m}\bar{o} < \bar{m} < \bar{x}$, în caz contrar $\bar{m}\bar{o} > \bar{m} > \bar{x}$.

Definiția 3.3.29. Numim coeficienți ai lui Pearson relativ la distribuția statistică X , rapoartele:

$$\begin{aligned} \bar{s} &= \frac{\bar{x} - \bar{m}\bar{o}}{\bar{\sigma}} \quad (\text{coeficient de asimetrie}), \\ \bar{\beta}_1 &= \frac{\bar{\mu}_3^2}{\bar{\mu}_2^3} \quad (\text{skewness}), \\ \bar{\beta}_2 &= \frac{\bar{\mu}_4}{\bar{\mu}_2^2} \quad (\text{kurtosis}). \end{aligned} \tag{3.3.1}$$

Definiția 3.3.30. Numim coeficienți ai lui Fisher relativ la distribuția statistică a lui X , valorile numerice

$$\begin{aligned} \bar{\gamma}_1 &= \sqrt{\bar{\beta}_1} = \frac{\bar{\mu}_3}{\bar{\sigma}^3} \quad (\text{asimetria}), \\ \bar{\gamma}_2 &= \bar{\beta}_2 - 3 = \frac{\bar{\mu}_4}{\bar{\mu}_2^2} - 3 = \frac{\bar{\mu}_4}{\bar{\sigma}^4} - 3 \quad (\text{excesul}). \end{aligned} \tag{3.3.2}$$

Observația 3.3.31. Pentru curba lui Gauss avem că $\bar{\beta}_1 = 0$, deci $\bar{\gamma}_1 = 0$. Când $\bar{\gamma}_1 < 0$, maximul curbei distribuției este deplasat spre dreapta în raport cu valoarea medie, iar când $\bar{\gamma}_1 > 0$, maximul este deplasat spre stânga.

De asemenea, pentru curba lui Gauss avem că $\bar{\beta}_2 = 3$, deci $\bar{\gamma}_2 = 0$. Astfel, dacă $\bar{\beta}_2 > 3$ ($\bar{\gamma}_2 > 0$) atunci curba distribuției este mai îngustată decât curba lui Gauss, distribuția numindu-se *leptocurtică*. Dacă $\bar{\beta}_2 < 3$ ($\bar{\gamma}_2 < 0$) curba distribuției este mai turtită, distribuția numindu-se *platicurtică*.

O parte din acești parametri statistici există în sistemul de bază Matlab, iar alții se găsesc în *Statistics toolbox*. Facem prezentarea lor în cele ce urmează.

3.3.7 Funcția moment

Apelul funcției moment se poate face prin

```
s=moment(x,k)
```

și are ca efect calculul, pentru componentele vectorului x , a momentului centrat de ordin k , iar dacă x este matrice, atunci se operează pe fiecare coloană în parte.

3.3.8 Funcțiile var și std

Apelurile funcțiilor `var` și `std` se pot face prin:

```
v=var(x)
s=std(x)
v=var(x,w)
v=var(x,1)
```

Primele două forme calculează, pentru componentele vectorului x , dispersia definită prin

$$v = \frac{1}{m-1} \sum_{i=1}^m (x_i - \bar{x})^2,$$

unde m este lungimea vectorului, și respectiv abaterea standard $s = \sqrt{v}$.

Forma a doua pentru `var`, când parametrul opțional w este prezent, acesta trebuie să fie un vector cu ponderi pozitive a căror sumă este 1, și are ca efect calculul lui v după formula

$$v = \sum_{i=1}^m w_i (x_i - \bar{x})^2, \quad \bar{x} = \sum_{i=1}^m w_i x_i.$$

Forma a treia pentru `var`, când parametrul opțional 1 este prezent, are ca efect calculul lui v după formula momentului centrat de ordinul 2.

În toate cazurile, dacă x este matrice, se operează pe fiecare coloană în parte.

3.3.9 Funcțiile range, iqr, mad, skewness și kurtosis

Apelurile funcțiilor `range`, `iqr`, `mad`, `skewness`, `kurtosis` se fac prin

```

r=range(x)
iq=iqr(x)
md=mad(x)
sk=skewness(x)
k=kurtosis(x)

```

și calculează pentru componentele vectorului x , respectiv amplitudinea, abaterea cuartilică, abaterea absolută, asimetria (cu formula (3.3.2)) și excesul (cu formula (3.3.1)).

Dacă x este matrice, se operează pe fiecare coloană în parte.

3.3.10 Funcțiile **max** și **min**

Prezentăm numai modul de apel pentru funcția **max**, deoarece funcția **min** se utilizează la fel.

Pentru **max** avem următoarele forme de apel:

```

r=max(x)
r=max(x,y)
[r,k]=max(x)
[r,k]=max(x,y)

```

Prin acestea se calculează valoarea maximă a componentelor vectorului x , iar în k se obțin și indicii corespunzători componentelor. Dacă este prezent parametrul opțional y , care este un vector de aceeași lungime cu cea a lui x , atunci r este un vector ce va conține valorile maxime, prin compararea celor doi vectori pe componente. Dacă x este matrice, respectiv când apare y matrice de aceeași dimensiune, atunci comparațiile sunt efectuate pe coloane în primul caz, iar în al doilea caz se face compararea element cu element.

3.3.11 Funcțiile **sort** și **sortrows**

Funcția **sort** are următoarele forme de apel

```

y=sort(x)
y=sort(x,k)
[y,s]=sort(x)
[y,s]=sort(x,k)

```

și **sortrows** aceleași forme de apel.

Funcția **sort** ordonează crescător elementele fiecărei coloane a matricei x , când $k=1$ (care este valoarea implicită), respectiv ordonarea crescătoare pe linii, când $k=2$. Vectorul s păstrează indicii elementelor dinainte de ordonare. Dacă x este vector, atunci parametrul opțional lipsește și este efectuată ordonarea crescătoare a componentelor vectorului, eventual cu păstrarea indicilor componentelor dinainte de ordonare în s .

Funcția `sortrows` ordonează lexicografic liniile matricei x , având în vedere toate elementele liniilor matricei (când parametrul k lipsește) sau numai elementele ce aparțin coloanelor precizate prin vectorul k , cu posibilitatea păstrării în s a indicilor liniilor înainte de ordonare.

3.3.12 Funcțiile `sum`, `prod`, `cumsum` și `cumprod`

Apelurile funcțiilor `sum`, `prod`, `cumsum`, `cumprod`, se pot realiza prin:

```
y=sum(x)
y=sum(x,k)
y=prod(x)
y=prod(x,k)
y=cumsum(x)
y=cumsum(x,k)
y=cumprod(x)
y=cumprod(x,k)
```

Execuția acestor instrucțiuni calculează respectiv sumele, sumele parțiale, produsele, produsele parțiale pentru fiecare coloană a matricei x , când $k=1$ (valoare implicită), iar când $k=2$, se execută aceste operații pe linii. Dacă x este vector, atunci parametrul opțional k lipsește, iar operațiile se execută asupra componentelor vectorului.

3.3.13 Funcția `diff`

Apelul funcției `diff` se realizează prin:

```
y=diff(x)
y=diff(x,r)
y=diff(x,r,k)
```

În urma executării unei astfel de instrucțiuni, se calculează diferențele de ordin r ale elementelor liniilor consecutive ale matricei x , când $k=1$ (valoare implicită), respectiv diferențele de ordin r ale elementelor coloanelor consecutive ale matricei x , când $k=2$, adică

$$\begin{aligned} Y_{i,j} &= x_{i+1,j} - x_{i,j}, & i = \overline{1, m-1}, j = \overline{1, n}, & (k=1) \\ Y_{i,j} &= x_{i,j+1} - x_{i,j}, & i = \overline{1, m}, j = \overline{1, n-1}, & (k=2). \end{aligned}$$

Diferențele de ordin doi pentru x sunt diferențele de ordinul întâi pentru y . Dacă x este vector, atunci parametrul k lipsește, iar operația se execută asupra componentelor vectorului.

3.3.14 Funcțiile `nanmean`, `nanmedian`, `nanstd`, `nanmin`, `nanmax` și `nansum`

Aceste funcții din *Statistics toolbox* sunt definite ca și funcțiile corespunzătoare `mean`, `median`, `std`, `min`, `max`, `sum`, cu excepția faptului că sunt excluse din calcule valorile NaN.

3.3.15 Corecțiile lui Sheppard

Dacă datele statistice relative la caracteristica X de tip continuu au fost grupate, în unele cazuri, se cere să fie aplicate *corecțiile lui Sheppard* pentru momentele $\bar{\nu}_k$. Dacă se notează cu $\bar{\nu}'_k$ momentul corectat de ordinul k , atunci

$$\bar{\nu}'_k = \frac{1}{k+1} \sum_{i=0}^{\lfloor \frac{k}{2} \rfloor} \binom{k+1}{2i+1} \left(\frac{d}{2}\right)^{2i} \bar{\nu}'_{k-2i}, \quad k = 1, 2, \dots$$

unde d este amplitudinea claselor.

Folosind formula precedentă scriem primele patru momente corectate:

$$\bar{\nu}'_1 = \bar{\nu}_1, \quad \bar{\nu}'_2 = \bar{\nu}_2 - \frac{d^2}{12}, \quad \bar{\nu}'_3 = \bar{\nu}_3 - \frac{d^2}{4}\bar{\nu}_1, \quad \bar{\nu}'_4 = \bar{\nu}_4 - \frac{d^2}{2}\bar{\nu}_2 + \frac{7d^4}{240}.$$

Având în vedere că

$$\bar{\mu}_2 = \bar{\nu}_2 - \bar{\nu}_1^2, \quad \bar{\mu}_3 = \bar{\nu}_3 - 3\bar{\nu}_1\bar{\nu}_2 + 2\bar{\nu}_1^3, \quad \bar{\mu}_4 = \bar{\nu}_4 - 4\bar{\nu}_1\bar{\nu}_3 + 6\bar{\nu}_1^2\bar{\nu}_2 - 3\bar{\nu}_1^4,$$

obținem, de asemenea, momentele centrate corectate

$$\bar{\mu}'_2 = \bar{\mu}_2 - \frac{d^2}{12}, \quad \bar{\mu}'_3 = \bar{\mu}_3, \quad \bar{\mu}'_4 = \bar{\mu}_4 - \frac{d^2}{2}\bar{\mu}_2 + \frac{7d^4}{240}.$$

Corecțiile lui Sheppard se aplică, de regulă, când avem distribuție unimodală și pentru $N > 1000$, $d > \frac{1}{20}\bar{\omega}$.

Exemplul 3.3.32. Să revenim la datele statistice din Exemplul 3.2.10 și să calculăm parametrii statistici definiți până aici.

Printre parametrii statistici care măsoară tendința îi calculăm pe următorii.

Media (aritmetică), care este dată prin

$$\bar{x}_a = \bar{x} = \frac{1}{N} \sum_{k=1}^N x'_k = \frac{1}{70} (1.318 + 3.128 + \dots + 1.966) = 2.27077.$$

Mediana, care se obține după ordonarea crescătoare a datelor statistice. Anume, avem că

$$\begin{aligned}
 &0.842 < 1.155 < 1.179 < 1.318 < 1.355 < 1.403 < 1.493 < 1.546 < 1.548 < 1.553 \\
 &< 1.560 < 1.583 < 1.628 < 1.631 < 1.681 < 1.690 < 1.708 < \boxed{1.802} < 1.855 < 1.864 \\
 &< 1.875 < 1.931 < 1.937 < 1.946 < 1.960 < 1.962 < 1.966 < 1.977 < 2.015 < 2.072 \\
 &< 2.128 < 2.206 < 2.214 < 2.219 < \boxed{\boxed{2.230 < 2.256}} < 2.262 < 2.298 < 2.304 < 2.316 \\
 &< 2.337 < 2.345 < 2.444 < 2.455 < 2.465 < 2.493 < 2.500 < 2.502 < 2.517 < 2.525 \\
 &< 2.636 < 2.637 < \boxed{2.641} < 2.676 < 2.726 < 2.758 < 2.807 < 2.838 < 2.879 < 3.006 \\
 &< 3.093 < 3.128 < 3.156 < 3.281 < 3.394 < 3.426 < 3.460 < 3.537 < 3.852 < 3.972
 \end{aligned}$$

Datele statistice dublu încadrate fiind la mijlocul acestei ordonări (număr par de date statistice) obținem mediana ca fiind

$$\bar{m} = \frac{2.230 + 2.256}{2} = 2.243.$$

Cuartila inferioară se obține din aceeași ordonare a datelor statistice, luând data statistică cu numărul de ordine 18, prima încadrată, adică $\bar{Q}_1 = 1.802$, iar cuartila superioară, luând data statistică cu numărul de ordine 53, a doua încadrată, adică $\bar{Q}_3 = 2.641$.

Printre parametrii statistici care măsoară gradul de împrăștiere îi calculăm pe cei care urmează.

Intervalul de variație este

$$\bar{\omega} = x_{\max} - x_{\min} = 3.972 - 0.842 = 3.130.$$

Abaterea cuartilă sau intervalul intercuartilic se obține prin

$$\bar{Q}_3 - \bar{Q}_1 = 2.641 - 1.802 = 0.839.$$

Deoarece $\bar{Q}_3 - \bar{Q}_1 = 0.839 < 1.565 = \frac{\bar{\omega}}{2}$, avem că distribuția statistică este intens concentrată.

Abaterea cuartilă relativă este

$$\bar{Q}_r = \frac{\bar{Q}_3 - \bar{Q}_1}{\bar{m}} = \frac{0.839}{2.243} = 0.374.$$

Abaterea medie (absolută) are valoarea dată prin

$$\bar{\delta} = \frac{1}{N} \sum_{k=1}^N |x'_k - \bar{x}| =$$

$$\begin{aligned}
&= \frac{1}{70} (|1.318 - 2.271| + |3.128 - 2.271| + \dots + |1.966 - 2.271|) \\
&= 0.5277.
\end{aligned}$$

Primele patru momente sunt

$$\begin{aligned}
\bar{\nu}_1 &= \bar{x} = 2.27077, \\
\bar{\nu}_2 &= \frac{1}{N} \sum_{k=1}^N x_k'^2 = \frac{1}{70} (1.318^2 + 3.128^2 + \dots + 1.966^2) = 5.59435, \\
\bar{\nu}_3 &= \frac{1}{N} \sum_{k=1}^N x_k'^3 = \frac{1}{70} (1.318^3 + 3.128^3 + \dots + 1.966^3) = 14.808, \\
\bar{\nu}_4 &= \frac{1}{N} \sum_{k=1}^N x_k'^4 = \frac{1}{70} (1.318^4 + 3.128^4 + \dots + 1.966^4) = 41.7278.
\end{aligned}$$

Folosind legăturile dintre momentele obișnuite și momentele centrate se obțin acestea din urmă cu formulele

$$\begin{aligned}
\bar{\mu}_2 &= \bar{\nu}_2 - \bar{\nu}_1^2 = 0.437946, \\
\bar{\mu}_3 &= \bar{\nu}_3 - 3\bar{\nu}_1\bar{\nu}_2 + 2\bar{\nu}_1^3 = 0.115571, \\
\bar{\mu}_4 &= \bar{\nu}_4 - 4\bar{\nu}_1\bar{\nu}_3 + 6\bar{\nu}_1^2\bar{\nu}_2 - 3\bar{\nu}_1^4 = 0.540172.
\end{aligned}$$

Se pot calcula apoi coeficienții lui Pearson

$$\bar{\beta}_1 = \frac{\bar{\mu}_3}{\bar{\mu}_2^{3/2}} = 0.159012 \quad (\text{skewness}), \quad \bar{\beta}_2 = \frac{\bar{\mu}_4}{\bar{\mu}_2^2} = 2.81637 \quad (\text{kurtosis}),$$

respectiv coeficienții lui Fisher

$$\bar{\gamma}_1 = \sqrt{\bar{\beta}_1} = 0.3988 \quad (\text{asimetria}), \quad \bar{\gamma}_2 = \bar{\beta}_2 - 3 = -0.18363 \quad (\text{excesul}).$$

Să recalculăm parametrii statistici pe baza datelor statistice sistematizate (grupe).

Media (aritmetică) va fi

$$\bar{x} = \frac{1}{N} \sum_{k=1}^n f_k x_k = \frac{1}{70} (1 \times 0.85 + 2 \times 1.15 + \dots + 2 \times 3.85) = 2.29.$$

Pentru a calcula mediana folosim formula

$$\bar{m} = a_{j-1} + d_j \frac{\frac{N}{2} - F_{j-1}}{f_j}.$$

Aici F_{j-1} este cea mai mare frecvență cumulată (ascendentă) ce nu depășește pe $\frac{N}{2}$, deci $F_5 = 31 < \frac{70}{2} < 46 = F_6$. Astfel se obține $[a_5; a_6) = [2.2; 2.5)$, intervalul median, care are amplitudinea $d_6 = d = 0.3$ și frecvența absolută $f_6 = 15$. Prin urmare rezultă că

$$\overline{m} = 2.2 + 0.3 \cdot \frac{35 - 31}{15} = 2.28.$$

Pentru a calcula cuartila inferioară se determină prima dată intervalul cuartilic inferior $[a_{i-1}; a_i)$ astfel încât F_{i-1} să fie cea mai mare frecvență cumulată ascendentă ce nu depășește pe $\frac{N}{4} = 17.5$. Aceasta este $F_3 = 12$. Prin urmare $[a_3; a_4) = [1.6; 1.9)$, iar $d_4 = d = 0.3$. Se folosește apoi formula

$$\overline{Q}_1 = a_{i-1} + d_i \cdot \frac{\frac{N}{4} - F_{i-1}}{f_i} = 1.6 + 0.3 \cdot \frac{17.5 - 12}{9} = 1.783.$$

În mod analog, pentru cuartila superioară

$$\overline{Q}_3 = 2.5 + 0.3 \cdot \frac{3 \times 17.5 - 46}{10} = 2.695.$$

Intervalul modal, cel cu frecvența maximă, este $[a_5; a_6) = [2.2; 2.5)$. Prin urmare modul va fi calculat cu formula

$$\overline{mo} = a_5 + d_6 \cdot \frac{\Delta f_5}{\Delta f_5 - \Delta f_6} = 2.2 + 0.3 \cdot \frac{15 - 10}{(15 - 10) - (10 - 15)} = 2.35.$$

Abaterea (absolută) este

$$\begin{aligned} \overline{\delta} &= \frac{1}{N} \sum_{k=1}^n f_k |x_k - \overline{x}| \\ &= \frac{1}{70} (1 \times |0.85 - 2.29| + 2 \times |1.15 - 2.29| + \dots + 2 \times |3.85 - 2.29|) \\ &= 0.5297. \end{aligned}$$

Primele patru momente se calculează prin

$$\begin{aligned} \overline{\nu}_1 &= \overline{x} = 2.29, \\ \overline{\nu}_2 &= \frac{1}{N} \sum_{k=1}^n f_k x_k^2 = \frac{1}{70} (1 \times 0.85^2 + 2 \times 1.15^2 + \dots + 1 \times 3.85^2) = 5.68792, \\ \overline{\nu}_3 &= \frac{1}{N} \sum_{k=1}^n f_k x_k^3 = \frac{1}{70} (1 \times 0.85^3 + 2 \times 1.15^3 + \dots + 1 \times 3.85^3) = 15.1467, \end{aligned}$$

$$\bar{\nu}_4 = \frac{1}{N} \sum_{k=1}^n f_k x_k^4 = \frac{1}{70} (1 \times 0.85^4 + 2 \times 1.15^4 + \dots + 1 \times 3.85^4) = 42.8038,$$

iar apoi se obțin momentele centrate

$$\begin{aligned}\bar{\mu}_2 &= \bar{\nu}_2 - \bar{\nu}_1^2 = 0.443829, \\ \bar{\mu}_3 &= \bar{\nu}_3 - 3\bar{\nu}_1\bar{\nu}_2 + 2\bar{\nu}_1^3 = 0.088591, \\ \bar{\mu}_4 &= \bar{\nu}_4 - 4\bar{\nu}_1\bar{\nu}_3 + 6\bar{\nu}_1^2\bar{\nu}_2 - 3\bar{\nu}_1^4 = 0.526822.\end{aligned}$$

Coefficienții lui Pearson sunt

$$\bar{\beta}_1 = \frac{\bar{\mu}_3^2}{\bar{\mu}_2^3} = 0.08977 \quad (\text{skewness}), \quad \bar{\beta}_2 = \frac{\bar{\mu}_4}{\bar{\mu}_2^2} = 2.67444 \quad (\text{kurtosis}),$$

iar coeficienții lui Fisher sunt

$$\bar{\gamma}_1 = \sqrt{\bar{\beta}_1} = 0.2996 \quad (\text{asimetria}), \quad \bar{\gamma}_2 = \bar{\beta}_2 - 3 = -0.3255 \quad (\text{excesul}).$$

Programul 3.3.33. Să scriem un program, care să calculeze o parte din acești parametri statistici din exemplul precedent. Ne vom opri numai la acei parametri, care sunt definiți prin funcții Matlab, dar evident că odată determinați acești parametri, se pot calcula ușor și ceilalți parametri din exemplul precedent.

```
load x.dat -ascii% fisierul x.dat contine datele
                        % primare din exemplul precedent
ma=mean(x);           fprintf('ma=      %6.4f\n',ma)
mg=geomean(x);        fprintf('mg=      %6.4f\n',mg)
mh=harmmean(x);       fprintf('mh=      %6.4f\n',mh)
m=median(x);          fprintf('mediana=%6.4f\n',m)
Q=prctile(x,[25,75]);
                        fprintf('Q1=      %6.4f\n',Q(1)),
                        fprintf('Q3=      %6.4f\n',Q(2))
md=mad(x);            fprintf('md=      %6.4f\n',md)
mu=moment(x,2);       fprintf('mu2=     %6.4f\n',mu)
mu=moment(x,3);       fprintf('mu3=     %6.4f\n',mu)
mu=moment(x,4);       fprintf('mu4=     %6.4f\n',mu)
r=range(x);           fprintf('range=   %6.4f\n',r)
iq=iqr(x);            fprintf('iq=      %6.4f\n',iq)
s=std(x);              fprintf('s=      %6.4f\n',s)
v=var(x);              fprintf('v=      %6.4f\n',v)
sk=skewness(x);       fprintf('sk=     %6.4f\n',sk)
k=kurtosis(x);        fprintf('k=      %6.4f\n',k)
```

Prin executarea programului se obțin rezultatele

```
ma=      2.2708
mg=      2.1725
```

```

mh=      2.0704
mediana=2.2430
Q1=      1.8020
Q3=      2.6410
md=      0.5277
mu2=     0.4379
mu3=     0.1156
mu4=     0.5402
range=   3.1300
iq=      0.8390
s=       0.6666
v=       0.4443
sk=      0.3988
k=       2.8164

```

3.3.16 Funcția `grpstats`

Funcția `grpstats`, din *Statistics toolbox*, face posibilă determinarea unor parametri statistici pentru date clasificate (grupate) după o anumită variabilă de grupare.

Apelul funcției se face prin

```

ma=grpstats(x,g)
[ma,s,fr,nume]=grpstats(x,g)

```

Parametrul `g` este un vector numeric sau de tip caracter, de aceeași lungime cu numărul liniilor matricei `x`, și care permite gruparea datelor conținute de `x`. O astfel de grupare este formată din acele date din `x`, pentru fiecare coloană în parte, pentru care valorile corespunzătoare din `g` sunt aceleași.

Executarea primei instrucțiuni are ca efect calculul mediilor aritmetice pentru fiecare grupă a fiecărei coloane din matricea `x`.

A doua instrucțiune, față de prima instrucțiune, mai calculează abaterile standard `s`, pentru fiecare grupă din fiecare coloană a matricei `x`, precum și numărul datelor `fr` și etichetele `nume`, pentru fiecare grupă.

Programul 3.3.34. Programul care urmează generează 100 de numere aleatoare ce urmează legea uniformă discretă, cu $N = 4$, în vectorul `g`, care va reprezenta un vector de grupare.

În matricea `x` cu 100 de linii și 3 coloane, vor fi generate numere aleatoare ce urmează legea normală $\mathcal{N}(\mu, \sigma)$, unde parametrul $\sigma = 1$, iar parametrul μ , este respectiv 1, 2 și 3 pentru cele trei coloane.

Folosind vectorul de grupare `g`, se calculează valorile medii, abaterile standard, frecvențele absolute și se precizează numele grupelor, pentru coloanele matricei `x`.

```

g=unidrnd(4,100,1);
t=1:3;
t=t(ones(100,1),:);

```

```

x=normrnd(t,1);
[ma,s,fr,nume]=grpstats(x,g);
fprintf('    Mediile pe grupe \n')
disp(ma)
fprintf('    Abaterile standard pe grupe \n')
disp(s)
fprintf('    Frecventele grupelor \n')
disp(fr)
fprintf('    Etichetele grupelor \n')
disp(nume')

```

În urma executării programului se obțin următoarele rezultate

```

Mediile pe grupe
1.2486    1.8376    3.0692
0.6895    2.0323    3.2620
1.0393    1.9893    3.0223
1.0034    1.8834    3.0419

```

```

Abaterile standard pe grupe
0.2024    0.1701    0.2033
0.2076    0.1793    0.2763
0.2051    0.1723    0.1574
0.1875    0.1670    0.1859

```

```

Frecventele grupelor
24    24    24
23    23    23
27    27    27
26    26    26

```

```

Etichetele grupelor
'1'    '2'    '3'    '4'

```

3.3.17 Funcția boxplot

Sistemul Matlab, prin *Statistics toolbox*, dispune de funcția `boxplot`, care se poate apela cu una din următoarele forme:

```

boxplot(x)
boxplot(x,cr)
boxplot(x,cr,'s')
boxplot(x,cr,'s',v)
boxplot(x,cr,'s',v,w)

```

Efectul executării acestor instrucțiuni este reprezentarea grafică prin dreptunghiuri cu prelungiri de segmente pentru fiecare coloană a matricei `x`. Laturile de jos și de sus ale dreptunghiurilor au ordonatele date de cuartilele inferioare și superioare corespunzătoare, medianele fiind și ele marcate la rândul lor prin segmente paralele situate în interioarele dreptunghiurilor, poziționate prin ordonatele date de valorile acestora.

În mijloacele laturilor de jos și de sus ale dreptunghiurilor sunt trasate spre exterior segmente verticale de lungimi date prin

$$\min [\bar{Q}_1 - x_{(1)}, w], \quad \text{respectiv} \quad \min [x_{(n)} - \bar{Q}_3, w],$$

unde avem valoarea implicită $w = \frac{3}{2} (\bar{Q}_3 - \bar{Q}_1)$. Datele care se află înafara acestor prelungiri se numesc *outliers* (erori grosolane) și sunt marcate prin simbolul precizat prin parametrul *s* (implicit *s*=+). Dacă nu există astfel de date într-o parte sau alta, atunci segmentul din partea de jos se termină cu un marcaj.

Dacă *cr*=1, atunci dreptunghiurile sunt crestate pe laturile verticale, având mijloacele crestăturilor în dreptul medianelor, iar deschiderile crestăturilor prezintă estimatii robuste pentru valorile medii teoretice. Valoarea *cr*=0 (valoare implicită) implică faptul că dreptunghiurile nu sunt crestate.

Când *v*=0 (valoare implicită este *v*=1) reprezentările prezentate până aici sunt rotite cu 90° , adică sunt prezentate pe orizontală.

Remarcăm faptul că dacă *x* este un vector, atunci poate apare un parametru suplimentar *g*, imediat după parametrul *x*. Parametrul *g* este un vector numeric sau de tip caracter, de aceeași lungime cu *x*, și care permite gruparea datelor conținute de *x*. O astfel de grupare este formată din acele date din *x*, pentru care valorile corespunzătoare din *g* sunt aceleași. Funcția *boxplot* în acest caz va produce dreptunghiuri pentru fiecare grupare în parte.

Programul 3.3.35. Programul următor generează *n* vectori aleatori, care urmează legea normală bidimensională, în matricea *X* de tipul $n \times 2$, după care, folosind funcția *boxplot*, reprezintă grafic datele pentru prima variabilă și respectiv pentru a doua variabilă.

```
clear, mu(1)=input('m1='); mu(2)=input('m2=');
v(1,1)=input('sigma1^2=');
v(2,2)=input('sigma2^2=');
v(1,2)=input('Cov(X,Y)='); v(2,1) =v(1,2);
if det(v) <= 0
    error('Matricea v nu e pozitiv definita!')
end
n=input('n='); X=mvnrnd(mu,v,n); boxplot(X)
xlabel('Numarul variabilelor'),
ylabel('Valorile variabilelor')
```

În urma executării programului, cu valorile *mu*₁=10, *mu*₂=10, $v = \begin{pmatrix} 2 & -2 \\ -2 & 3 \end{pmatrix}$ și *n*=50, se obține Figura 3.11.

3.4 Corelație și regresie

Corelația (într-un sens larg) poate fi înțeleasă ca legătura care există între o caracteristică dependentă și una sau mai multe caracteristici independente, iar *regresia* este

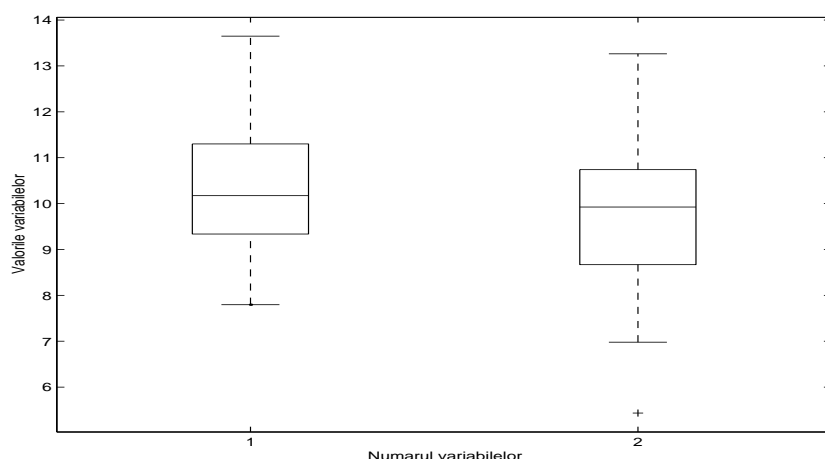


Figura 3.11: Reprezentare de date cu boxplot

$X \setminus Y$	y_1	y_2	\dots	y_n	
x_1	f_{11}	f_{12}	\dots	f_{1n}	$f_{1.}$
x_2	f_{21}	f_{22}	\dots	f_{2n}	$f_{2.}$
\vdots	\vdots	\vdots		\vdots	\vdots
x_m	f_{m1}	f_{m2}	\dots	f_{mn}	$f_{m.}$
	$f_{.1}$	$f_{.2}$	\dots	$f_{.n}$	$f_{..} = N$

Tabelul 3.2: Tabel de corelație

metoda prin care se stabilește această legătură.

Să considerăm caracteristicile cantitative X și Y relative la colectivitatea \mathcal{C} . Datele statistice primare sunt (x'_k, y'_k) , $k = \overline{1, N}$, care după grupare sunt prezentate în Tabelul de corelație 3.2. unde f_{ij} reprezintă frecvența absolută a clasei (x_i, y_j) . De asemenea, avem că

$$\sum_{j=1}^n f_{ij} = f_{i.}, \quad \sum_{i=1}^m f_{ij} = f_{.j}, \quad \sum_{i=1}^m f_{i.} = \sum_{j=1}^n f_{.j} = f_{..} = N.$$

Definiția 3.4.1. Numim moment de ordinul (k_1, k_2) al distribuției statistice a caracteristicii bidimensionale (X, Y) , valoarea numerică

$$\bar{\nu}_{k_1 k_2} = \frac{1}{N} \sum_{i=1}^N x_i'^{k_1} y_i'^{k_2} = \frac{1}{N} \sum_{i=1}^m \sum_{j=1}^n f_{ij} x_i^{k_1} y_j^{k_2} = \sum_{i=1}^m \sum_{j=1}^n p_{ij} x_i^{k_1} y_j^{k_2},$$

unde $p_{ij} = \frac{f_{ij}}{N}$ este frecvența relativă a clasei (x_i, y_j) .

Definiția 3.4.2. Numim moment centrat de ordinul (k_1, k_2) al distribuției statistice bidimensionale (X, Y) , valoarea numerică

$$\begin{aligned}\bar{\mu}_{k_1 k_2} &= \frac{1}{N} \sum_{i=1}^N (x'_i - \bar{x})^{k_1} (y'_i - \bar{y})^{k_2} = \frac{1}{N} \sum_{i=1}^m \sum_{j=1}^n f_{ij} (x_i - \bar{x})^{k_1} (y_j - \bar{y})^{k_2} \\ &= \sum_{i=1}^m \sum_{j=1}^n p_{ij} (x_i - \bar{x})^{k_1} (y_j - \bar{y})^{k_2},\end{aligned}$$

unde

$$\bar{x} = \bar{\nu}_{10} = \frac{1}{N} \sum_{i=1}^m f_{i.} x_i, \quad \bar{y} = \bar{\nu}_{01} = \frac{1}{N} \sum_{j=1}^n f_{.j} y_j.$$

Observația 3.4.3. Remarcăm, printre momentele centrate, dispersiile pentru distribuțiile statistice ale caracteristicilor X și respectiv Y , anume

$$\bar{\sigma}_X^2 = \bar{\mu}_{20} = \frac{1}{N} \sum_{i=1}^m f_{i.} (x_i - \bar{x})^2, \quad \bar{\sigma}_Y^2 = \bar{\mu}_{02} = \frac{1}{N} \sum_{j=1}^n f_{.j} (y_j - \bar{y})^2.$$

Definiția 3.4.4. Numim coeficient de corelație (al lui Pearson) al distribuției statistice bidimensionale (X, Y) , raportul

$$\bar{r} = \frac{\bar{\mu}_{11}}{\sqrt{\bar{\mu}_{20}} \sqrt{\bar{\mu}_{02}}} = \frac{\bar{\nu}_{11} - \bar{\nu}_{10} \bar{\nu}_{01}}{\bar{\sigma}_X \bar{\sigma}_Y} = \frac{\bar{\nu}_{11} - \bar{x} \bar{y}}{\bar{\sigma}_X \bar{\sigma}_Y}.$$

Observația 3.4.5. Folosind datele statistice negrupate avem formulele de calcul pentru coeficientul \bar{r} de corelație

$$\bar{r} = \frac{\sum_{i=1}^N (x'_i - \bar{x}) (y'_i - \bar{y})}{\sqrt{\sum_{i=1}^N (x'_i - \bar{x})^2} \sqrt{\sum_{i=1}^N (y'_i - \bar{y})^2}}$$

sau

$$\bar{r} = \frac{\sum_{i=1}^N x'_i y'_i - \frac{1}{N} \left(\sum_{i=1}^N x'_i \right) \left(\sum_{i=1}^N y'_i \right)}{\sqrt{\sum_{i=1}^N x_i'^2 - \frac{1}{N} \left(\sum_{i=1}^N x'_i \right)^2} \sqrt{\sum_{i=1}^N y_i'^2 - \frac{1}{N} \left(\sum_{i=1}^N y'_i \right)^2}}.$$

Observația 3.4.6. Folosind inegalitatea Cauchy–Schwarz–Buniakovski

$$\left[\sum_{i=1}^N (x'_i - \bar{x}) (y'_i - \bar{y}) \right]^2 \leq \left[\sum_{i=1}^N (x'_i - \bar{x})^2 \right] \left[\sum_{i=1}^N (y'_i - \bar{y})^2 \right]$$

se stabilește că $|\bar{r}| \leq 1$. De asemenea, dacă $|\bar{r}| = 1$, atunci între cele două caracteristici există o legătură liniară, și invers. Când $\bar{r} = 0$, spunem că cele două caracteristici sunt (liniar) *necorelate*. Prin urmare, coeficientul de corelație poate fi folosit pentru a stabili dacă între cele două caracteristici există sau nu o legătură liniară. În cazul în care caracteristica bidimensională (X, Y) urmează legea normală bidimensională $\bar{r} = 0$ implică faptul că cele două caracteristici sunt independente.

Exemplul 3.4.7. Să calculăm câțiva din parametri definiți pentru datele statistice de la Exemplul 3.2.20.

Pentru a calcula coeficientul de corelație, calculăm prima dată valorile medii ale caracteristicilor X și Y . Remarcăm, pentru aceasta, că distribuțiile statistice ale caracteristicilor X și Y sunt respectiv

$$X \begin{pmatrix} 126 & 127 & 128 & 129 & 130 & 131 & 132 & 133 & 134 & 135 & 136 & 137 \\ 6 & 9 & 21 & 29 & 70 & 152 & 64 & 38 & 18 & 8 & 6 & 4 \end{pmatrix},$$

$$Y \begin{pmatrix} 24 & 25 & 26 & 27 & 28 & 29 & 30 & 31 & 32 & 33 & 34 & 35 \\ 2 & 9 & 16 & 30 & 84 & 106 & 89 & 45 & 24 & 11 & 7 & 2 \end{pmatrix}.$$

Obținem astfel că

$$\bar{x} = \frac{1}{425} (6 \cdot 126 + 9 \cdot 127 + \dots + 4 \cdot 137) = 131.054 \text{ cm},$$

$$\bar{y} = \frac{1}{425} (2 \cdot 24 + 9 \cdot 25 + \dots + 2 \cdot 35) = 29.2447 \text{ kg}.$$

De asemenea, calculând dispersiile caracteristicilor se obține

$$\begin{aligned} \bar{\sigma}_X^2 &= \frac{1}{N} \left[\sum_{i=1}^m f_i x_i^2 - \frac{1}{N} \left(\sum_{i=1}^m f_i x_i \right)^2 \right] \\ &= \frac{1}{425} \left[(6 \cdot 126^2 + 9 \cdot 127^2 + \dots + 4 \cdot 137^2) - \frac{1}{425} (6 \cdot 126 + \dots + 4 \cdot 137)^2 \right] \\ &= \frac{1}{425} \left[7300920 - \frac{1}{425} \cdot 55698^2 \right] = 3.45354, \end{aligned}$$

$$\begin{aligned}
\bar{\sigma}_Y^2 &= \frac{1}{N} \left[\sum_{j=1}^n f_{.j} y_j^2 - \frac{1}{N} \left(\sum_{j=1}^n f_{.j} y_j \right)^2 \right] \\
&= \frac{1}{425} \left[(2 \cdot 24^2 + 9 \cdot 25^2 + \dots + 2 \cdot 35^2) - \frac{1}{425} (2 \cdot 24 + 9 \cdot 25 + \dots + 2 \cdot 35)^2 \right] \\
&= \frac{1}{425} \left[364907 - \frac{1}{425} \cdot 12429^2 \right] = 3.35188.
\end{aligned}$$

Covarianța dintre X și Y se calculează cu formula

$$\begin{aligned}
\text{Cov}(X, Y) &= \frac{1}{N} \left[\sum_{i=1}^m \sum_{j=1}^n f_{ij} x_i y_j - \frac{1}{N} \left(\sum_{i=1}^m f_{i.} x_i \right) \left(\sum_{j=1}^n f_{.j} y_j \right) \right] \\
&= \frac{1}{425} \left[1629961 - \frac{1}{425} \cdot 55698 \cdot 12429 \right] = 2.56323.
\end{aligned}$$

Astfel se ajunge la coeficientul de corelație

$$\bar{r} = \frac{\text{Cov}(X, Y)}{\sqrt{\bar{\sigma}_X^2 \bar{\sigma}_Y^2}} = \frac{2.56323}{\sqrt{3.45354 \times 3.35188}} = 0.753374.$$

Programul 3.4.8. Programul Matlab, care urmează, efectuează aceste calcule

```

f=[1,1,3,1,zeros(1,8);
  0,3,4,1,0,1,zeros(1,6);
  1,4,5,7,1,2,1,zeros(1,5);
  0,1,2,6,9,6,4,1,zeros(1,4);
  zeros(1,2),1,7,36,17,6,2,1,zeros(1,3);
  zeros(1,2),1,6,27,56,39,18,4,1,zeros(1,2);
  zeros(1,3),2,10,16,26,7,2,1,zeros(1,2);
  zeros(1,4),1,7,10,11,6,2,1,0;
  zeros(1,5),1,2,4,7,3,1,0; zeros(1,6),1,2,3,1,1,0;
  zeros(1,8),1,2,2,1; zeros(1,9),1,2,1];
x=[126:137]; y=[24:35]; fx=sum(f'); fy=sum(f);
xa=sum(fx.*x)/(sum(fx)); ya=sum(fy.*y)/(sum(fy));
vx=(sum(fx.*x.^2)-(sum(fx.*x))^2/sum(fx))/sum(fx);
vy=(sum(fy.*y.^2)-(sum(fy.*y))^2/sum(fy))/sum(fy);
cov=(x*f*y'-sum(fx.*x)*sum(fy.*y)/sum(fx))/sum(fx);
r=cov/(sqrt(vx*vy));
fprintf(' xa=%8.3f\n ya=%8.4f\n',xa,ya)
fprintf(' vx=%8.5f\n vy=%8.5f\n',vx,vy)
fprintf(' cov=%7.5f\n r=%9.6f',cov,r)

```

iar rezultatele sunt următoarele:

```

xa= 131.054
ya= 29.2447
vx= 3.45354

```

$v_Y = 3.35188$
 $\text{cov} = 2.56323$
 $r = 0.753374$

Definiția 3.4.9. Numim valoare medie condiționată a distribuției statistice a caracteristicii Y în raport cu $X = x_j$, valoarea numerică

$$\bar{y}_i = \bar{y}(x_i) = \frac{1}{f_{i\cdot}} \sum_{j=1}^n f_{ij} y_j = \sum_{j=1}^n \frac{f_{ij}}{f_{i\cdot}} y_j, \quad i = \overline{1, m},$$

și respectiv valoare medie condiționată a distribuției statistice a caracteristicii X în raport cu $Y = y_j$, valoarea numerică

$$\bar{x}_j = \bar{x}(y_j) = \frac{1}{f_{\cdot j}} \sum_{i=1}^m f_{ij} x_i = \sum_{i=1}^m \frac{f_{ij}}{f_{\cdot j}} x_i, \quad j = \overline{1, n}.$$

Definiția 3.4.10. Curba de ecuație $y = f(x)$ pe care se situează punctele de coordonate (x_i, \bar{y}_i) , $i = \overline{1, m}$, se numește curba de regresie a lui Y în raport cu X , iar curba de ecuație $x = g(y)$ pe care se situează punctele de coordonate (\bar{x}_j, y_j) , $j = \overline{1, n}$, se numește curba de regresie a lui X în raport cu Y .

Definiția 3.4.11. Numim dispersie condiționată a distribuției statistice a caracteristicii Y în raport cu $X = x_i$, valoarea numerică

$$\bar{\sigma}_{Y|x_i}^2 = \frac{1}{f_{i\cdot}} \sum_{j=1}^n f_{ij} (y_j - \bar{y}_i)^2 = \sum_{j=1}^n \frac{f_{ij}}{f_{i\cdot}} (y_j - \bar{y}_i)^2, \quad i = \overline{1, m},$$

și respectiv dispersie condiționată a distribuției statistice a caracteristicii X în raport cu $Y = y_j$, valoarea numerică

$$\bar{\sigma}_{X|y_j}^2 = \frac{1}{f_{\cdot j}} \sum_{i=1}^m f_{ij} (x_i - \bar{x}_j)^2 = \sum_{i=1}^m \frac{f_{ij}}{f_{\cdot j}} (x_i - \bar{x}_j)^2, \quad j = \overline{1, n}.$$

Definiția 3.4.12. Numim dispersia condiționată a distribuției statistice a lui Y în raport cu distribuția statistică a lui X , valoarea numerică

$$\bar{\sigma}_{Y|X}^2 = \frac{1}{N} \sum_{i=1}^m f_{i\cdot} \bar{\sigma}_{Y|x_i}^2 = \sum_{i=1}^m p_{i\cdot} \bar{\sigma}_{Y|x_i}^2$$

și respectiv dispersia condiționată a distribuției statistice a lui X în raport cu distribuția statistică a lui Y , valoarea numerică

$$\bar{\sigma}_{X|Y}^2 = \frac{1}{N} \sum_{j=1}^n f_{\cdot j} \bar{\sigma}_{X|y_j}^2 = \sum_{j=1}^n p_{\cdot j} \bar{\sigma}_{X|y_j}^2,$$

unde $p_{i.} = \frac{f_{i.}}{N}$ este frecvența relativă a clasei x_i , iar $p_{.j} = \frac{f_{.j}}{N}$ este frecvența relativă a clasei y_j .

Proprietatea 3.4.13. Dispersiile condiționate definite mai înainte satisfac relațiile

$$\overline{\sigma}_Y^2 = \overline{\sigma}_{Y|X}^2 + \overline{\sigma}_{\overline{Y}|X}^2, \quad \overline{\sigma}_X^2 = \overline{\sigma}_{X|Y}^2 + \overline{\sigma}_{\overline{X}|Y}^2,$$

unde

$$\overline{\sigma}_{\overline{Y}|X}^2 = \frac{1}{N} \sum_{i=1}^m f_{i.} (\overline{y}_i - \overline{y})^2, \quad \overline{\sigma}_{\overline{X}|Y}^2 = \frac{1}{N} \sum_{j=1}^n f_{.j} (\overline{x}_j - \overline{x})^2,$$

adică sunt dispersiile valorilor medii condiționate.

Demonstrație. Se pornește de la identitatea

$$(y_j - \overline{y})^2 = (y_j - \overline{y}_i)^2 + (\overline{y}_i - \overline{y})^2 + 2(y_j - \overline{y}_i)(\overline{y}_i - \overline{y}),$$

care se înmulțește cu $\frac{f_{ij}}{f_{i.}}$ și se însumează după $j = \overline{1, n}$. Astfel se obține

$$\begin{aligned} \sum_{j=1}^n \frac{f_{ij}}{f_{i.}} (y_j - \overline{y})^2 &= \sum_{j=1}^n \frac{f_{ij}}{f_{i.}} (y_j - \overline{y}_i)^2 + \sum_{j=1}^n \frac{f_{ij}}{f_{i.}} (\overline{y}_i - \overline{y})^2 + \\ &\quad + 2 \sum_{j=1}^n \frac{f_{ij}}{f_{i.}} (y_j - \overline{y}_i) (\overline{y}_i - \overline{y}) \\ &= \overline{\sigma}_{Y|x_i}^2 + (\overline{y}_i - \overline{y})^2 + 2(\overline{y}_i - \overline{y})(\overline{y}_i - \overline{y}_i) = \overline{\sigma}_{Y|x_i}^2 + (\overline{y}_i - \overline{y})^2. \end{aligned}$$

Reținem extremitățile acestui șir de egalități, le înmulțim cu $\frac{f_{i.}}{N}$ și le însumăm după $i = \overline{1, m}$. Rezultă în acest mod

$$\begin{aligned} \overline{\sigma}_Y^2 &= \frac{1}{N} \sum_{j=1}^n f_{.j} (y_j - \overline{y})^2 = \frac{1}{N} \sum_{i=1}^m \sum_{j=1}^n f_{ij} (y_j - \overline{y})^2 \\ &= \frac{1}{N} \sum_{i=1}^m f_{i.} \overline{\sigma}_{Y|x_i}^2 + \frac{1}{N} \sum_{i=1}^m f_{i.} (\overline{y}_i - \overline{y})^2 = \overline{\sigma}_{Y|X}^2 + \overline{\sigma}_{\overline{Y}|X}^2, \end{aligned}$$

de unde se obține prima relație. Analog se obține și cealaltă relație. \square

Definiția 3.4.14. Numim raport de corelație al distribuției statistice a caracteristicii Y față de distribuția statisticii X , valoarea numerică

$$\overline{\eta}_{Y|X} = \sqrt{1 - \frac{\overline{\sigma}_{Y|X}^2}{\overline{\sigma}_Y^2}} = \sqrt{\frac{\overline{\sigma}_{\overline{Y}|X}^2}{\overline{\sigma}_Y^2}}$$

și în mod analog avem

$$\bar{\eta}_{X|Y} = \sqrt{1 - \frac{\bar{\sigma}_{X|Y}^2}{\bar{\sigma}_X^2}} = \sqrt{\frac{\bar{\sigma}_{X|Y}^2}{\bar{\sigma}_X^2}}.$$

Observația 3.4.15. Când $\bar{y}_1 = \bar{y}_2 = \dots = \bar{y}_m$, atunci $\bar{y} = \bar{y}_i$, prin urmare $\bar{\sigma}_{Y|X}^2 = 0$ ($\bar{\eta}_{Y|X} = 0$), ceea ce înseamnă absența dependenței în medie. Pe de altă parte, $\bar{\eta}_{Y|X} = 1$ când $\bar{\sigma}_{Y|x_i}^2 = 0$, $i = \overline{1, m}$, adică pentru fiecare clasă x_i valorile caracteristicii Y sunt aceleași.

În cazul în care, relativ la caracteristica X , există numai două clase avem

$$\bar{y} = \frac{f_{1.}}{N} \bar{y}_1 + \frac{f_{2.}}{N} \bar{y}_2.$$

Dacă se înlocuiește această expresie în formula

$$\bar{\sigma}_{Y|X}^2 = \frac{1}{N} \left[f_{1.} (\bar{y}_1 - \bar{y})^2 + f_{2.} (\bar{y}_2 - \bar{y})^2 \right]$$

se obține că

$$\bar{\sigma}_{Y|X}^2 = \frac{f_{1.} f_{2.}}{N^2} (\bar{y}_1 - \bar{y}_2)^2,$$

deci

$$\bar{\eta}_{Y|X}^2 = \frac{f_{1.} f_{2.} (\bar{y}_1 - \bar{y}_2)^2}{N^2 \bar{\sigma}_Y^2}.$$

Aplicația 3.4.16. (*Determinarea curbelor de regresie*). Prezentăm în continuare metoda celor mai mici pătrate pentru determinarea ecuațiilor curbelor de regresie.

Presupunem că din reprezentarea în plan a punctelor (x_i, \bar{y}_i) , $i = \overline{1, m}$, curba de regresie a lui Y în raport cu X este de forma $y = y(x) = f(x; a_1, a_2, \dots, a_s)$. Vom determina parametrii a_k , $k = \overline{1, s}$, astfel încât

$$\begin{aligned} S(a_1, a_2, \dots, a_s) &= \sum_{i=1}^N \left[y'_i - y(x'_i) \right]^2 = \sum_{i=1}^m \sum_{j=1}^n f_{ij} \left[y_j - y(x_i) \right]^2 \\ &= \sum_{i=1}^m \sum_{j=1}^n f_{ij} \left[y_j - f(x_i; a_1, a_2, \dots, a_s) \right]^2 \end{aligned}$$

să fie minimă.

Punctul de minim $(\bar{a}_1, \bar{a}_2, \dots, \bar{a}_s)$ al funcției S se obține prin rezolvarea *sistemului normal de ecuații*, rezultat din

$$\frac{\partial S}{\partial a_k} = -2 \sum_{i=1}^m \sum_{j=1}^n f_{ij} \left[y_j - f(x_i; a_1, a_2, \dots, a_s) \right] \frac{\partial f(x_i; a_1, a_2, \dots, a_s)}{\partial a_k} = 0$$

pentru $k = \overline{1, s}$. Ecuația curbei de regresie va fi

$$y = f(x; \bar{a}_1, \bar{a}_2, \dots, \bar{a}_s).$$

La fel se determină și ecuația curbei de regresie a lui X în raport cu Y .

Aplicația 3.4.17. (*Drepte de regresie*). Considerăm cazul liniar, adică ecuația curbei de regresie este $y = y(x) = ax + b$. Expresia minimizată este

$$S(a, b) = \sum_{i=1}^m \sum_{j=1}^n f_{ij} (y_j - ax_i - b)^2,$$

care conduce la sistemul normal de ecuații

$$\begin{cases} \left(\sum_{i=1}^m \sum_{j=1}^n f_{ij} x_i^2 \right) a + \left(\sum_{i=1}^m \sum_{j=1}^n f_{ij} x_i \right) b = \sum_{i=1}^m \sum_{j=1}^n f_{ij} x_i y_j \\ \left(\sum_{i=1}^m \sum_{j=1}^n f_{ij} x_i \right) a + \left(\sum_{i=1}^m \sum_{j=1}^n f_{ij} \right) b = \sum_{i=1}^m \sum_{j=1}^n f_{ij} y_j, \end{cases}$$

sau

$$\begin{cases} \bar{\nu}_{20}a + \bar{\nu}_{10}b = \bar{\nu}_{11} \\ \bar{\nu}_{10}a + \bar{\nu}_{00}b = \bar{\nu}_{01}. \end{cases}$$

Soluția acestui sistem liniar este

$$\begin{aligned} \bar{a} &= \frac{\bar{\nu}_{11} - \bar{\nu}_{10}\bar{\nu}_{01}}{\bar{\nu}_{20} - \bar{\nu}_{10}^2} = \frac{\bar{\nu}_{11} - \bar{\nu}_{10}\bar{\nu}_{01}}{\sqrt{\bar{\nu}_{20} - \bar{\nu}_{10}^2} \sqrt{\bar{\nu}_{02} - \bar{\nu}_{01}^2}} \frac{\sqrt{\bar{\nu}_{02} - \bar{\nu}_{01}^2}}{\sqrt{\bar{\nu}_{20} - \bar{\nu}_{10}^2}} = \bar{r} \frac{\bar{\sigma}_Y}{\bar{\sigma}_X}, \\ \bar{b} &= \bar{\nu}_{01} - \bar{a}\bar{\nu}_{10} = \bar{y} - \bar{a}\bar{x}. \end{aligned}$$

Se obține în acest fel ecuația dreptei de regresie a lui Y în raport cu X ca fiind

$$y - \bar{y} = \bar{r} \frac{\bar{\sigma}_Y}{\bar{\sigma}_X} (x - \bar{x}).$$

În mod analog se ajunge la ecuația dreptei de regresie a lui X în raport cu Y

$$(3.4.1) \quad x - \bar{x} = \bar{r} \frac{\bar{\sigma}_X}{\bar{\sigma}_Y} (y - \bar{y}).$$

Punctul de intersecție al celor două drepte de regresie are coordonatele (\bar{x}, \bar{y}) și se numește *centrul de greutate* al distribuției statistice a caracteristicii bidimensionale (X, Y) .

Dacă se notează

$$\bar{a}_{Y|X} = \bar{r} \frac{\bar{\sigma}_Y}{\bar{\sigma}_X}$$

coeficientul unghiular al dreptei de regresie a lui Y în raport cu X (numit *coeficientul de regresie* al lui Y în raport cu X) și cu

$$\bar{a}_{X|Y} = \bar{r} \frac{\bar{\sigma}_X}{\bar{\sigma}_Y}$$

coeficientul de regresie al lui X în raport cu Y , atunci

$$|\bar{r}| = \sqrt{\bar{a}_{Y|X} \bar{a}_{X|Y}}.$$

De asemenea se constată că $\text{sign}(\bar{a}_{X|Y}) = \text{sign}(\bar{a}_{Y|X})$. Se arată că unghiul α format de cele două drepte de regresie este dat prin relația

$$\text{tg } \alpha = \frac{1 - \bar{r}^2}{\bar{r}^2} \frac{\bar{\sigma}_X \bar{\sigma}_Y}{\bar{\sigma}_X^2 + \bar{\sigma}_Y^2}.$$

Folosind această relație se pot trage următoarele concluzii:

a) dacă $|\bar{r}| = 1$ atunci $\alpha = 0$, deci dreptele de regresie se confundă, cu specificația că pentru $\bar{r} = -1$ dreptele au panta (coeficientul unghiular) negativă, iar pentru $\bar{r} = 1$ panta este pozitivă.

b) dacă X și Y sunt independente atunci $\bar{r} = 0$, deci $\alpha = \frac{\pi}{2}$ (dreptele de regresie sunt perpendiculare).

Observația 3.4.18. Alte tipuri de curbe de regresie mai des întâlnite și care pot fi liniarizate sunt:

1. $y = ab^x$ (exponențială), care prin logaritmare se liniarizează după cum urmează $\log y = \log a + x \log b$, luând $z = \log y$, $A = \log a$, $B = \log b$;
2. $y = \frac{a}{x} + b$ (hiperbolică), care se liniarizează dacă se notează $z = \frac{1}{x}$;
3. $\frac{1}{y} = \frac{a}{x} + b$ sau $y = \frac{1}{\frac{a}{x} + b}$, care se liniarizează dacă se notează $u = \frac{1}{x}$, $v = \frac{1}{y}$;

4. $y = a \log x + b$ (logaritmică), care se liniarizează dacă se notează $z = \log x$;
5. $y = be^{ax}$ (exponențială), care prin logaritmare se liniarizează $\ln y = \ln b + ax$, luând $z = \ln y$;
6. $y = be^{\frac{a}{x}}$, care prin logaritmare se liniarizează $\ln y = \ln b + \frac{a}{x}$, luând $u = \frac{1}{x}$, $v = \ln y$;
7. $y = bx^a$, care prin logaritmare se liniarizează, $\log y = \log b + a \log x$, luând $u = \log x$, $v = \log y$;
8. $\frac{1}{y} = ae^{-x} + b$ sau $y = \frac{1}{ae^{-x} + b}$, care se liniarizează dacă se ia $u = e^{-x}$ și $v = \frac{1}{y}$.

Această ultimă legătură este un caz particular al *legăturii logistice* definită prin

$$y = \frac{1}{ae^{-cx} + b}.$$

Printre legăturile ce nu pot fi liniarizate amintim:

$$y = ax^b + c \log x,$$

$$y = ax^b e^{cx},$$

$$y = a + bx + ce^{dx},$$

a doua putând fi adusă la forma polinomială.

Exemplul 3.4.19. Să considerăm din nou datele statistice de la Exemplul 3.2.20, să calculăm noile caracteristici numerice definite și să determinăm ecuațiile dreptelor de regresie.

Valorile medii condiționate și dispersiile condiționate ale lui Y în raport cu X se calculează cu formulele

$$\bar{y}_i = \frac{1}{f_{i.}} \sum_{j=1}^n f_{ij} y_j, \quad i = \overline{1, m},$$

și

$$\sigma_{Y|x_i}^2 = \frac{1}{f_{i.}} \sum_{j=1}^n f_{ij} (y_j - \bar{y}_i)^2, \quad i = \overline{1, m}.$$

În mod analog se obțin valorile medii condiționate și dispersiile condiționate ale lui X în raport cu Y .

Aceste valori calculate sunt trecute în tabelele următoare.

i	\bar{y}_i	$\bar{\sigma}_{Y x_i}^2$		j	\bar{x}_j	$\bar{\sigma}_{X y_j}^2$
1	25.67	0.889		1	127.00	1.000
2	26.11	1.432		2	127.56	0.691
3	26.62	2.141		3	127.81	1.902
4	28.14	1.843		4	129.43	2.112
5	28.43	1.045		5	130.46	0.844
6	29.32	1.363	și	6	130.94	1.204
7	29.56	1.340		7	131.44	1.280
8	30.63	1.706		8	132.00	1.778
9	31.67	1.444		9	133.12	2.276
10	31.88	1.359		10	134.09	2.992
11	33.50	0.917		11	135.43	1.959
12	34.00	0.500		12	136.50	0.250

Dispersia condiționată a lui Y în raport cu X va fi

$$\bar{\sigma}_{Y|X}^2 = \frac{1}{N} \sum_{i=1}^m f_i \bar{\sigma}_{Y|x_i}^2 = \frac{1}{425} (6 \times 0.889 + 9 \times 1.432 + \dots + 4 \times 0.500) = 1.39279.$$

În mod analog dispersia condiționată a lui X în raport cu Y este

$$\bar{\sigma}_{X|Y}^2 = \frac{1}{N} \sum_{j=1}^n f_j \bar{\sigma}_{X|y_j}^2 = \frac{1}{425} (2 \times 1.000 + 9 \times 0.691 + \dots + 2 \times 0.250) = 1.40291.$$

Dispersia condiționată a lui Y în raport cu X se poate calcula și cu formula

$$\bar{\sigma}_{Y|X}^2 = \bar{\sigma}_Y^2 - \bar{\sigma}_{\bar{Y}|X}^2,$$

unde

$$\bar{\sigma}_{\bar{Y}|X}^2 = \frac{1}{N} \sum_{i=1}^m f_i (\bar{y}_i - \bar{y})^2$$

este dispersia valorilor medii. Anume, avem că

$$\bar{\sigma}_{\bar{Y}|X}^2 = \frac{1}{425} \left[6 (25.67 - 29.2447)^2 + \dots + 4 (34.00 - 29.2447)^2 \right] = 1.95909,$$

de unde

$$\bar{\sigma}_{Y|X}^2 = 3.35188 - 1.95909 = 1.39279.$$

În mod analog,

$$\bar{\sigma}_{X|Y}^2 = \bar{\sigma}_X^2 - \bar{\sigma}_{\bar{X}|Y}^2,$$

unde

$$\overline{\sigma}_{X|Y}^2 = \frac{1}{N} \sum_{j=1}^n f_{.j} (\bar{x}_j - \bar{x})^2,$$

adică

$$\overline{\sigma}_{X|Y}^2 = \frac{1}{425} \left[2(127 - 131.054)^2 + \dots + 2(136.5 - 131.054)^2 \right] = 2.05063.$$

În acest fel se obține că

$$\overline{\sigma}_{X|Y}^2 = 3.45354 - 2.05063 = 1.40291.$$

Raportul de corelație al lui Y față de X este

$$\bar{\eta}_{Y|X} = \sqrt{1 - \frac{\overline{\sigma}_{Y|X}^2}{\overline{\sigma}_Y^2}} = \sqrt{\frac{\overline{\sigma}_{Y|X}^2}{\overline{\sigma}_Y^2}} = \sqrt{\frac{1.95909}{3.35188}} = 0.764509,$$

iar raportul de corelație al lui X față de Y este

$$\bar{\eta}_{X|Y} = \sqrt{1 - \frac{\overline{\sigma}_{X|Y}^2}{\overline{\sigma}_X^2}} = \sqrt{\frac{\overline{\sigma}_{X|Y}^2}{\overline{\sigma}_X^2}} = \sqrt{\frac{2.05063}{3.45354}} = 0.770568.$$

Ecuatiile dreptelor de regresie respectiv a lui Y în raport cu X și a lui X în raport cu Y au ecuațiile

$$y - \bar{y} = \bar{r} \frac{\overline{\sigma}_Y}{\overline{\sigma}_X} (x - \bar{x}), \quad x - \bar{x} = \bar{r} \frac{\overline{\sigma}_X}{\overline{\sigma}_Y} (y - \bar{y}).$$

Având în vedere că

$$\begin{aligned} \bar{r} \frac{\overline{\sigma}_Y}{\overline{\sigma}_X} &= 0.753374 \sqrt{\frac{3.35188}{3.45354}} = 0.742203, \\ \bar{r} \frac{\overline{\sigma}_X}{\overline{\sigma}_Y} &= 0.753374 \sqrt{\frac{3.45354}{3.35188}} = 0.764713, \end{aligned}$$

avem ecuațiile

$$y - 29.2447 = 0.742203 (x - 131.054), \quad x - 131.054 = 0.764713 (y - 29.2447).$$

Centrul de greutate al distribuției statistice (X, Y) este punctul de coordonate $(\bar{x}; \bar{y}) = (131.05; 29.24)$, adică punctul de intersecție a celor două drepte de regresie.

Coeficientul de regresie al lui Y în raport cu X este coeficientul unghiular al drepte de regresie a lui Y în raport cu X , adică

$$\bar{a}_{Y|X} = \bar{r} \frac{\bar{\sigma}_Y}{\bar{\sigma}_X} = 0.742203.$$

În mod analog, coeficientul de regresie al lui X în raport cu Y este

$$\bar{a}_{X|Y} = \bar{r} \frac{\bar{\sigma}_X}{\bar{\sigma}_Y} = 0.764713.$$

Pentru determinarea unghiului α format de cele două drepte de regresie avem formula

$$\operatorname{tg} \alpha = \frac{1 - \bar{r}^2}{\bar{r}^2} \frac{\bar{\sigma}_X \bar{\sigma}_Y}{\bar{\sigma}_X^2 + \bar{\sigma}_Y^2},$$

adică

$$\operatorname{tg} \alpha = \frac{1 - 0.5675}{0.5675} \cdot \frac{\sqrt{3.45354 \times 3.35188}}{3.45354 + 3.35188} = 0.381,$$

de unde $\alpha \cong 21^0 = 0.364 \text{ rad.}$

Programul 3.4.20. Calculele de mai sus se pot face cu programul Matlab, care urmează:

```
f=[1,1,3,1,zeros(1,8);
    0,3,4,1,0,1,zeros(1,6);
    1,4,5,7,1,2,1,zeros(1,5);
    0,1,2,6,9,6,4,1,zeros(1,4);
    zeros(1,2),1,7,36,17,6,2,1,zeros(1,3);
    zeros(1,2),1,6,27,56,39,18,4,1,zeros(1,2);
    zeros(1,3),2,10,16,26,7,2,1,zeros(1,2);
    zeros(1,4),1,7,10,11,6,2,1,0;
    zeros(1,5),1,2,4,7,3,1,0; zeros(1,6),1,2,3,1,1,0;
    zeros(1,8),1,2,2,1; zeros(1,9),1,2,1];
x=[126:137]; y=[24:35]; fx=sum(f'); fy=sum(f);
xa=sum(fx.*x)/(sum(fx)); ya=sum(fy.*y)/(sum(fy));
vx=(sum(fx.*x.^2)-(sum(fx.*x))^2/sum(fx))/sum(fx);
vy=(sum(fy.*y.^2)-(sum(fy.*y))^2/sum(fy))/sum(fy);
cov=(x*f*y'-sum(fx.*x)*sum(fy.*y)/sum(fx))/sum(fx);
r=cov/(sqrt(vx*vy));
yb=f*y'./fx'; xb=f'*x'./fy';
for i=1:12
    syb(i)=f(i,:)*(y'-yb(i)).^2/fx(i);
end
for j=1:12
    sxb(j)=f(:,j)*(x'-xb(j)).^2/fy(j);
end
syx=sum(fx.*syb)/sum(fx);
```

```

sxy=sum(fy.*sxb)/sum(fx);
syxb=sum(fx'.*(yb-ya).^2)/sum(fx);
sxyb=sum(fy'.*(xb-xa).^2)/sum(fx);
eyx=sqrt(syxb/vy); exy=sqrt(sxyb/vx);
ayx=r*sqrt(vy/vx); axy=r*sqrt(vx/vy);
tang=(1-r^2)/r^2*sqrt(vx*vy)/(vx+vy);
arctan=atan(tang); arc=180*arctan/pi;
fprintf(' i | yb | syb | j | xb | sxb\n')
for i=1:12
    fprintf('%3d | %6.2f | %6.3f | %7d | %6.3f | %6.3f\n',...
        i,yb(i),syb(i),i,xb(i),sxb(i))
end
fprintf('\n syx=%8.5f, sxy=%8.5f',syx,sxy)
fprintf('\n eyx=%8.6f, exy=%8.6f',eyx,exy)
fprintf('\n ayx=%8.6f, axy=%8.6f',ayx,axy)
fprintf('\n tg=%9.3f, alpha=%6.0f',tang,arc)

```

În urma executării programului precedent, se obțin următoarele rezultate:

i	yb	syb	j	xb	sxb
1	25.67	0.889	1	127.000	1.000
2	26.11	1.432	2	127.556	0.691
3	26.62	2.141	3	127.813	1.902
4	28.14	1.843	4	129.433	2.112
5	28.43	1.045	5	130.464	0.844
6	29.32	1.363	6	130.943	1.204
7	29.56	1.340	7	131.438	1.280
8	30.63	1.706	8	132.000	1.778
9	31.67	1.444	9	133.125	2.276
10	31.88	1.359	10	134.091	2.992
11	33.50	0.917	11	135.429	1.959
12	34.00	0.500	12	136.500	0.250

```

syx= 1.39279, sxy= 1.40291
eyx=0.764509, exy=0.770568
ayx=0.742203, axy=0.764713
tg= 0.381, alpha= 21

```

3.4.1 Funcțiile cov și corrcoef

Aceste două funcții sunt specifice sistemului Matlab de bază și pot fi apelate prin:

```

C=cov(x)
C=cov(x,y)
r=corrcoef(x)

```

Funcția cov calculează matricea covarianțelor pentru matricea x, unde fiecare coloană este considerată a fi o variabilă:

$$C_{ij} = \frac{1}{m-1} \sum_{k=1}^m (x_{ki} - \bar{x}_i)(x_{kj} - \bar{x}_j), \quad i, j = \overline{1, m}.$$

Dacă \mathbf{x} este vector este returnată varianța vectorului, caz în care se poate folosi și parametrul opțional \mathbf{y} , care este vector de aceeași dimensiune cu \mathbf{x} , când este calculată covarianța dintre \mathbf{x} și \mathbf{y} , pe lângă varianțele lui \mathbf{x} respectiv \mathbf{y} , adică

$$C = \frac{1}{m-1} \sum_{k=1}^m (x_k - \bar{x})(y_k - \bar{y}).$$

Prin urmare, în acest caz, $\text{cov}(\mathbf{x}, \mathbf{y})$ produce același rezultat ca și $\text{cov}([\mathbf{x} \ \mathbf{y}])$.

Remarcăm faptul că instrucțiunea $\text{diag}(\text{cov}(\mathbf{x}))$ generează vectorul varianțelor fiecărei coloane a matricei \mathbf{x} , iar $\text{sqrt}(\text{diag}(\text{cov}(\mathbf{x})))$ este vectorul abaterilor standard pentru fiecare coloană a matricei \mathbf{x} .

Funcția corrcoef calculează matricea coeficienților de corelație pentru matricea \mathbf{x} , unde fiecare coloană este considerată o variabilă:

$$r_{ij} = \frac{C_{ij}}{\sqrt{C_{ii}C_{jj}}}, \quad i, j = \overline{1, m}.$$

Programul 3.4.21. Programul următor generează n vectori aleatori, care urmează legea normală bidimensională, în matricea \mathbf{X} de tipul $n \times 2$, după care, folosind funcțiile cov și corrcoef , calculează matricele covarianțelor și a coeficienților de corelație, pentru cele două coloane ale matricei \mathbf{X} .

```
clear, mu(1)=input('m1='); mu(2)=input('m2=');
v(1,1)=input('sigma1^2=');
v(2,2)=input('sigma2^2=');
v(1,2)=input('Cov(X,Y)='); v(2,1) =v(1,2);
if det(v) <= 0
    error('Matricea v nu e pozitiv definita!')
end
n=input('n='); X=mvnrnd(mu,v,n);
v=cov(X); r=corrcoef(X);
fprintf('Matricea covariantelor\n'), disp(v)
fprintf('Matricea coeficientilor de corelatie\n'), disp(r)
```

În urma executării programului, cu valorile $\mu_1=100$, $\mu_2=-100$, $\mathbf{v} = \begin{pmatrix} 5 & 4 \\ 4 & 4 \end{pmatrix}$ și $n=1000$, se obțin rezultatele

```
Matricea covariantelor
    5.1327    4.0419
    4.0419    4.0487

Matricea coeficientilor de corelatie
    1.0000    0.8867
    0.8867    1.0000
```

3.4.2 Funcțiile `lsline`, `refline` și `gline`

Funcțiile `lsline`, `refline` și `gline` aparțin la *Statistics toolbox*, apelurile lor făcându-se prin:

```
lsline
refline
refline(m)
refline(m,n)
gline
gline(f)
```

Funcția `lsline` reprezintă grafic în figura curentă dreptele de regresie, obținute cu metoda celor mai mici pătrate, pentru fiecare curbă a figurii, care a fost reprezentată grafic cu instrucțiunea `plot`, dar nu prin linie continuă, întreruptă sau de tipul punct–linie.

Funcția `refline` reprezintă grafic în figura curentă dreapta de ecuație $y=mx+n$. Dacă lipsește parametrul n , se impune ca parametrul m să fie un vector având două componente, $m(1)$ reprezentând panta dreptei, iar $m(2)=n$. Adică în acest caz dreapta, ce se va reprezenta grafic, va avea ecuația $y=m(1)x+m(2)$. Dacă lipsesc ambele argumente, m și n , efectul funcției `refline` coincide cu cel al funcției `lsline`.

Funcția `gline` reprezintă grafic în figura precizată prin parametrul f , iar dacă acesta lipsește în figura curentă, segmentul de dreaptă prin marcarea cu ajutorul *mouse*-lui a capetelor segmentului.

Programul 3.4.22. Vom prezenta un program Matlab, care generează N vectori aleatori ce urmează legea normală bidimensională, reprezintă grafic norul statistic pentru aceste date, precum și cele două drepte de regresie, a lui Y în raport cu X , respectiv a lui X în raport cu Y .

Dreapta de regresie a lui Y în raport cu X se va reprezenta grafic cu funcția `lsline`, iar dreapta de regresie a lui X în raport cu Y se va face cu funcția `refline`, având în vedere parametrii m și n ai dreptei de regresie definită prin (3.4.1), anume

$$m = \frac{1}{\bar{r}} \frac{\bar{\sigma}_Y}{\bar{\sigma}_X}, \quad n = \bar{y} - m\bar{x}.$$

```
clear,clf, mu(1)=input('m1='); mu(2)=input('m2=');
v(1,1)=input('sigma1^2=');
v(2,2)=input('sigma2^2=');
v(1,2)=input('Cov(X,Y)='); v(2,1) =v(1,2);
if det(v) <= 0
    error('Matricea v nu e pozitiv definita!')
end
N=input('N='); Z=mvnrnd(mu,v,N);
X=Z(:,1); Y=Z(:,2);
```

```
ma=mean(Z); s=std(Z); r=corrcoef([X,Y]);
m=s(2)/(r(1,2)*s(1)); n=ma(2)-m*ma(1);
plot(X,Y,'.'), lsline, refline(m,n)
text(ma(1),ma(2),'o Centrul de greutate')
```

Pentru datele de intrare $N=20$, $\mu=(0,0)$, $\mathbf{v} = \begin{pmatrix} 1 & -0.8 \\ -0.8 & 1 \end{pmatrix}$, se obțin graficele din Figura 3.12.

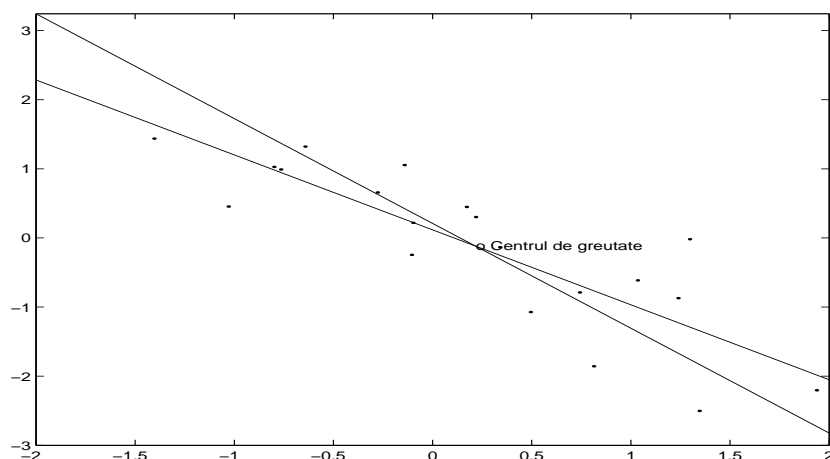


Figura 3.12: Drepte de regresie

3.4.3 Funcțiile `polyfit` și `refcurve`

Funcția `polyfit` este specifică sistemului de bază Matlab, iar funcția `refcurve` este conținută de *Statistics toolbox*, fiind folosite în ajustarea datelor prin funcții polinomiale.

Apelurile acestor funcții se pot face prin:

```
p=polyfit(x,y,n)
refcurve(p)
```

Funcția `polyfit` determină coeficienții polinomului de ajustare de grad n , cu metoda celor mai mici pătrate, corespunzător datelor statistice $(x_i, y_i)_{i=1,m}$. Prima componentă a vectorului p este coeficientul monomului de grad n . Dacă $n=1$, se obțin coeficienții dreptei de regresie a variabilei Y în raport cu variabila X .

Funcția `refcurve` reprezintă grafic în figura curentă funcția polinomială precizată prin coeficienții p , cu precizarea că $p(1)$ reprezintă coeficientul monomului de grad maxim. Dacă parametrul p are două componente, funcția este echivalentă cu funcția `refline`.

Calculul valorilor unei funcții polinomiale se face cu ajutorul funcției `polyval`, pe care o prezentăm în continuare, împreună cu alte funcții specifice sistemului Matlab de bază, din acest grup de funcții.

3.4.4 Funcțiile `polyval`, `polyvalm`, `poly` și `roots`

Forme de apel pentru funcțiile `polyval`, `polyvalm`, `poly` și `roots` sunt:

```
y=polyval(p,x)
y=polyvalm(p,x)
p=poly(a)
p=roots(a)
```

Funcția `polyval` calculează valorile y ale polinomului precizat prin coeficienții conținuți de vectorul p , pe fiecare punct al matricei sau vectorului x .

Funcția `polyvalm` calculează matricea pătratică y , care reprezintă valoarea matriceală a polinomului precizat prin coeficienții conținuți de vectorul p , pentru matricea pătratică x , adică

$$y = p_1 x^n + p_2 x^{n-1} + \dots + p_{n+1},$$

dacă p are $n + 1$ componente.

Funcția `poly` calculează coeficienții polinomului caracteristic al matricei pătratică a :

$$P(\lambda) = \begin{vmatrix} \lambda - a_{11} & a_{12} & \dots & a_{1m} \\ a_{21} & \lambda - a_{22} & \dots & a_{2m} \\ \dots & \dots & \dots & \dots \\ a_{m1} & a_{m2} & \dots & \lambda - a_{mm} \end{vmatrix},$$

iar când a este un vector calculează coeficienții polinomului care are ca rădăcini elementele vectorului, adică ai lui

$$P(x) = \prod_{i=1}^n (x - a_i).$$

Funcția `roots` calculează rădăcinile polinomului cu coeficienții dați prin vectorul p .

Programul 3.4.23. Programul următor va genera matricea Z , care conține N vectori aleatori ce urmează legea normală bidimensională. Se transformă apoi componenta a doua a vectorilor aleatori, prin ridicarea acestora la puterea a treia, iar apoi se determină polinomul de ajustare de grad trei, prin folosirea datelor transformate. Polinomul de ajustare obținut se reprezintă grafic împreună cu norul statistic al datelor transformate.

```

clf, clear, mu(1)=input('m1='); mu(2)=input('m2=');
v(1,1)=input('sigma1^2=');
v(2,2)=input('sigma2^2=');
v(1,2)=input('Cov(X,Y)='); v(2,1) =v(1,2);
if det(v) <= 0
    error('Matricea v nu e pozitiv definita!')
end
N=input('N='); Z=mvnrnd(mu,v,N);
X=Z(:,1); Y=Z(:,2); p=polyfit(X,Y.^3,3);
scatter(X,Y.^3,7,'filled'), hold on,
m=mu(1); s=sqrt(v(1,1));
x=m-3*s:0.01:m+3*s;
y=polyval(p,x); plot(x,y)

```

Pentru datele de intrare $N=25$, $\mu=(0,0)$, $v = \begin{pmatrix} 1 & -0.9 \\ -0.9 & 1 \end{pmatrix}$, se obține graficul din Figura 3.13.

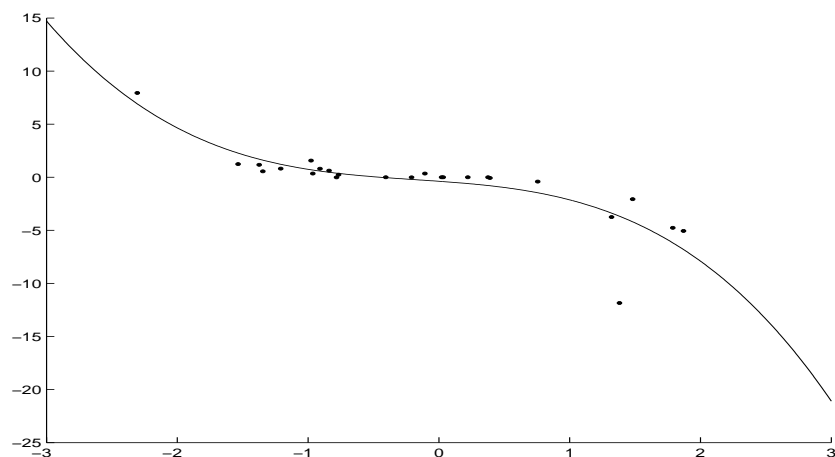


Figura 3.13: Polinom de ajustare de grad 3

3.4.5 Funcția polytool

Funcția demonstrativă `polytool` se lansează prin:

```

>>polytool(x,y)
>>polytool(x,y,n)

```

În urma lansării se produce un grafic interactiv (demonstrativ) privind polinomul de ajustare pentru datele conținute în vectorii coloană x și y , având gradul precizat

prin parametrul n (valoarea implicită este $n=1$). Schimbarea gradului polinomului de ajustare se face prin introducerea noii valori în fereastra din partea de sus.

Pe figură este reprezentat norul statistic pentru datele x și y , împreună cu graficul polinomului de ajustare (cu linie continuă).

Prin introducerea valorii unui argument x în fereastra de jos sau prin deplasarea drepte verticale pe poziția abscisei x , pe axa ordonatelor se obține valoarea corespunzătoare.

Pe lângă aceste elemente, graficul demonstrativ conține și alte elemente, care se referă la noțiuni ce urmează a fi prezentate în capitolele ce urmează, cum ar fi *probabilitate de încredere* și *intervale de încredere*.

De asemenea, există posibilitatea ajustării datelor cu alte metode decât metoda celor mai mici pătrate. Alegerea unei astfel de metodă se face prin activarea butonului `method` din meniu.

3.4.6 Funcția `nlinfit`

Pentru ajustarea neliniară, *Statistics toolbox* dispune de funcția `nlinfit`, bazată pe metoda celor mai mici pătrate. Apelul funcției se poate face prin:

```
a=nlinfit(x,y,numef,a0)
[a,r]=nlinfit(x,y,f,a0)
```

unde `numef` reprezintă numele funcției, care definește legătura căutată dintre variabila dependentă, cu valorile conținute de vectorul coloană y , variabilele independente, cu valorile conținute în coloanele matricei x . Matricea x trebuie să aibă același număr de linii ca și lungimea vectorului y .

Parametrul `a0` reprezintă valorile de pornire (inițiale) ale parametrilor a , care urmează să fie calculați.

Funcția `numef` poate fi o funcție Matlab proprie, care are linia de definiție

```
function z=numef(a,x)
```

În urma executării primei instrucțiuni, se obțin coeficienții a ai funcției `numef`, iar dacă se folosește a doua instrucțiune se mai obțin și erorile $r=y-\text{numef}(a,x)$.

Programul 3.4.24. Vom considera N vectori aleatori ce urmează legea normală bi-dimensională, păstrați în matricea Z de tipul $(N, 2)$. Asupra elementelor coloanei a doua se efectuează transformarea dată prin funcția `exp`. Dacă a doua coloană a lui Z astfel transformată o considerăm că reprezintă valorile unei variabile dependente Y , iar prima coloană considerăm că reprezintă valorile unei variabile independente X , vrem să determinăm legătura de forma

$$Y = be^{aX}.$$

Procedăm în două moduri, cu ajutorul funcției `nlinfit` (linie continuă pe grafic), respectiv prin liniarizarea acestei relații (linie întreruptă pe grafic), adică prin logaritmare: $\ln Y = \ln b + aX$.

```

clf, mu(1)=input('m1='); mu(2)=input('m2=');
v(1,1)=input('sigma1^2=');
v(2,2)=input('sigma2^2=');
v(1,2)=input('Cov(X,Y)='); v(2,1) =v(1,2);
if det(v) <= 0
    error('Matricea v nu e pozitiv definita!')
end
N=input('N='); Z=mvnrnd(mu,v,N);
X=Z(:,1); Y=exp(Z(:,2));
a=nlinfit(X,Y,'expo',[1,1]);
p=polyfit(X,log(Y),1);
a1=p(1); b1=exp(p(2));
scatter(X,Y,7,'filled'), hold on
m=mu(1); s=sqrt(v(1,1));
x=m-3*s:0.01:m+3*s;
y=expo(a,x); y1=b1*exp(a1*x);
plot(x,y,'-',x,y1,'--')

```

Considerând datele de intrare $N=20$, $\mu=(0,0)$, $v = \begin{pmatrix} 1 & -0.9 \\ -0.9 & 1 \end{pmatrix}$, se obțin graficele din Figura 3.14.

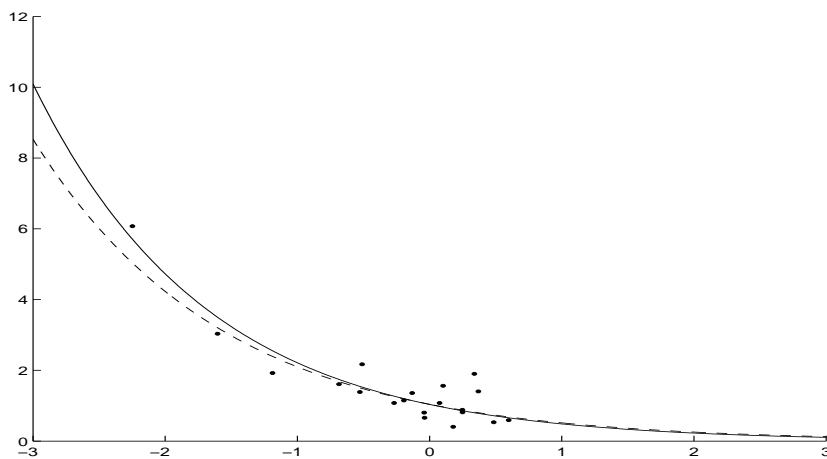


Figura 3.14: Ajustare neliniară

3.4.7 Funcția nlintool

Funcția demonstrativă polytool se lansează prin:

```
>>nlintool(x,y,numef,a0)
```

În urma lansării se produce un grafic interactiv (demonstrativ) privind ajustarea datelor conținute în vectorul coloană y și matricea x , având tipul legăturii precizat prin parametrul `numef`, care poate fi o funcție Matlab proprie.

Pe figură este reprezentat graficul de ajustare (cu linie continuă).

Prin introducerea valorii unui argument x în fereastra de jos sau prin deplasarea dreptei verticale pe poziția abscisei x , pe axa ordonatelor se obține valoarea corespunzătoare pentru y .

Elementele obținute, prin lansarea acestei funcții, pot fi salvate în spațiul de lucru, `workspace`, parțial sau în totalitate, prin utilizarea butonului `export`.

Pe lângă aceste elemente, graficul demonstrativ conține și alte elemente, care se referă la noțiuni ce urmează a fi prezentate în capitolele ce urmează, cum ar fi *intervale de încredere*.

3.4.8 Coeficienții Spearman și Kendall

Definiția 3.4.25. Fie (u_k, v_k) , $k = \overline{1, N}$, rangurile datelor statistice primare (x'_k, y'_k) , $k = \overline{1, N}$, obținute prin ordonarea crescătoare după prima, respectiv a doua componentă. Numim coeficient de corelație al rangurilor sau coeficientul lui Spearman, valoarea numerică

$$\bar{s} = \bar{r}(U, V),$$

unde U și V sunt noile caracteristici care definesc rangurile datelor statistice, respectiv pentru X și Y .

Proprietatea 3.4.26. Dacă notăm $d_k = |u_k - v_k|$ diferența rangurilor aceluiași individ, atunci

$$\bar{s} = 1 - \frac{6}{N(N^2 - 1)} \sum_{k=1}^N d_k^2.$$

Demonstrație. Conform definiției lui \bar{s} avem că

$$\bar{s} = \frac{\sum_{i=1}^N (u_i - \bar{u})(v_i - \bar{v})}{\sqrt{\sum_{i=1}^N (u_i - \bar{u})^2} \sqrt{\sum_{i=1}^N (v_i - \bar{v})^2}}.$$

Deoarece $\{u_1, u_2, \dots, u_N\} = \{v_1, v_2, \dots, v_N\} = \{1, 2, \dots, N\}$, rezultă că

$$\bar{u} = \bar{v} = \frac{1}{N} \sum_{k=1}^N k = \frac{N+1}{2}$$

și

$$\sum_{k=1}^N (u_k - \bar{u})^2 = \sum_{k=1}^N (v_k - \bar{v})^2 = \sum_{k=1}^N \left(k - \frac{N+1}{2}\right)^2$$

$$= \frac{N(N+1)(2N+1)}{6} - \frac{N(N+1)^2}{2} + \frac{N(N+1)^2}{4} = \frac{N(N^2-1)}{12},$$

deci numitorul din expresia lui \bar{s} este calculat.

Pentru calcularea numărătorului pornim de la identitatea

$$\begin{aligned} 2(u_k - \bar{u})(v_k - \bar{v}) &= (u_k - \bar{u})^2 + (v_k - \bar{v})^2 - [(u_k - \bar{u}) - (v_k - \bar{v})]^2 \\ &= (u_k - \bar{u})^2 + (v_k - \bar{v})^2 - d_k^2, \end{aligned}$$

care se însumează pentru $k = \overline{1, N}$. Obținem în acest fel

$$\begin{aligned} 2 \sum_{k=1}^N (u_k - \bar{u})(v_k - \bar{v}) &= \sum_{k=1}^N (u_k - \bar{u})^2 + \sum_{k=1}^N (v_k - \bar{v})^2 - \sum_{k=1}^N d_k^2 \\ &= \frac{N(N^2-1)}{12} \cdot 2 - \sum_{k=1}^N d_k^2, \end{aligned}$$

adică

$$\sum_{k=1}^N (u_k - \bar{u})(v_k - \bar{v}) = \frac{N(N^2-1)}{12} - \frac{1}{2} \sum_{k=1}^N d_k^2.$$

Dacă se înlocuiesc aceste expresii în formula lui \bar{s} se obține relația din enunțul proprietății. \square

Corolarul 3.4.27. Coeficientul lui Spearman verifică relațiile $-1 \leq \bar{s} \leq 1$.

Demonstrație. Valoarea maximă $\bar{s} = 1$ se obține când $d_k = 0$, $k = \overline{1, N}$, adică toate rangurile corespunzătoare celor două caracteristici coincid. Valoarea minimă $\bar{s} = -1$ se va obține când diferențele rangurilor sunt maxime, adică rangurile sunt inverse pentru cele două caracteristici. În acest caz avem

$$\sum_{k=1}^N d_k^2 = 2 \left[(N-1)^2 + (N-3)^2 + \dots \right] = 2 \cdot \frac{N(N^2-1)}{6},$$

de unde se obține, într-adevăr, $\bar{s} = -1$.

Același rezultat se obține și din faptul că \bar{s} este definit cu ajutorul coeficientului (pentru ranguri) Pearson, despre care știe că ia valori în intervalul $[-1, 1]$. \square

Observația 3.4.28. Pe baza corolarului precedent, vom spune că pentru $\bar{s} = 1$ cele două clasamente pentru caracteristicile X și Y sunt identice, pentru $\bar{s} = -1$ sunt inverse unul celuilalt, iar pentru $\bar{s} = 0$ sunt independente.

Observația 3.4.29. Când există două sau mai multe date statistice primare care au aceeași valoare, atunci rangurile acestora se consideră toate egale (*ex-aequo*) cu media aritmetică a rangurilor pe care le ocupă aceste date în ordonarea crescătoare. corespunzător vom face modificarea lui \bar{s} .

Dacă în ordonarea după caracteristica X avem u grupe de date statistice care coincid, fiecare grupă conținând g_i , $i = \overline{1, u}$, date statistice și dacă în ordonarea după caracteristica Y avem v grupe de date statistice ce coincid, fiecare grupă conținând h_j , $j = \overline{1, v}$, date statistice, atunci \bar{s} se modifică după cum urmează

$$\bar{s} = \frac{\frac{1}{6} [N(N^2 - 1) - (t_x + t_y)] - \sum_{k=1}^N d_k^2}{\sqrt{\frac{1}{6} N(N^2 - 1) - 2t_x} \sqrt{\frac{1}{6} N(N^2 - 1) - 2t_y}},$$

unde

$$t_x = \frac{1}{12} \sum_{i=1}^u g_i (g_i^2 - 1) \quad \text{și} \quad t_y = \frac{1}{12} \sum_{j=1}^v h_j (h_j^2 - 1).$$

Definiția 3.4.30. Numim coeficientul lui Kendall relativ la distribuția statistică a caracteristicii bidimensionale (X, Y) , raportul

$$\bar{k} = \frac{2\bar{t}}{N(N-1)}, \quad \text{unde} \quad \bar{t} = \sum_{\substack{i,j=1 \\ i < j}}^N \text{sign} \{ (x'_j - x'_i) (y'_j - y'_i) \}.$$

Proprietatea 3.4.31. Coeficientul lui Kendall satisface relațiile $-1 \leq \bar{k} \leq 1$.

Demonstrație. Valoarea maximă a lui \bar{k} se obține când \bar{t} este maxim, adică atunci când $(x'_j - x'_i) (y'_j - y'_i) > 0$, pentru fiecare $i < j$. În acest caz avem că

$$\bar{t} = \sum_{\substack{i,j=1 \\ i < j}}^N 1 = \binom{N}{2} = \frac{N(N-1)}{2},$$

deci $\bar{k} = 1$.

Valoarea minimă a lui \bar{k} se obține când \bar{t} este minim, $(x'_j - x'_i) (y'_j - y'_i) < 0$, pentru fiecare $i < j$, caz în care

$$\bar{t} = - \sum_{\substack{i,j=1 \\ i < j}}^N 1 = - \binom{N}{2} = - \frac{N(N-1)}{2},$$

deci $\bar{k} = -1$. Prin urmare avem $-1 \leq \bar{k} \leq 1$. □

Observația 3.4.32. Din proprietatea precedentă, avem că pentru $\bar{k} = 1$ cele două clasamente ale celor două caracteristici X și Y sunt identice, pentru $\bar{k} = -1$ sunt inverse unul celuilalt, iar pentru $\bar{k} = 0$ sunt independente.

Observația 3.4.33. Pentru calculul rapid al lui \bar{k} , se poate proceda după cum urmează. Se ordonează datele primare (x'_k, y'_k) , $k = \overline{1, N}$, în mod crescător după prima componentă. Fie aceasta $(x'_{(k)}, y'_{(k)})$, $k = \overline{1, N}$, cu $x'_{(1)} \leq x'_{(2)} \leq \dots \leq x'_{(N)}$. Se calculează apoi numărul

$$\bar{t} = \sum_{\substack{u, v=1 \\ u < v}} \text{sign} \left(y'_{(v)} - y'_{(u)} \right),$$

obținându-se astfel \bar{k} .

Remarcăm de asemenea că toate formulele relative la coeficientul lui Kendall rămân adevărate dacă se lucrează cu rangurile (u_i, v_i) , corespunzătoare datelor statistice primare (x'_i, y'_i) , $i = \overline{1, N}$, așa cum s-au introdus pentru definirea coeficientului lui Spearman.

Observația 3.4.34. Când în cele două clasamente sunt valori egale (*ex-aequo*) se procedează ca la Observația 3.4.29, adică se înlocuiesc toate rangurile pentru valorile egale prin media aritmetică a rangurilor pe care le ocupă în ordonare. Corespunzător se face modificarea lui \bar{k} prin formula

$$\bar{k} = \frac{\bar{t}}{\sqrt{\frac{N(N-1)}{2} - u_x} \sqrt{\frac{N(N-1)}{2} - u_y}},$$

unde $u_x = \frac{1}{2} \sum_{i=1}^u g_i (g_i - 1)$ și $u_y = \frac{1}{2} \sum_{j=1}^v h_j (h_j - 1)$, cu g_i , h_j definiți la Observația 3.4.29. Pentru calculul lui \bar{t} pentru fiecare pereche *ex-aequo* valoarea lui \bar{t} rămâne neschimbată.

Exemplul 3.4.35. Echipele diviziei naționale de fotbal sunt clasificate după numărul de goluri primite (X) și după numărul cartonașelor galbene primite (Y) în primele 10 meciuri ale campionatului. Datele statistice sunt trecute în tabelul următor

Echipa	E ₁	E ₂	E ₃	E ₄	E ₅	E ₆	E ₇	E ₈	E ₉	E ₁₀	E ₁₁	E ₁₂	E ₁₃	E ₁₄	E ₁₅	E ₁₆	E ₁₇	E ₁₈
X	15	13	10	13	12	11	16	11	13	9	10	14	15	9	8	9	10	11
Y	13	9	11	11	12	8	7	6	9	8	8	10	11	9	5	10	13	14

Să calculăm coeficienții lui Spearman respectiv al lui Kendall. Pentru aceasta ordonăm crescător echipele după caracteristica X , după care folosind regula *ex-aequo* stabilim rangurile datelor statistice.

Calculul este efectuat în tabelul următor.

Echipa	X	Y	Rang(X)	Rang(Y)	d	d^2	-	+
E ₁₅	8	5	1	1	0	0	0	17
E ₁₀	9	8	3	5	2	4	2	12
E ₁₄	9	9	3	8	5	25	4	9
E ₁₆	9	10	3	10.5	7.5	56.25	6	7
E ₃	10	11	6	13	7	49	7	4
E ₁₁	10	8	6	5	1	1	2	9
E ₁₇	10	13	6	16.5	10.5	110.25	9	1
E ₆	11	8	9	5	4	16	2	8
E ₈	11	6	9	2	7	49	0	9
E ₁₈	11	14	9	18	9	81	8	0
E ₅	12	12	11	15	4	16	6	1
E ₂	13	9	13	8	5	25	1	4
E ₄	13	11	13	13	0	0	3	1
E ₉	13	9	13	8	5	25	1	3
E ₁₂	14	10	15	10.5	4.5	20.25	1	2
E ₁	15	13	16.5	16.5	0	0	2	0
E ₁₃	15	11	16.5	13	3.5	12.25	1	0
E ₇	16	7	18	3	15	225	0	0
Total						715	55	87

Numărul grupelor ex-aequo este cinci atât pentru X cât și pentru Y . Avem astfel

$$t_x = \frac{1}{12} [3(9-1) + 3(9-1) + 3(9-1) + 3(9-1) + 2(4-1)] = \frac{17}{2},$$

$$t_y = \frac{1}{12} [3(9-1) + 3(9-1) + 2(4-1) + 3(9-1) + 2(4-1)] = 7,$$

$$u_x = \frac{1}{2} [3(3-1) + 3(3-1) + 3(3-1) + 3(3-1) + 2(2-1)] = 13,$$

$$u_y = \frac{1}{2} [3(3-1) + 3(3-1) + 2(2-1) + 3(3-1) + 2(2-1)] = 11,$$

$$\bar{t} = 87 - 55 = 32.$$

Dacă se înlocuiesc aceste valori în formule se obțin coeficienții

$$\bar{s} = \frac{\frac{1}{6} [18 \cdot 19 \cdot 17 - (\frac{17}{2} + 7)] - 715}{\sqrt{\frac{1}{6} \cdot 18 \cdot 17 \cdot 19 - 2 \cdot \frac{17}{2}} \sqrt{\frac{1}{6} \cdot 18 \cdot 17 \cdot 19 - 2 \cdot 7}} = 0.263,$$

$$\bar{k} = \frac{3}{2} \sqrt{\frac{18 \cdot 17}{2}} - 13 \sqrt{\frac{18 \cdot 17}{2}} - 11 = 0.227.$$

Programul 3.4.36. Programul Matlab, care urmează, calculează cei doi coeficienți, Spearman și Kendall, fără tratarea cazurilor ex-aequo.

```
x=[15,13,10,13,12,11,16,11,13,9,10,14,15,9,8,9,10,11];
y=[13,9,11,11,12,8,7,6,9,8,8,10,11,9,5,10,13,14];
[xord,ind]=sort(x); [yord,jnd]=sort(y);
for i=1:18
    rx(ind(i))=i; ry(jnd(i))=i;
end
d=abs(rx-ry); d2=d.^2;
sp=1-6*sum(d2)/18/(18^2-1);
ke=0;
for i=1:18
    for j=1:18
        if i<j
            ke=ke+sign((rx(i)-rx(j))*(ry(i)-ry(j)));
        end
    end
end
ke=2*ke/(18*17);
fprintf(' x | y |rang(x)|rang(y)| d | d^2 \n')
for i=1:18
    fprintf('%3d |%4d |%5d |%5d |%4d |%4d\n',...
        x(i),y(i),rx(i),ry(i),d(i),d2(i))
end
fprintf('\n sp=%5.3f, ke=%5.3f',sp,ke)
```

În urma executării acestui program, se obțin următoarele rezultate:

x	y	rang(x)	rang(y)	d	d ²
15	13	16	16	0	0
13	9	12	7	5	25
10	11	5	12	7	49
13	11	13	13	0	0
12	12	11	15	4	16
11	8	8	4	4	16
16	7	18	3	15	225
11	6	9	2	7	49
13	9	14	8	6	36
9	8	2	5	3	9
10	8	6	6	0	0
14	10	15	10	5	25
15	11	17	14	3	9
9	9	3	9	6	36
8	5	1	1	0	0
9	10	4	11	7	49
10	13	7	17	10	100
11	14	10	18	8	64

sp=0.269, ke=0.203

Observația 3.4.37. (*Formula lui Daniels*). Coeficientul \bar{r} de corelație (al lui Pearson), coeficientul \bar{s} al lui Spearman și coeficientul \bar{k} al lui Kendall se pot exprima prin formula unică

$$\bar{d} = \frac{\sum_{i=1}^N \sum_{j=1}^N a_{ij} b_{ij}}{\sqrt{\sum_{i=1}^N \sum_{j=1}^N a_{ij}^2} \sqrt{\sum_{i=1}^N \sum_{j=1}^N b_{ij}^2}}.$$

Dacă se ia

$$a_{ij} = x'_i - x'_j, \quad b_{ij} = y'_i - y'_j$$

se obține coeficientul de corelație \bar{r} .

Într-adevăr, prin calcule algebrice simple se poate scrie

$$\sum_{i=1}^N \sum_{j=1}^N (x'_i - x'_j) (y'_i - y'_j) = 2N \sum_{i=1}^N x'_i y'_i - 2 \left(\sum_{i=1}^N x'_i \right) \left(\sum_{i=1}^N y'_i \right),$$

iar

$$\begin{aligned} \sum_{i=1}^N \sum_{j=1}^N (x'_i - x'_j)^2 &= 2N \sum_{i=1}^N x_i'^2 - 2 \left(\sum_{i=1}^N x'_i \right)^2, \\ \sum_{i=1}^N \sum_{j=1}^N (y'_i - y'_j)^2 &= 2N \sum_{i=1}^N y_i'^2 - 2 \left(\sum_{i=1}^N y'_i \right)^2, \end{aligned}$$

astfel că

$$\bar{d} = \frac{\sum_{i=1}^N x'_i y'_i - \frac{1}{N} \left(\sum_{i=1}^N x'_i \right) \left(\sum_{i=1}^N y'_i \right)}{\sqrt{\sum_{i=1}^N x_i'^2 - \frac{1}{N} \left(\sum_{i=1}^N x'_i \right)^2} \sqrt{\sum_{i=1}^N y_i'^2 - \frac{1}{N} \left(\sum_{i=1}^N y'_i \right)^2}} = \bar{r}.$$

Dacă se ia

$$a_{ij} = u_i - u_j, \quad b_{ij} = v_i - v_j$$

se obține, ca mai înainte, coeficientul \bar{s} al lui Spearman, iar dacă se consideră

$$a_{ij} = \text{sign} (x'_i - x'_j), \quad b_{ij} = \text{sign} (y'_i - y'_j)$$

se obține coeficientul \bar{k} al lui Kendall.

Definiția 3.4.38. Fie colectivitatea C cercetată din punct de vedere a p caracteristici X_1, X_2, \dots, X_p și fie datele statistice primare relative la aceste caracteristici scrise în matricea

$$\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \dots & \dots & \dots & \dots \\ x_{N1} & x_{N2} & \dots & x_{Np} \end{pmatrix},$$

cu matricea \mathbf{R} a rangurilor corespunzătoare pentru fiecare caracteristică dată prin

$$\mathbf{R} = \begin{pmatrix} r_{11} & r_{12} & \dots & r_{1p} \\ r_{21} & r_{22} & \dots & r_{2p} \\ \dots & \dots & \dots & \dots \\ r_{N1} & r_{N2} & \dots & r_{Np} \end{pmatrix}.$$

Numim coeficientul de concordanță al lui Kendall raportul

$$\overline{w} = \frac{12\overline{s}}{p^2 (N^3 - N)},$$

unde

$$\overline{s} = \sum_{i=1}^N \left(r_{i.} - \frac{r_{..}}{N} \right)^2, \quad r_{i.} = \sum_{j=1}^p r_{ij}, \quad r_{..} = \sum_{i=1}^N r_{i.}.$$

Observația 3.4.39. Fiecare coloană a tabloului rangurilor este o permutare a numerelor $1, 2, \dots, N$, deci suma elementelor din fiecare coloană este $\frac{N(N+1)}{2}$ și prin urmare

$$r_{..} = p \frac{N(N+1)}{2}.$$

Când cele p clasamente sunt identice, deci coloanele din matricea rangurilor sunt identice, totalurile $r_{1.}, r_{2.}, \dots, r_{N.}$ formează o permutare a elementelor $p, 2p, \dots, Np$. În acest caz valoarea lui \overline{s} va fi maximă, anume

$$\overline{s} = \sum_{i=1}^N \left(pi - \frac{p(N+1)}{2} \right)^2 = \frac{p^2 (N^3 - N)}{12}.$$

prin urmare $\overline{w} \leq 1$.

Când $r_{1.} = r_{2.} = \dots = r_{N.}$ atunci se obține valoarea minimă $\overline{w} = 0$, caz în care avem independență între caracteristicile cercetate.

Observația 3.4.40. Cazurile ex-aequo se tratează ca și pentru două caracteristici, iar formula modificată este

$$\overline{w} = \frac{12\overline{s}}{p^2 (N^3 - N) - p \sum_{j=1}^p (t_j^3 - t_j)}$$

unde t_j este numărul datelor ex-aequo pentru caracteristica X_j .

Exemplul 3.4.41. Se consideră 12 familii pentru care sunt cercetate următoarele caracteristici X_1, X_2, X_3, X_4 reprezentând cheltuielile (în sute de lei) zilnice pentru pâine, legume, fructe, carne respectiv. S-au obținut următoarele date statistice

X_1	X_2	X_3	X_4
332	428	354	1437
293	559	388	1527
372	767	562	1948
406	563	341	1509
386	608	396	1501
438	843	689	2345
534	660	367	1620
460	699	483	1856
385	789	621	2366
655	776	423	1848
584	995	548	2056
515	1097	887	2630

Vrem să calculăm coeficientul de concordanță al lui Kendall.

Pentru a calcula coeficientul de concordanță al lui Kendall, se determină rangurile datelor statistice pentru fiecare caracteristică în parte. Matricea rangurilor este

$$\mathbf{R} = \left(r_{ij} \right)_{\substack{i=\overline{1,12} \\ j=\overline{1,4}}} = \begin{pmatrix} 2 & 1 & 2 & 1 \\ 1 & 2 & 4 & 4 \\ 3 & 7 & 9 & 8 \\ 6 & 3 & 1 & 3 \\ 5 & 4 & 5 & 2 \\ 7 & 10 & 11 & 10 \\ 10 & 5 & 3 & 5 \\ 8 & 6 & 7 & 7 \\ 4 & 9 & 10 & 11 \\ 12 & 8 & 6 & 6 \\ 11 & 11 & 8 & 9 \\ 9 & 12 & 12 & 12 \end{pmatrix}.$$

De asemenea, avem că sumele rangurilor de pe fiecare linie a matricei rangurilor sunt 6, 11, 27, 13, 16, 38, 23, 28, 34, 32, 39, 45, deci

$$\bar{s} = \sum_{i=1}^N \left(r_{i.} - \frac{r_{..}}{N} \right)^2 = \left(6 - \frac{312}{12} \right)^2 + \left(11 - \frac{312}{12} \right)^2 + \cdots + \left(45 - \frac{312}{12} \right)^2$$

$$=(6-26)^2 + (11-26)^2 + \dots + (45-26)^2 = 1682.$$

Coeficientul de concordanță al lui Kendall va fi

$$\overline{w} = \frac{12 \overline{s}}{p^2(N^3 - N)} = \frac{12 \cdot 1682}{4^2(12^3 - 12)} = \frac{841}{1441} \cong 0.7351.$$

Programul 3.4.42. Un program Matlab, care efectuează aceste calcule, este prezentat în continuare

```
x=[332,428,354,1437;293,559,388,1527;372,767,562,1948;
  406,563,341,1509;386,608,396,1501;438,843,689,2345;
  534,660,367,1620;460,699,483,1856;385,789,621,2366;
  655,776,423,1848;584,995,548,2056;515,1097,887,2630];
N=12; p=4;
[x1,i1]=sort(x(:,1)); [x2,i2]=sort(x(:,2));
[x3,i3]=sort(x(:,3)); [x4,i4]=sort(x(:,4));
for i=1:12
    r1(i1(i))=i; r2(i2(i))=i; r3(i3(i))=i; r4(i4(i))=i;
end
r=[r1',r2',r3',r4']; sb=sum((sum(r')-sum(sum(r))/12).^2);
w=12*sb/(p^2*(N^3-N));
fprintf('    Matricea rangurilor\n')
for i=1:12
    fprintf('%6d | %3d | %3d | %3d\n',r1(i),r2(i),r3(i),r4(i))
end
fprintf('\n          w=%5.3f',w)
```

iar rezultatele obținute, în urma executării programului, sunt

Matricea rangurilor

2		1		2		1
1		2		4		4
3		7		9		8
6		3		1		3
5		4		5		2
7		10		11		10
10		5		3		5
8		6		7		7
4		9		10		11
12		8		6		6
11		11		8		9
9		12		12		12

w=0.735

Capitolul 4

Teoria selecției

Studiul statistic al unei colectivități C din punct de vedere al uneia sau mai multor caracteristici prin considerarea tuturor indivizilor colectivității de multe ori este imposibil de efectuat sau foarte costisitor. Dacă ne gândim, de exemplu, la studiul calității produselor unei fabrici de conserve, ne dăm seama că nu are sens considerarea tuturor acestor produse pentru efectuarea controlului, ceea ce ar duce la distrugerea întregii producții. De asemenea, dacă se are în vedere recensământul populației unei țări, ne dăm seama de costul ridicat al acestei operații, motiv pentru care astfel de recensăminte sunt destul de rare.

Și numai din cele două exemple amintite înainte, ne dăm seama că ar fi potrivit ca să fie considerată pentru studiu numai o parte a colectivității cercetate, iar apoi rezultatele obținute relative la această parte să fie extinse la întreaga colectivitate.

4.1 Tipuri de selecție

Definiția 4.1.1. Numim eșantion (selecție, sondaj) relativ la colectivitatea C o submulțime de indivizi \mathcal{E} a lui C , care urmează să fie cercetați din punct de vedere a uneia sau mai multor caracteristici, iar numărul indivizilor din eșantionul \mathcal{E} se numește volumul eșantionului.

Observația 4.1.2. Modurile de obținere a eșantionului \mathcal{E} ne conduc la *metode nealeatoare* și respectiv *metode aleatoare* de selecție.

Dintre metodele nealeatoare amintim:

- *selecția sistematică*, când indivizii care intră în eșantion sunt considerați după o anumită regulă, de exemplu din 10 în 10;

- *selecție tipică*, când, cunoscându-se informații anterioare referitoare la colectivitate, sunt considerați indivizi cu valori medii apropiate de valoarea medie a întregii colectivități;
- *selecție stratificată*, când colectivitatea este clasificată (stratificată) după anumite criterii, cunoscându-se proporția indivizilor pentru fiecare strat. Eșantionul se ia astfel încât să fie respectate aceste proporții pentru fiecare strat.

Pentru metodele aleatoare, fiecare individ al colectivității C poate să intre în eșantion cu aceeași probabilitate (*selecție cu probabilități egale*) sau cu probabilități diferite. Metodele aleatoare de selecție sunt:

- *repetate (bernoulliene)*, când individul, ce intră în eșantion, după ce a fost cercetat, este reintrodus în colectivitate;
- *nerepetate*, când individul ce intră în eșantion, după ce a fost cercetat, nu este reintrodus în colectivitate.

Observația 4.1.3. Dacă volumul colectivității este mult mai mare decât volumul eșantionului, atunci o selecție nerepetată poate fi considerată ca fiind de tip repetat. Aceasta are la bază rezultatul privind comportarea asimptotică a legii hipergeometrice ca și o lege binomială (a se vedea Observația 2.4.4).

4.2 Funcții de selecție

În cele ce urmează vom considera că avem de fiecare dată o selecție repetată.

Definiția 4.2.1. Fie colectivitatea C cercetată din punct de vedere al caracteristicii X . Numim date de selecție relative la caracteristica X datele statistice x_1, x_2, \dots, x_n privind indivizii care intră în eșantion.

Observația 4.2.2. Datele de selecție pot fi considerate ca fiind valorile unor variabile aleatoare X_1, X_2, \dots, X_n , numite *variabile de selecție* și care în cazul unei selecții repetate sunt variabile aleatoare independente, identic repartizate cu X .

Definiția 4.2.3. Numim funcție de selecție sau statistică variabila aleatoare

$$Z_n = h_n(X_1, X_2, \dots, X_n),$$

unde $h_n: \mathbb{R}^n \rightarrow \mathbb{R}$ este o funcție măsurabilă, iar $z_n = h_n(x_1, x_2, \dots, x_n)$ se numește valoarea funcției de selecție.

4.2.1 Media de selecție

Definiția 4.2.4. Numim medie de selecție funcția de selecție

$$\bar{X} = \frac{1}{n} \sum_{k=1}^n X_k, \quad \text{iar} \quad \bar{x} = \frac{1}{n} \sum_{k=1}^n x_k$$

se numește valoarea mediei de selecție.

Proprietatea 4.2.5. Fie caracteristica X având valoarea medie $m = E(X)$ și dispersia $\sigma^2 = Var(X)$, atunci

$$E(\bar{X}) = m \quad \text{și} \quad Var(\bar{X}) = \frac{1}{n} \sigma^2.$$

Demonstrație. Folosind proprietățile valorii medii și ale dispersiei și având în vedere că selecția este repetată avem succesiv

$$\begin{aligned} E(\bar{X}) &= \frac{1}{n} \sum_{k=1}^n E(X_k) = \frac{1}{n} \sum_{k=1}^n E(X) = \frac{1}{n} n m = m, \\ Var(\bar{X}) &= \frac{1}{n^2} \sum_{k=1}^n Var(X_k) = \frac{1}{n^2} \sum_{k=1}^n Var(X) = \frac{1}{n^2} n \sigma^2 = \frac{1}{n} \sigma^2. \end{aligned}$$

□

Observația 4.2.6. În cazul în care caracteristica X urmează legea normală $\mathcal{N}(\mu, \sigma)$, atunci \bar{X} , fiind o combinație liniară de variabile aleatoare independente, ce urmează fiecare aceeași lege normală, va urma de asemenea legea normală (a se vedea Observația 2.5.6). Pe baza proprietății precedente \bar{X} va urma așadar legea normală $\mathcal{N}\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$.

Proprietatea 4.2.7. Fie caracteristica X având valoarea medie $m = E(X)$ și dispersia $\sigma^2 = Var(X)$, atunci statistica

$$(4.2.1) \quad Z_n = \frac{\bar{X} - m}{\frac{\sigma}{\sqrt{n}}}$$

converge în repartiție la legea normală $\mathcal{N}(0, 1)$, când $n \rightarrow \infty$, iar când X urmează legea normală $\mathcal{N}(m, \sigma)$ afirmația are loc pentru orice valoare a lui n .

Proprietatea nu reprezintă altceva decât Teorema 2.7.8.

Funcția 4.2.8. Pentru ilustrarea proprietății precedente, scriem o funcție Matlab, care generează o matrice de tipul (n, m) , ale cărei elemente sunt numere aleatoare ce urmează o lege de probabilitate precizată. Pentru fiecare coloană a matricei generate, se calculează valoarea z_n , folosind formula (4.2.1), după care se reprezintă grafic histograma valorilor z_n , împreună cu densitatea de probabilitate a legii normale $\mathcal{N}(\bar{z}, \bar{\sigma})$.

```
function hist2(n,m,lege)
% Comportare asimptotica a mediei de selectie
% lege - legea de probabilitate
% (n,m) - tipul matricei generate
switch lege
case 'unid'
    N=input('N=');
    X=unidrnd(N,n,m); [med,var]=unidstat(N);
case 'bino'
    nn=input('n='); p=input('p=');
    X=binornd(nn,p,n,m); [med,var]=binostat(nn,p);
case 'hyge'
    M=input('M='); K=input('K='); nn=input('n=');
    X=hygernd(M,K,nn,n,m); [med,var]=hygestat(M,K,nn);
case 'poiss'
    la=input('lambda=');
    X=poissrnd(la,n,m); [med,var]=poisstat(la);
case 'nbin'
    r=input('r='); p=input('p=');
    X=nbinrnd(r,p,n,m); [med,var]=nbinstat(r,p);
case 'geo'
    p=input('p=');
    X=geornd(p,n,m); [med,var]=geostat(p);
case 'unif'
    a=input('a:'); b=input('b(a<b):');
    X=unifrnd(a,b,n,m); [med,var]=unifstat(a,b);
case 'norm'
    mu=input('mu='); s=input('sigma=');
    X=normrnd(mu,s,n,m); [med,var]=normstat(mu,s);
case 'logn'
    mu=input('mu='); s=input('sigma=');
    X=lognrnd(mu,s,n,m); [med,var]=lognstat(mu,s);
case 'gam'
    a=input('a='); b=input('b=');
    X=gamrnd(a,b,n,m); [med,var]=gamstat(a,b);
case 'exp'
    mu=input('mu=');
    X=exprnd(mu,n,m); [med,var]=expstat(mu);
case 'beta'
    a=input('a='); b=input('b=');
    X=betarnd(a,b,n,m); [med,var]=betastat(a,b);
case 'weib'
    a=input('a='); b=input('b=');
```

```

X=weibrnd(a,b,n,m); [med,var]=weibstat(a,b);
case 'rayl'
    b=input('b=');
    X=raylrnd(b,n,m); [med,var]=raylstat(b);
case 't'
    nn=input('n=');
    X=trnd(nn,n,m); [med,var]=tstat(nn);
case 'chi2'
    nn=input('n=');
    X=chi2rnd(nn,n,m); [med,var]=chi2stat(nn);
case 'f'
    mm=input('m='); nn=input('n=');
    X=frnd(mm,nn,n,m); [med,var]=fstat(mm,nn);
otherwise
    error('Lege necunoscuta')
end
Z=mean(X); Z=(Z-med)/sqrt(var/n);
hist(Z,fix(1+10/3*log10(m)), hold on
x=-3:0.01:3; f=normpdf(x,0,1);
plot(x,f,'k-'), colormap spring

```

Apelul funcției hist2 se face prin

```
>>hist2(n,m,'lege')
```

unde valorile pentru n, m și lege, fie că sunt precizate în acest apel, fie sunt precizate înainte. De exemplu, comanda

```
>>hist2(25,35,'unif')
```

are ca efect apelul funcției pentru legea uniformă, iar pe ecran se va cere introducerea parametrilor a și b ($a < b$), după care pe ecran va fi reprezentat graficul din Figura 4.1, în cazul în care $a = -1$ și $b = 1$. Să remarcăm totuși că la un nou apel, cu aceeași parametri, graficul diferă, deoarece sunt generate alte numere aleatoare.

4.2.2 Momente de selecție

Definiția 4.2.9. Numim moment de selecție de ordin k funcția de selecție

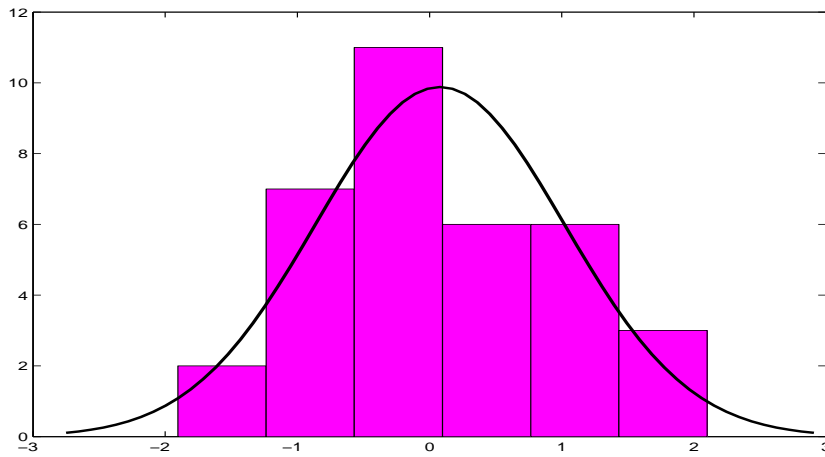
$$\bar{\nu}_k = \frac{1}{n} \sum_{i=1}^n X_i^k, \quad \text{iar} \quad \bar{\nu}_k = \frac{1}{n} \sum_{i=1}^n x_i^k,$$

se numește valoarea momentului de selecție de ordin k .

Observația 4.2.10. Se vede imediat că $\bar{\nu}_1 = \bar{X}$.

Proprietatea 4.2.11. Fie caracteristica X pentru care există momentul teoretic de ordin $2k$, $\nu_{2k} = E(X^{2k})$, atunci

$$E(\bar{\nu}_k) = \nu_k \quad \text{și} \quad Var(\bar{\nu}_k) = \frac{1}{n} (\nu_{2k} - \nu_k^2),$$

Figura 4.1: Legea $\mathcal{N}(0, 1)$

iar pentru $n \rightarrow \infty$, avem că

$$Z_n = \frac{\bar{\nu}_k - \nu_k}{\frac{\sqrt{\nu_{2k} - \nu_k^2}}{\sqrt{n}}}$$

converge în repartiție la legea normală $\mathcal{N}(0, 1)$.

Demonstrație. Deoarece selecția este repetată putem scrie succesiv

$$\begin{aligned} E(\bar{\nu}_k) &= \frac{1}{n} \sum_{i=1}^n E(X_i^k) = \frac{1}{n} \sum_{i=1}^n E(X^k) = \frac{1}{n} n \nu_k = \nu_k, \\ \text{Var}(\bar{\nu}_k) &= \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i^k) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X^k) = \frac{1}{n^2} n (\nu_{2k} - \nu_k^2) \\ &= \frac{1}{n} (\nu_{2k} - \nu_k^2). \end{aligned}$$

Ultima afirmație rezultă prin aplicarea Teoremei 2.7.8 șirului $(X_n^k)_{n \geq 1}$ de variabile aleatoare independente și identic repartizate. Anume, deoarece $E(X_n^k) = \nu_k$ și $\text{Var}(X_n^k) = \nu_{2k} - \nu_k^2$, avem

$$\begin{aligned} Z_n &= \sum_{i=1}^n \frac{X_i^k - \nu_k}{\sqrt{\nu_{2k} - \nu_k^2} \sqrt{n}} = \frac{1}{\sqrt{\nu_{2k} - \nu_k^2} \sqrt{n}} \left(\sum_{i=1}^n X_i^k - n \nu_k \right) \\ &= \frac{1}{\sqrt{\nu_{2k} - \nu_k^2} \sqrt{n}} (n \bar{\nu}_k - n \nu_k) = \frac{\bar{\nu}_k - \nu_k}{\frac{\sqrt{\nu_{2k} - \nu_k^2}}{\sqrt{n}}}, \end{aligned}$$

dar Z_n converge în repartiție la legea normală $\mathcal{N}(0, 1)$, ceea ce încheie demonstrația. \square

Definiția 4.2.12. Numim moment centrat de selecție de ordin k funcția de selecție

$$\bar{\mu}_k = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^k, \quad \text{iar} \quad \bar{\mu}_k = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^k,$$

se numește valoarea momentului centrat de selecție de ordin k .

Observația 4.2.13. Se vede imediat că $\bar{\mu}_1 = 0$, iar momentul centrat de ordinul doi se poate scrie $\bar{\mu}_2 = \bar{\nu}_2 - \bar{\nu}_1^2$.

Proprietatea 4.2.14. Fie caracteristica X pentru care există momentul teoretic ν_4 , atunci pentru momentul centrat de ordinul doi avem

$$E(\bar{\mu}_2) = \frac{n-1}{n} \sigma^2 \quad \text{și} \quad Var(\bar{\mu}_2) = \frac{n-1}{n^3} \left[(n-1) \mu_4 - (n-3) \sigma^4 \right],$$

unde $\sigma^2 = Var(X)$, și de asemenea

$$Cov(\bar{X}, \bar{\mu}_2) = \frac{n-1}{n^2} \mu_3.$$

Demonstrație. Pentru prima relație scriem succesiv

$$\begin{aligned} E(\bar{\mu}_2) &= E(\bar{\nu}_2) - E(\bar{\nu}_1^2) = \nu_2 - \frac{1}{n^2} E\left(\sum_{k=1}^n X_k^2 + 2 \sum_{\substack{i,j=1 \\ i < j}}^n X_i X_j\right) \\ &= \nu_2 - \frac{1}{n^2} \left[\sum_{k=1}^n E(X_k^2) + 2 \sum_{\substack{i,j=1 \\ i < j}}^n E(X_i) E(X_j) \right] \\ &= \nu_2 - \frac{1}{n^2} [n\nu_2 + n(n-1)\nu_1^2] = \nu_2 - \frac{1}{n} \nu_2 - \frac{n-1}{n} \nu_1^2 \\ &= \frac{n-1}{n} (\nu_2 - \nu_1^2) = \frac{n-1}{n} \sigma^2. \end{aligned}$$

Dacă se rețin extremitățile șirului de egalități avem prima relație.

Pentru a calcula dispersia, considerăm variabilele aleatoare reduse, notate prin $Y_k = X_k - E(X_k)$, $k = \overline{1, n}$, care sunt independente și identic repartizate și pentru care $E(Y_k) = 0$, $Var(Y_k) = \sigma^2$, $k = \overline{1, n}$.

Se arată ușor că

$$\bar{\mu}_2 = \frac{1}{n} \sum_{k=1}^n (Y_k - \bar{Y})^2, \quad \text{unde} \quad \bar{Y} = \frac{1}{n} \sum_{k=1}^n Y_k.$$

Pe de altă parte avem că

$$Var(\bar{\mu}_2) = E(\bar{\mu}_2^2) - \left(E(\bar{\mu}_2)\right)^2 = E(\bar{\mu}_2^2) - \left(\frac{n-1}{n}\right)^2 \sigma^4.$$

A rămas, prin urmare, de calculat $E(\bar{\mu}_2^2)$. Pentru aceasta avem

$$\begin{aligned} E(\bar{\mu}_2^2) &= \frac{1}{n^2} E\left[\left(\sum_{k=1}^n (Y_k - \bar{Y})^2\right)^2\right] = \frac{1}{n^2} E\left[\left(\sum_{k=1}^n Y_k^2 - 2\bar{Y} \sum_{k=1}^n Y_k + n\bar{Y}^2\right)^2\right] \\ &= \frac{1}{n^2} E\left[\left(\sum_{k=1}^n Y_k^2 - n\bar{Y}^2\right)^2\right] = \frac{1}{n^2} E\left[\left(\sum_{k=1}^n Y_k^2\right)^2 - 2n\bar{Y}^2 \sum_{k=1}^n Y_k^2 + n^2\bar{Y}^4\right], \end{aligned}$$

de unde se obține

$$E(\bar{\mu}_2^2) = \frac{1}{n^2} E\left[\left(\sum_{k=1}^n Y_k^2\right)^2\right] - \frac{2}{n} E\left(\bar{Y}^2 \sum_{k=1}^n Y_k^2\right) + E(\bar{Y}^4).$$

Calculăm pe rând termenii din membrul drept al aceste relații. Astfel avem

$$E\left[\left(\sum_{k=1}^n Y_k^2\right)^2\right] = \sum_{k=1}^n E(Y_k^4) + 2 \sum_{\substack{i,j=1 \\ i < j}}^n E(Y_i^2) E(Y_j^2) = n\mu_4 + n(n-1)\mu_2^2,$$

apoi

$$\begin{aligned} E\left(\bar{Y}^2 \sum_{k=1}^n Y_k^2\right) &= \frac{1}{n^2} E\left[\left(\sum_{i=1}^n Y_i^2 + 2 \sum_{\substack{i,j=1 \\ i < j}}^n Y_i Y_j\right) \left(\sum_{k=1}^n Y_k^2\right)\right] \\ &= \frac{1}{n^2} E\left[\sum_{k=1}^n Y_k^4 + 2 \sum_{\substack{i,j=1 \\ i < j}}^n Y_i^2 Y_j^2 + 2 \sum_{\substack{i,j,k=1 \\ i < j}}^n Y_i Y_j Y_k^2\right] \\ &= \frac{1}{n^2} [n\mu_4 + n(n-1)\mu_2^2] = \frac{1}{n}\mu_4 + \frac{n-1}{n}\mu_2^2, \end{aligned}$$

deoarece $E(Y_i Y_j Y_k^2) = 0$, pentru orice $i, j, k = \overline{1, n}$, $i \neq j$. Pentru ultimul termen avem

$$\begin{aligned} E(\bar{Y}^4) &= \frac{1}{n^4} E\left[\sum_{k=1}^n Y_k^4 + 6 \sum_{\substack{i,j=1 \\ i < j}}^n Y_i^2 Y_j^2 + \dots\right] \\ &= \frac{1}{n^4} [n\mu_4 + 3n(n-1)\mu_2^2] = \frac{1}{n^3}\mu_4 + \frac{3(n-1)}{n^3}\mu_2^2. \end{aligned}$$

Termenii lui $E(\bar{Y}^4)$, care nu au fost luați în considerare, sunt nuli deoarece conțin ca factor pe $E(Y_i) = 0$. Obținem în acest fel că

$$\begin{aligned} E(\bar{\mu}_2^2) &= \frac{1}{n^2} [n\mu_4 + n(n-1)\mu_2^2] - \frac{2}{n} \left[\frac{1}{n}\mu_4 + \frac{n-1}{n}\mu_2^2 \right] + \frac{1}{n^3}\mu_4 + \frac{3(n-1)}{n^3}\mu_2^2 \\ &= \frac{(n-1)^2}{n^3}\mu_4 + \frac{(n-1)(n^2-2n+3)}{n^3}\mu_2^2, \end{aligned}$$

deci pentru dispersie avem succesiv

$$\begin{aligned} Var(\bar{\mu}_2) &= \frac{(n-1)^2}{n^3}\mu_4 + \frac{(n-1)(n^2-2n+3)}{n^3}\mu_2^2 - \frac{(n-1)^2}{n^2}\mu_2^2 \\ &= \frac{(n-1)^2}{n^3}\mu_4 - \frac{(n-1)(n-3)}{n^3}\mu_2^2 = \frac{n-1}{n^3} [(n-1)\mu_4 - (n-3)\mu_2^2], \end{aligned}$$

adică

$$Var(\bar{\mu}_2) = \frac{n-1}{n^3} [(n-1)\mu_4 - (n-3)\mu_2^2].$$

Pentru ultima relație putem considera $E(X) = 0$, deci $E(\bar{X}) = 0$, caz în care avem

$$\begin{aligned} Cov(\bar{X}, \bar{\mu}_2) &= E(\bar{X} \bar{\mu}_2) = \frac{1}{n^2} E \left[\left(\sum_{k=1}^n X_k \right) \left(\sum_{k=1}^n X_k^2 - n\bar{X}^2 \right) \right] \\ &= \frac{1}{n^2} E \left[\left(\sum_{k=1}^n X_k \right) \left(\sum_{i=1}^n X_i^2 \right) \right] - E(\bar{X}^3) = \frac{1}{n^2} \sum_{k=1}^n E(X_k^3) - \frac{1}{n^3} \sum_{k=1}^n E(X_k^3), \end{aligned}$$

deoarece $E(X_i X_j) = E(X_i X_j^2) = 0$, când $i \neq j$. Așadar

$$Cov(\bar{X}, \bar{\mu}_2) = \frac{\mu_3}{n} - \frac{\mu_3}{n^2} = \frac{n-1}{n^2} \mu_3.$$

□

Observația 4.2.15. Din proprietatea precedentă avem că

$$Var(\bar{\mu}_2) = \frac{1}{n} (\mu_4 - \mu_2^2) + O\left(\frac{1}{n^2}\right), \quad \text{deci} \quad Var(\bar{\mu}_2) \cong \frac{1}{n} (\mu_4 - \mu_2^2).$$

Deoarece $Cov(\bar{X}, \bar{\mu}_2) = \frac{n-1}{n^2} \mu_3$, avem că statisticile \bar{X} și $\bar{\mu}_2$ sunt necorelate pentru $n \rightarrow \infty$, iar dacă X are distribuția simetrică ($\mu_3 = 0$), atunci \bar{X} și $\bar{\mu}_2$ sunt necorelate pentru orice n .

De asemenea se poate arăta că statistica

$$Z_n = \frac{\bar{\mu}_2 - \mu_2}{\frac{\sqrt{\mu_4 - \mu_2^2}}{\sqrt{n}}}$$

converge în repartiție la legea normală $\mathcal{N}(0, 1)$, când $n \rightarrow \infty$.

Definiția 4.2.16. Numim dispersie de selecție funcția de selecție

$$\bar{\sigma}^2 = \frac{1}{n-1} \sum_{k=1}^n (X_k - \bar{X})^2, \text{ iar valoarea numerică } \bar{\sigma}^2 = \frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x})^2,$$

se numește valoarea dispersiei de selecție.

Observația 4.2.17. Între momentul centrat de selecție de ordinul doi și dispersia de selecție există relația

$$(4.2.2) \quad \bar{\sigma}^2 = \frac{n}{n-1} \bar{\mu}_2.$$

Ca urmare avem că

$$\begin{aligned} E(\bar{\sigma}^2) &= \mu_2 = \sigma^2, \quad Cov(\bar{X}, \bar{\sigma}^2) = \frac{1}{n} \mu_3, \\ Var(\bar{\sigma}^2) &= \frac{1}{n(n-1)} \left[(n-1) \mu_4 - (n-3) \mu_2^2 \right]. \end{aligned}$$

Proprietatea 4.2.18. Fie caracteristica X pentru care există momentul centrat teoretic $\mu_k = E\left[\left(X - E(X)\right)^k\right]$, atunci $E(\bar{\mu}_k) = \mu_k + O\left(\frac{1}{n}\right)$.

Demonstrație. Folosind formula binomului avem că

$$\begin{aligned} \bar{\mu}_k &= \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^k = \frac{1}{n} \sum_{i=1}^n \left[X_i^k - \binom{k}{1} X_i^{k-1} \bar{X} + \cdots + (-1)^k \bar{X}^k \right] \\ &= \bar{\nu}_k - \binom{k}{1} \bar{X} \bar{\nu}_{k-1} + \cdots + (-1)^k \bar{X}^k. \end{aligned}$$

Fără a restrânge generalitatea, se poate considera că $E(X) = 0$, deci au loc egalitățile $\mu_i = \nu_i$, $i = \overline{1, k}$. Putem scrie atunci

$$E(\bar{\mu}_k) = E(\bar{\nu}_k) - \binom{k}{1} E(\bar{X} \bar{\nu}_{k-1}) + \cdots + (-1)^k E(\bar{X}^k)$$

$$= \mu_k - \binom{k}{1} E(\bar{X} \bar{\nu}_{k-1}) + \cdots + (-1)^k E(\bar{X}^k).$$

Pe de altă parte, din inegalitatea lui Schwarz, se obține că

$$E^2(\bar{X}^i \bar{\nu}_{k-i}) \leq E(\bar{X}^{2i}) E(\bar{\nu}_{k-i}^2), \quad i = \overline{2, k}.$$

Dacă se are în vedere Proprietatea 4.2.11, rezultă că

$$\begin{aligned} E(\bar{\nu}_{k-i}^2) &= Var(\bar{\nu}_{k-i}) + E^2(\bar{\nu}_{k-i}) = \frac{1}{n} [\nu_{2(k-i)} - \nu_{k-i}^2] + \nu_{k-i}^2 \\ &= \mu_{k-i}^2 + \frac{1}{n} [\mu_{2(k-i)} - \mu_{k-i}^2]. \end{aligned}$$

Se poate arăta, de asemenea, că $E(\bar{X}^{2i}) = O\left(\frac{1}{n^i}\right)$. Astfel din inegalitatea lui Schwarz se obține

$$E^2(\bar{X}^i \bar{\nu}_{k-i}) \leq O\left(\frac{1}{n^i}\right) \quad \text{sau} \quad E(\bar{X}^i \bar{\nu}_{k-i}) \leq O\left(\frac{1}{n^{\frac{i}{2}}}\right), \quad i = \overline{2, k}.$$

În cazul $i = 1$, se calculează direct

$$E(\bar{X} \bar{\nu}_{k-1}) = \frac{1}{n^2} E\left[\left(\sum_{i=1}^n X_i\right)\left(\sum_{j=1}^n X_j^{k-1}\right)\right] = \frac{1}{n^2} \sum_{i=1}^n E(X_i^k) = \frac{1}{n} \mu_k.$$

Luând în considerare aceste evaluări, se obține în final

$$E(\bar{\mu}_k) = \mu_k - \frac{k}{n} \mu_k + O\left(\frac{1}{n}\right) = \mu_k + O\left(\frac{1}{n}\right).$$

□

Observația 4.2.19. Printr-o cale analogă, se poate obține că

$$Var(\bar{\mu}_k) = \frac{\mu_{2k} - 2k\mu_{k-1}\mu_{k+1} - \mu_k^2 + k^2\mu_k\mu_{k-1}^2}{n} + O\left(\frac{1}{n^2}\right).$$

4.2.3 Coeficient de corelație de selecție

Definiția 4.2.20. Fie caracteristica bidimensională (X, Y) și o selecție repetată de volum n , cu datele de selecție (x_k, y_k) , $k = \overline{1, n}$, și respectiv variabilele de selecție (X_k, Y_k) , $k = \overline{1, n}$. Numim coeficient de corelație de selecție funcția de selecție

$$\bar{r} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}},$$

iar valoarea numerică

$$\bar{r} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}},$$

se numește valoarea coeficientului de corelație de selecție.

Observația 4.2.21. Fisher a arătat că pentru o caracteristică bidimensională (X, Y) , care urmează legea normală $\mathcal{N}(m_1, m_2; \sigma_1, \sigma_2; r)$, densitatea de probabilitate a coeficientului de corelație de selecție este

$$\begin{aligned} f_n(x) &= \frac{n-2}{\pi} \sqrt{(1-r^2)^{n-1}} \sqrt{(1-x^2)^{n-4}} \int_0^1 \frac{u^{n-2}}{(1-rxu)^{n-1}} \frac{du}{\sqrt{1-u^2}} \\ &= \frac{2^{n-3}}{\pi (n-3)!} \sqrt{(1-r^2)^{n-1}} \sqrt{(1-x^2)^{n-4}} \times \\ &\quad \times \sum_{k=0}^{\infty} F^2\left(\frac{n+k-1}{2}\right) \frac{(2rx)^k}{k!}, \quad x \in (-1, 1). \end{aligned}$$

De asemenea, dacă se consideră transformarea (lui Fisher) $Z = \frac{1}{2} \ln \frac{1+\bar{r}}{1-\bar{r}}$ și dacă $\rho = \frac{1}{2} \ln \frac{1+r}{1-r}$, atunci Z urmează aproximativ legea normală $\mathcal{N}\left(\rho + \frac{r}{2(n-1)}, \frac{1}{\sqrt{n-3}}\right)$.

În cazul particular $r = 0$, avem densitatea de probabilitate a lui \bar{r} exprimată prin

$$f_n(x) = \frac{1}{\sqrt{\pi}} \frac{\Gamma\left(\frac{n-1}{2}\right)}{\Gamma\left(\frac{n-2}{2}\right)} \sqrt{(1-x^2)^{n-4}}, \quad x \in (-1, 1),$$

iar statistica

$$(4.2.3) \quad T = \sqrt{n-2} \frac{\bar{r}}{\sqrt{1-\bar{r}^2}}$$

urmează legea Student cu $n-2$ grade de libertate.

Programul 4.2.22. Vom scrie un program, care generează de N ori câte n vectori aleatori, ce urmează legea normală bidimensională, având coeficientul de corelație dintre cele două componente nul, adică $r=0$. Folosind aceste date se vor calcula N numere aleatoare după regula precizată prin statistica (4.2.3), iar histograma corespunzătoare acestor noi date va fi reprezentată grafic împreună cu densitatea de probabilitate a legii Student cu $n-2$ grade de libertate.

```
clear,clf,
mu(1)=input('m1='); mu(2)=input('m2=');
v(1,1)=input('sigma1^2=');
v(2,2)=input('sigma2^2=');
v(1,2)=0; v(2,1)=0;
```

```

n=input('n='); N=input('N=');
for k=1:N
Z=mvnrnd(mu,v,n); r=corrcoef(Z); r12=r(1,2);
t(k)=sqrt(n-2)*r12/sqrt(1-r12^2);
end
x=-3:0.01:3; f=tpdf(x,n-2);
nn=fix(1+10/3*log10(N));
[fr,cl]=hist(t,nn); h=cl(2)-cl(1);
bar(cl,fr/(h*N),1),hold on, plot(x,f,'k-')
colormap spring

```

Pentru $N=50, n=20, \mu=(5, 10), v=\begin{pmatrix} 2 & 0 \\ 0 & 3 \end{pmatrix}$, se obțin graficele din Figura 4.2.

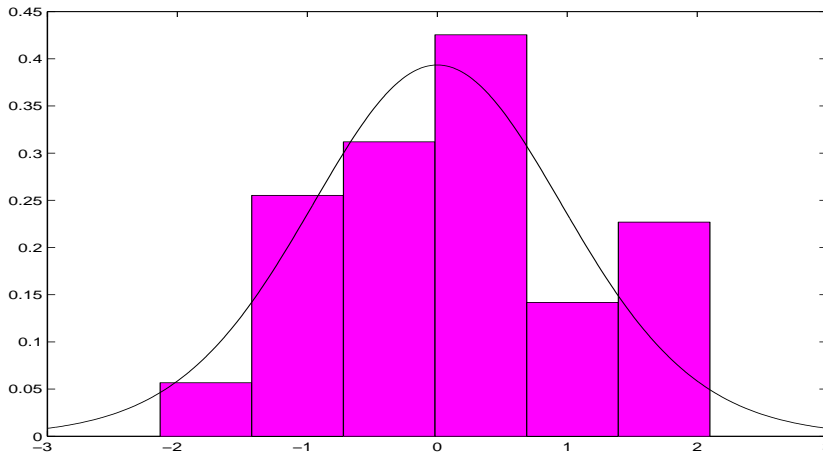


Figura 4.2: Legea $\mathcal{T}(18)$

Lema 4.2.23 (Fisher). Dacă variabilele aleatoare X_1, X_2, \dots, X_n sunt independente, fiecare urmând legea normală $\mathcal{N}(0, 1)$ și dacă se consideră matricea ortogonală $A = (a_{ij})_{i,j=\overline{1,n}}$, atunci variabilele aleatoare

$$Y_i = \sum_{k=1}^n a_{ik} X_k, \quad i = \overline{1, n},$$

sunt independente, fiecare urmând legea normală $\mathcal{N}(0, 1)$.

Demonstrație. Pentru a arăta că variabilele aleatoare $Y_i, i = \overline{1, n}$, sunt independente determinăm funcția caracteristică φ a vectorului aleator (Y_1, Y_2, \dots, Y_n) . Pornind de

la definiția funcției caracteristice avem

$$\begin{aligned}\varphi(t_1, t_2, \dots, t_n) &= E \left[\exp \left(i \sum_{k=1}^n t_k Y_k \right) \right] = E \left[\exp \left(i \sum_{k=1}^n t_k \sum_{j=1}^n a_{kj} X_j \right) \right] \\ &= E \left[\exp \left(i \sum_{j=1}^n \left(\sum_{k=1}^n t_k a_{kj} \right) X_j \right) \right].\end{aligned}$$

Variabilele aleatoare X_j , $j = \overline{1, n}$, fiind independente se obține

$$\varphi(t_1, t_2, \dots, t_n) = \prod_{j=1}^n E \left[\exp \left(i \sum_{k=1}^n t_k a_{kj} \right) X_j \right].$$

Fiecare factor din membrul drept este valoarea funcției caracteristice a legii normale $\mathcal{N}(0, 1)$, pe punctele, respectiv $\sum_{k=1}^n t_k a_{kj}$, $j = \overline{1, n}$, adică

$$\begin{aligned}\varphi(t_1, t_2, \dots, t_n) &= \prod_{j=1}^n \exp \left[-\frac{1}{2} \left(\sum_{k=1}^n t_k a_{kj} \right)^2 \right] = \exp \left[-\frac{1}{2} \sum_{j=1}^n \left(\sum_{k=1}^n t_k a_{kj} \right)^2 \right] \\ &= \exp \left[-\frac{1}{2} \sum_{j=1}^n \left(\sum_{k=1}^n t_k^2 a_{kj}^2 + 2 \sum_{\substack{i,s=1 \\ i < s}}^n t_i t_s a_{ij} a_{sj} \right) \right] \\ &= \exp \left[-\frac{1}{2} \left(\sum_{k=1}^n t_k^2 \sum_{j=1}^n a_{kj}^2 + 2 \sum_{\substack{i,s=1 \\ i < s}}^n t_i t_s \sum_{j=1}^n a_{ij} a_{sj} \right) \right].\end{aligned}$$

Folosind condițiile de ortonormalitate relative la coeficienții a_{ij} , se obține că

$$\varphi(t_1, t_2, \dots, t_n) = \exp \left(-\frac{1}{2} \sum_{k=1}^n t_k^2 \right),$$

adică faptul că Y_1, Y_2, \dots, Y_n sunt variabile aleatoare independente, ce urmează fiecare legea normală $\mathcal{N}(0, 1)$. \square

Proprietatea 4.2.24. Fie caracteristica X , ce urmează legea normală $\mathcal{N}(0, 1)$ și variabilele de selecție X_1, X_2, \dots, X_n , ce corespund unei selecții repetate de volum n , atunci statisticile

$$U_n = \sqrt{n} \bar{X} = \frac{1}{\sqrt{n}} \sum_{k=1}^n X_k, \quad V_n = \sum_{k=1}^n (X_k - \bar{X})^2,$$

sunt variabile aleatoare independente, ce urmează legea normală $\mathcal{N}(0, 1)$ și respectiv legea χ^2 cu $n - 1$ grade de libertate.

Demonstrație. Se consideră matricea ortonormată

$$A = \begin{pmatrix} \frac{1}{\sqrt{n}} & \frac{1}{\sqrt{n}} & \frac{1}{\sqrt{n}} & \frac{1}{\sqrt{n}} & \cdots & \frac{1}{\sqrt{n}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} & 0 & 0 & \cdots & 0 \\ \frac{1}{\sqrt{6}} & \frac{1}{\sqrt{6}} & -\frac{2}{\sqrt{6}} & 0 & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ \frac{1}{\sqrt{n(n-1)}} & \frac{1}{\sqrt{n(n-1)}} & \frac{1}{\sqrt{n(n-1)}} & \frac{1}{\sqrt{n(n-1)}} & \cdots & -\frac{n-1}{\sqrt{n(n-1)}} \end{pmatrix}$$

pentru care se aplică Lema 4.2.23.

Pe de o parte avem că

$$Y_1 = \sum_{k=1}^n a_{1k} X_k = \frac{1}{\sqrt{n}} \sum_{k=1}^n X_k = U_n$$

și

$$\begin{aligned} Y_2^2 + \cdots + Y_n^2 &= (Y_1^2 + Y_2^2 + \cdots + Y_n^2) - Y_1^2 = (X_1^2 + \cdots + X_n^2) - U_n^2 \\ &= \sum_{k=1}^n (X_k - \bar{X})^2 = V_n. \end{aligned}$$

Deoarece Y_1, Y_2, \dots, Y_n sunt independente rezultă că U_n și V_n sunt independente.

Pe de altă parte $U_n = Y_1$ urmează legea normală $\mathcal{N}(0, 1)$, iar V_n fiind suma pătratelor a $n - 1$ variabile aleatoare independente, fiecare urmând legea normală $\mathcal{N}(0, 1)$, obținem că V_n urmează legea χ^2 cu $n - 1$ grade de libertate (a se vedea legea χ^2 prezentată în Capitolul 2). \square

Proprietatea 4.2.25. Fie caracteristica X , ce urmează legea normală $\mathcal{N}(m, \sigma)$ și variabilele de selecție X_1, X_2, \dots, X_n , ce corespund unei selecții repetate de volum n , atunci statisticile

$$U_n = \frac{\bar{X} - m}{\frac{\sigma}{\sqrt{n}}}, \quad V_n = \frac{1}{\sigma^2} \sum_{k=1}^n (X_k - \bar{X})^2,$$

sunt variabile aleatoare independente, ce urmează respectiv legea normală $\mathcal{N}(0, 1)$ și legea χ^2 cu $n - 1$ grade de libertate.

Demonstrație. Se consideră variabilele aleatoare

$$Z_k = \frac{X_k - m}{\sigma}, \quad k = \overline{1, n},$$

care sunt variabile aleatoare independente ce urmează fiecare legea normală $\mathcal{N}(0, 1)$. Prin aplicarea Proprietății 4.2.24 pentru variabilele aleatoare Z_k , $k = \overline{1, n}$, se obține afirmația făcută. Într-adevăr avem că

$$\frac{1}{\sqrt{n}} \sum_{k=1}^n Z_k = \frac{1}{\sqrt{n}} \sum_{k=1}^n \frac{X_k - m}{\sigma} = U_n$$

și

$$\begin{aligned} \sum_{k=1}^n (Z_k - \bar{Z})^2 &= \sum_{k=1}^n \left(\frac{X_k - m}{\sigma} - \frac{1}{n} \sum_{i=1}^n \frac{X_i - m}{\sigma} \right)^2 = \frac{1}{\sigma^2} \sum_{k=1}^n \left(X_k - \frac{1}{n} \sum_{i=1}^n X_i \right)^2 \\ &= \frac{1}{\sigma^2} \sum_{k=1}^n (X_k - \bar{X})^2 = V_n. \end{aligned}$$

□

Proprietatea 4.2.26. Fie caracteristica X , ce urmează legea normală $\mathcal{N}(m, \sigma)$ și variabilele de selecție X_1, X_2, \dots, X_n , ce corespund unei selecții repetate de volum n , atunci statistica

$$(4.2.4) \quad T = \frac{\bar{X} - m}{\frac{\bar{\sigma}}{\sqrt{n}}} = \frac{\bar{X} - m}{\sqrt{\frac{\mu_2}{n-1}}},$$

urmează legea Student cu $n - 1$ grade de libertate.

Demonstrație. Cu notațiile de la Proprietatea 4.2.25, arătăm că

$$T = \frac{U_n}{\sqrt{\frac{V_n}{n-1}}}.$$

Într-adevăr, avem succesiv

$$\frac{U_n}{\sqrt{\frac{V_n}{n-1}}} = \frac{\bar{X} - m}{\frac{\sigma}{\sqrt{n}}} \cdot \frac{\sqrt{n-1}}{\frac{1}{\sigma} \sqrt{\sum_{k=1}^n (X_k - \bar{X})^2}} = \frac{\bar{X} - m}{\frac{\sqrt{\frac{1}{n-1} \sum_{k=1}^n (X_k - \bar{X})^2}}{\sqrt{n}}} = \frac{\bar{X} - m}{\frac{\bar{\sigma}}{\sqrt{n}}} = T.$$

Din teoria probabilităților se știe că raportul dintre o variabilă aleatoare, ce urmează legea normală $\mathcal{N}(0, 1)$ și radicalul unei variabile aleatoare ce urmează legea χ^2 , raportată la numărul gradelor de libertate, în cazul în care cele două variabile aleatoare sunt independente, este o variabilă aleatoare, ce urmează legea Student cu același număr al gradelor de libertate ca legea χ^2 considerată (a se vedea legile χ^2 și t prezentate în Capitolul 2). Având în vedere cine sunt U_n și V_n rezultă afirmația din enunțul proprietății. □

Programul 4.2.27. Pentru ilustrarea acestui rezultat, programul Matlab, care urmează, generează o matrice de tipul (n, m) de numere aleatoare, ce urmează legea normală $\mathcal{N}(\mu, \sigma)$, după care pentru fiecare coloană a matricei generate, se construiesc datele aleatoare conform statisticii (4.2.4). Pentru aceste date se reprezintă grafic histograma corespunzătoare, împreună cu densitatea de probabilitate a legii Student cu $n-1$ grade de libertate.

```
clear, clf, mu=input('mu='); sigma=input('s=');
n=input('n='); m=input('m=');
Z=normrnd(mu,sigma,n,m);
t=(mean(Z)-mu)./sqrt(var(Z)/n);
x=-3:0.01:3; f=tpdf(x,n-1);
nn=fix(1+10/3*log10(m));
[fr,clasa]=hist(t,nn);
h=clasa(2)-clasa(1);
bar(clasa,fr/(h*m),1)
hold on, plot(x,f,'k-')
colormap([.7,.7,.7])
```

Pentru $n=15$, $m=100$, $\mu=5$ și $\sigma=2$, se obțin graficele din Figura 4.3.

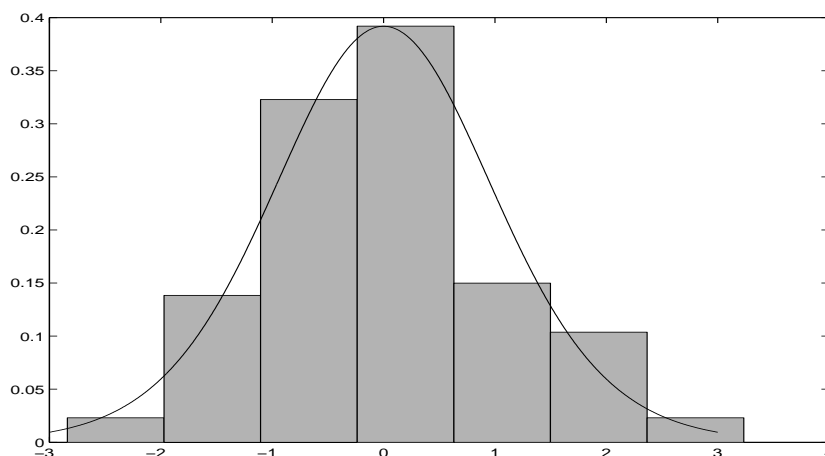


Figura 4.3: Legea $\mathcal{T}(14)$

Proprietatea 4.2.28. Fie caracteristicile independente X' și X'' , fiecare urmând legea normală, respectiv $\mathcal{N}(m', \sigma)$ și $\mathcal{N}(m'', \sigma)$ și variabilele de selecție $X'_1, \dots, X'_{n'}$, respectiv $X''_1, \dots, X''_{n''}$, ce corespund unei selecții repetate de volum n' pentru caracteristica X' și unei selecții repetate de volum n'' pentru caracteristica

X'' , atunci statistica

$$(4.2.5) \quad T = \frac{(\bar{X}' - \bar{X}'') - (m' - m'')}{\sqrt{(n' - 1)\bar{\sigma}'^2 + (n'' - 1)\bar{\sigma}''^2}} \sqrt{\frac{n' + n'' - 2}{\frac{1}{n'} + \frac{1}{n''}}},$$

unde

$$\bar{X}' = \frac{1}{n'} \sum_{k=1}^{n'} X'_k, \quad \bar{X}'' = \frac{1}{n''} \sum_{k=1}^{n''} X''_k,$$

$$\bar{\sigma}'^2 = \frac{1}{n' - 1} \sum_{k=1}^{n'} (X'_k - \bar{X}')^2, \quad \bar{\sigma}''^2 = \frac{1}{n'' - 1} \sum_{k=1}^{n''} (X''_k - \bar{X}'')^2,$$

urmează legea Student cu $n' + n'' - 2$ grade de libertate.

Demonstrație. Conform Observației 4.2.6, avem că mediile de selecție \bar{X}' și \bar{X}'' urmează fiecare legea normală respectiv $\mathcal{N}\left(m', \frac{\sigma}{\sqrt{n'}}\right)$ și $\mathcal{N}\left(m'', \frac{\sigma}{\sqrt{n''}}\right)$. Prin urmare statistica

$$U = \frac{(\bar{X}' - \bar{X}'') - (m' - m'')}{\sigma \sqrt{\frac{1}{n'} + \frac{1}{n''}}},$$

urmează legea normală $\mathcal{N}(0, 1)$.

Pe de altă parte, folosind Proprietatea 4.2.25, se obține că statistica

$$V = \frac{1}{\sigma^2} \sum_{k=1}^{n'} (X'_k - \bar{X}')^2 + \frac{1}{\sigma^2} \sum_{k=1}^{n''} (X''_k - \bar{X}'')^2,$$

urmează legea χ^2 cu $n' + n'' - 2$ grade de libertate, fiind suma a două variabile aleatoare independente, ce urmează legea χ^2 cu $n' - 1$ grade de libertate și respectiv $n'' - 1$ grade de libertate.

Ca și la demonstrația Proprietății 4.2.26, statistica $U / \sqrt{\frac{V}{n' + n'' - 2}}$ urmează legea Student cu $n' + n'' - 2$ grade de libertate. Mai rămâne de arătat că această statistică este chiar statistica T . Într-adevăr, avem succesiv

$$\begin{aligned} \frac{U}{\sqrt{\frac{V}{n' + n'' - 2}}} &= \frac{(\bar{X}' - \bar{X}'') - (m' - m'')}{\sigma \sqrt{\frac{1}{n'} + \frac{1}{n''}}} \times \\ &\quad \times \frac{\sqrt{n' + n'' - 2}}{\frac{1}{\sigma} \sqrt{\sum_{k=1}^{n'} (X'_k - \bar{X}')^2 + \sum_{k=1}^{n''} (X''_k - \bar{X}'')^2}} \end{aligned}$$

$$= \frac{(\bar{X}' - \bar{X}'') - (m' - m'')}{\sqrt{(n' - 1)\bar{\sigma}'^2 + (n'' - 1)\bar{\sigma}''^2}} \sqrt{\frac{n' + n'' - 2}{\frac{1}{n'} + \frac{1}{n''}}} = T,$$

ceea ce trebuie demonstrat. \square

Observația 4.2.29. Dacă se consideră caracteristicile independente X' și X'' , fiecare urmând legea normală $\mathcal{N}(m', \sigma')$ și respectiv $\mathcal{N}(m'', \sigma'')$ și dacă avem variabilele de selecție $X'_1, X'_2, \dots, X'_{n'}$, ce corespund unei selecții repetate de volum n' relativă la caracteristica X' și variabilele de selecție $X''_1, X''_2, \dots, X''_{n''}$, ce corespund unei selecții repetate de volum n'' relativă la caracteristica X'' , atunci statistica

$$Z = \frac{(\bar{X}' - \bar{X}'') - (m' - m'')}{\sqrt{\frac{\sigma'^2}{n'} + \frac{\sigma''^2}{n''}}},$$

urmează legea normală $\mathcal{N}(0, 1)$.

Proprietatea 4.2.30. Fie caracteristicile independente X' și X'' , fiecare urmând legea normală, respectiv $\mathcal{N}(m', \sigma')$ și $\mathcal{N}(m'', \sigma'')$ și variabilele de selecție $X'_1, \dots, X'_{n'}$, respectiv $X''_1, \dots, X''_{n''}$, ce corespund unei selecții repetate de volum n' pentru caracteristica X' și unei selecții repetate de volum n'' pentru caracteristica X'' , atunci statistica

$$(4.2.6) \quad F = \frac{\bar{\sigma}'^2}{\sigma'^2} \bigg/ \frac{\bar{\sigma}''^2}{\sigma''^2}$$

urmează legea Fisher–Snedecor cu $m = n' - 1$ și $n = n'' - 1$ grade de libertate.

Demonstrație. Din Proprietatea 4.2.25, avem că funcțiile de selecție

$$V' = (n' - 1) \frac{\bar{\sigma}'^2}{\sigma'^2}, \quad V'' = (n'' - 1) \frac{\bar{\sigma}''^2}{\sigma''^2}$$

urmează fiecare legea χ^2 cu $m = n' - 1$ și $n = n'' - 1$ grade de libertate.

Pe de altă parte, deoarece X' și X'' sunt independente, V' și V'' sunt independente. Din calculul probabilităților (a se vedea legea f prezentată în Capitolul 2) este cunoscut că raportul a două variabile aleatoare independente, ce urmează legea χ^2 , raportate fiecare la numărul gradelor de libertate corespunzător, este o variabilă aleatoare, ce urmează legea Fisher–Snedecor cu numărul gradelor de libertate date de numerele gradelor de libertate ale celor două legi χ^2 . Așadar avem că

$$F = \frac{\bar{\sigma}'^2}{\sigma'^2} \bigg/ \frac{\bar{\sigma}''^2}{\sigma''^2} = \frac{V'}{n' - 1} \bigg/ \frac{V''}{n'' - 1}$$

urmează legea Fisher–Snedecor cu $m = n' - 1$ și $n = n'' - 1$ grade de libertate. \square

Observația 4.2.31. Dacă F urmează legea Fisher–Snedecor cu m și n grade de libertate, atunci $\frac{1}{F}$ urmează legea Fisher–Snedecor cu n și m grade de libertate. Prin urmare, dacă notăm cu $f_{m,n;1-\alpha}$ cuantila de ordin $1 - \alpha$ a lui F , adică, numărul pentru care are loc $P(F \leq f_{m,n;1-\alpha}) = 1 - \alpha$, atunci avem că $\frac{1}{f_{n,m;\alpha}} = f_{m,n;1-\alpha}$, unde $f_{n,m;\alpha}$ este cuantila de ordin α a lui $\frac{1}{F}$. Într-adevăr avem că relația care definește pe $f_{n,m;\alpha}$, $P\left(\frac{1}{F} \leq f_{n,m;\alpha}\right) = \alpha$, este echivalentă cu $P\left(F \geq \frac{1}{f_{n,m;\alpha}}\right) = \alpha$ sau cu relația $P\left(F \leq \frac{1}{f_{n,m;\alpha}}\right) = 1 - \alpha$. Comparând ultima relație cu relația care definește cuantila $f_{m,n;1-\alpha}$, se obține relația dintre cele două cuantile mai sus precizată.

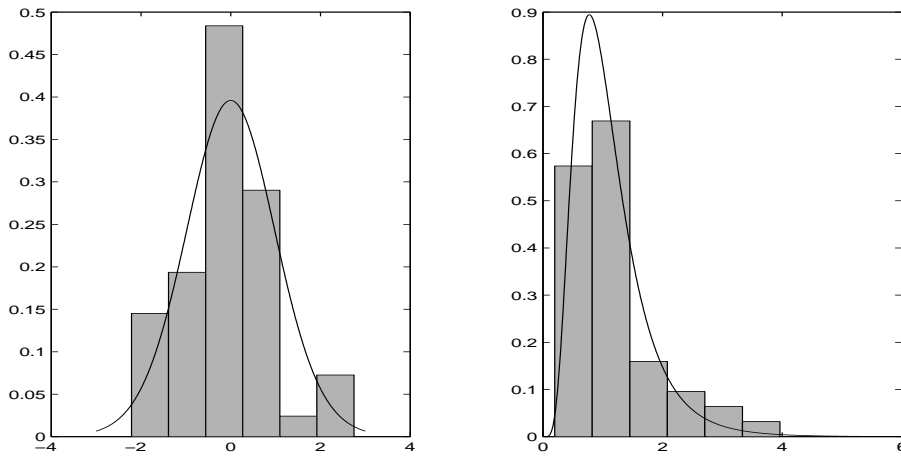
Programul 4.2.32. Pentru ilustrarea Proprietăților 4.2.28 și 4.2.30, vom genera două matrice de tipurile (n_1, m) și (n_2, m) , ale căror elemente sunt numere aleatoare, ce urmează respectiv legile normale $\mathcal{N}(\mu_1, \sigma_1)$ și $\mathcal{N}(\mu_2, \sigma_2)$. Folosind cele două matrice se vor determina, pentru fiecare pereche de coloane cu același număr de ordine, din cele două matrice, datele construite conform definițiilor statisticilor (4.2.5) și (4.2.6). Datele astfel obținute se reprezintă grafic folosind histogramele corespunzătoare, împreună cu densitatea legii Student cu n_1+n_2-2 grade de libertate, respectiv a legii Fisher–Snedecor cu n_1-1 și n_2-1 grade de libertate.

```
clear,clf
mu1=input('mu1='); s1=input('sigma1=');
mu2=input('mu2='); s2=input('sigma2=');
n1=input('n1='); n2=input('n2='); m=input('m=');
X1=normrnd(mu1,s1,n1,m); X2=normrnd(mu2,s2,n2,m);
t=(mean(X1)-mean(X2)-mu1+mu2);
t=t./sqrt((n1-1)*var(X1)+(n2-1)*var(X2));
t=t*sqrt((n1+n2-2)/(1/n1+1/n2));
f=(var(X1)/s1^2)./(var(X2)/s2^2);
x1=-3:0.01:3; x2=0:0.01:6;
f1=tpdf(x1,n1+n2-2); f2=fpdf(x2,n1-1,n2-1);
nn=fix(1+10/3*log10(m));
[fr1,c1]=hist(t,nn); [fr2,c2]=hist(f,nn);
h1=c1(2)-c1(1); h2=c2(2)-c2(1);
subplot(1,2,1)
bar(c1,fr1/(m*h1),1), hold on, plot(x1,f1,'k-')
hold off
subplot(1,2,2)
bar(c2,fr2/(m*h2),1), hold on, plot(x2,f2,'k-')
colormap([.7,.7,.7])
```

Pentru $n_1=15$, $n_2=20$, $m=50$, $\mu_1=5$, $\mu_2=10$, $\sigma_1=2$ și $\sigma_2=3$, se obțin graficele din Figura 4.4.

4.2.4 Funcție de repartiție de selecție

Definiția 4.2.33. Fie caracteristica X cercetată, datele de selecție x_1, x_2, \dots, x_n și variabilele de selecție X_1, X_2, \dots, X_n . Numim funcție de repartiție de selecție,

Figura 4.4: Legea $\mathcal{T}(33)$ și legea $\mathcal{F}(14, 19)$

funcția de selecție definită prin

$$\bar{F}_n(x) = \frac{\nu_n(x)}{n}, \quad x \in \mathbb{R},$$

unde

$$\nu_n(x) = \text{card}\{X_i \mid X_i \leq x, i = \overline{1, n}\},$$

iar

$$\bar{F}_n(x) = \frac{\text{card}\{x_i \mid x_i \leq x, i = \overline{1, n}\}}{n}, \quad x \in \mathbb{R},$$

se numește valoarea funcției de repartiție de selecție.

Observația 4.2.34. Valoarea funcției de repartiție selecție este o funcție în scară de o variabilă reală. Dacă datele de selecție sunt ordonate crescător, atunci

$$\bar{F}_n(x) = \begin{cases} 0, & \text{dacă } x < x_1, \\ \frac{k}{n}, & \text{dacă } x_{k-1} \leq x < x_k, \quad k = \overline{1, n}, \\ 1, & \text{dacă } x \geq x_n. \end{cases}$$

De asemenea, funcția de repartiție de selecție, se vede că este o variabilă aleatoare de tip discret și care are distribuția

$$\bar{F}_n(x) \left(\binom{n}{k} (F(x))^k (1 - F(x))^{n-k} \right)_{k=\overline{0, n}}.$$

Pentru demonstrația unui rezultat fundamental al statisticii matematice, descoperit de matematicianul rus Glivenko, vom prezenta două leme.

Lema 4.2.35. Fie șirul de evenimente $(A_n)_{n \geq 1}$ astfel încât $P(A_n) = 1$, $n \geq 1$, atunci, pentru evenimentul $A = \bigcap_{n=1}^{\infty} A_n$, avem $P(A) = 1$.

Demonstrație. Arătăm la început că proprietatea enunțată are loc pentru un număr finit m de evenimente. Demonstrația o facem prin inducție după m .

Pentru $m = 2$, avem $P(A_1 \cup A_2) \geq P(A_1) = 1$, de unde $P(A_1 \cup A_2) = 1$, care înlocuită în relația

$$P(A_1 \cup A_2) = P(A_1) + P(A_2) - P(A_1 \cap A_2),$$

conduce la $P(A_1 \cap A_2) = 1$.

Presupunem că proprietatea are loc pentru m și arătăm că este adevărată și pentru $m + 1$. Din nou avem că $P\left(\bigcup_{n=1}^{m+1} A_n\right) \geq P(A_{m+1}) = 1$, deci $P\left(\bigcup_{n=1}^{m+1} A_n\right) = 1$. Scriem formula lui Poincaré și ținem seama de ipoteza inducției, ceea ce conduce la

$$\begin{aligned} 1 = P\left(\bigcup_{n=1}^{m+1} A_n\right) &= \sum_{i=1}^{m+1} P(A_i) - \sum_{\substack{i,j=1 \\ i < j}}^{m+1} P(A_i \cap A_j) + \cdots + (-1)^m P\left(\bigcap_{n=1}^{m+1} A_n\right) \\ &= \binom{m+1}{1} - \binom{m+1}{2} + \cdots + (-1)^{m-1} \binom{m+1}{m} + (-1)^m P\left(\bigcap_{n=1}^{m+1} A_n\right). \end{aligned}$$

Reținem extremitățile și vom obține

$$(-1)^m P\left(\bigcap_{n=1}^{m+1} A_n\right) = \binom{m+1}{0} - \binom{m+1}{1} + \binom{m+1}{2} - \cdots + (-1)^m \binom{m+1}{m}$$

sau

$$(-1)^m P\left(\bigcap_{n=1}^{m+1} A_n\right) = (1 - 1)^{m+1} - (-1)^{m+1} \binom{m+1}{m+1},$$

adică $P\left(\bigcap_{n=1}^{m+1} A_n\right) = 1$.

Când șirul de evenimente este infinit, avem în vedere scrierea lui A ca următoarea intersecție de evenimente

$$A = A_1 \cap (A_1 \cap A_2) \cap (A_1 \cap A_2 \cap A_3) \cap \cdots$$

unde evenimentele intersecției satisfac relațiile

$$A_1 \supset (A_1 \cap A_2) \supset (A_1 \cap A_2 \cap A_3) \supset \cdots$$

Conform teoremei de continuitate din teoria probabilităților și având în vedere prima parte a demonstrației putem încheia prin $P(A) = \lim_{m \rightarrow \infty} P\left(\bigcap_{n=1}^m A_n\right) = 1$. \square

Din calculul probabilităților este cunoscută forma tare a teoremei lui Bernoulli sau teorema lui Borel și care o dăm sub forma unei leme.

Lema 4.2.36. *Se consideră n repetări independente ale unui experiment. La fiecare repetare evenimentul A apare cu probabilitatea p . Fie k frecvența absolută a apariției evenimentului A în cele n repetări, atunci*

$$P\left(\lim_{n \rightarrow \infty} \left| \frac{k}{n} - p \right| = 0\right) = 1,$$

adică frecvența relativă a apariției evenimentului A converge tare (aproape sigur) la probabilitatea p de apariție a evenimentului A .

Teorema 4.2.37 (Glivenko). *Fie caracteristica X care are funcția de repartiție teoretică F și fie o selecție repetată de volum n relativă la caracteristica X cu variabilele de selecție X_1, X_2, \dots, X_n și funcția de repartiție de selecție \bar{F}_n corespunzătoare, atunci*

$$P\left(\lim_{n \rightarrow \infty} \sup_{x \in \mathbb{R}} |\bar{F}_n(x) - F(x)| = 0\right) = 1,$$

adică funcția de repartiție de selecție converge aproape sigur la funcția de repartiție teoretică.

Demonstrație. Se consideră $r \in \mathbb{N}$ oarecare, dar fixat, și considerăm numerele $x_{r,k}$, $k = 0, r$, definite prin relația

$$x_{r,k} = \inf \left\{ x \in \mathbb{R} \mid F(x) \leq \frac{k}{r} \leq F(x+0) \right\}.$$

Geometric, punctul $x_{r,k}$ se obține după cum este arătat în Figura 4.5 și reprezintă cuantila de ordin $\frac{k}{r}$.

Dacă se consideră evenimentul $B_{r,k} = (X \leq x_{r,k})$, atunci avem

$$P(B_{r,k}) = P(X \leq x_{r,k}) = F(x_{r,k})$$

și

$$\bar{F}_n(x_{r,k}) = \frac{\nu_n(x_{r,k})}{n} = \frac{\text{card}\{X_i \mid X_i \leq x_{r,k}, i = \overline{1, n}\}}{n}.$$

Prin urmare $\bar{F}_n(x_{r,k})$ este frecvența relativă a apariției evenimentului $B_{r,k}$ și conform Lemei 4.2.36 rezultă că

$$P\left(\lim_{n \rightarrow \infty} |\bar{F}_n(x_{r,k}) - F(x_{r,k})| = 0\right) = 1,$$

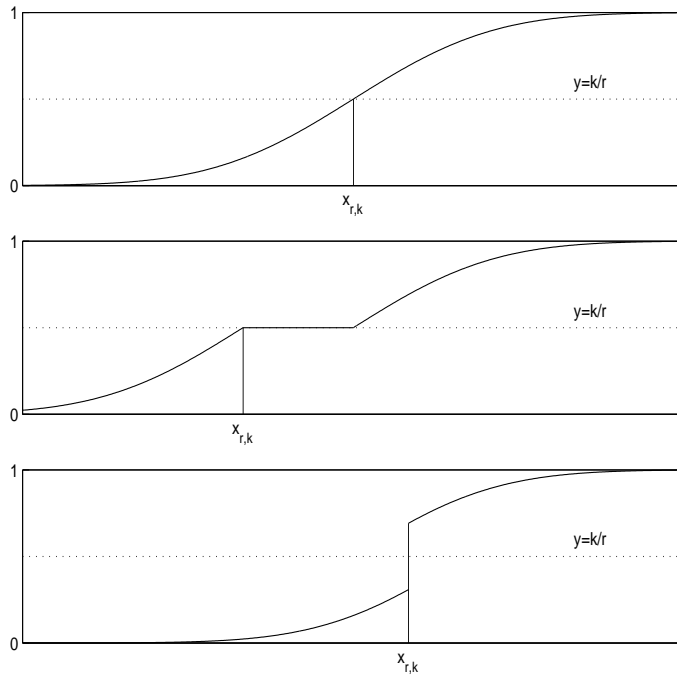


Figura 4.5: Determinarea cuantilei $x_{r,k}$

pentru fiecare $k = \overline{1, r}$. Se poate renunța la valoarea $k = 0$, deoarece $x_{r,0} = -\infty$ și $\bar{F}_n(-\infty) = F(-\infty) = 0$.

Notăm prin $A_{r,k}$ evenimentul

$$|\bar{F}_n(x_{r,k}) - F(x_{r,k})| \rightarrow 0, \quad \text{când } n \rightarrow \infty.$$

Deoarece $P(A_{r,k}) = 1$, $k = \overline{1, r}$, conform Lemei 4.2.35, rezultă că

$$P(A_r) = P\left(\bigcap_{k=1}^r A_{r,k}\right) = 1,$$

unde evenimentul A_r înseamnă

$$\max_{k=\overline{1, r}} |\bar{F}_n(x_{r,k}) - F(x_{r,k})| \rightarrow 0, \quad \text{când } n \rightarrow \infty.$$

În mod analog, considerând evenimentul $C_{r,j} = (X < x_{r,j})$, se obține că

$$P\left(\lim_{n \rightarrow \infty} |\bar{F}_n(x_{r,j} - 0) - F(x_{r,j} - 0)| = 0\right) = 1,$$

pentru fiecare $j = \overline{1, r}$. Dacă se notează prin $D_{r,j}$ evenimentul

$$|\bar{F}_n(x_{r,j} - 0) - F(x_{r,j} - 0)| \longrightarrow 0, \quad \text{când } n \rightarrow \infty.$$

și $D_r = \bigcap_{j=1}^r D_{r,j}$, care înseamnă

$$\max_{j=\overline{1,r}} |\bar{F}_n(x_{r,j} - 0) - F(x_{r,j} - 0)| \longrightarrow 0, \quad \text{când } n \rightarrow \infty,$$

atunci $P(D_r) = 1$.

Dacă se consideră evenimentul $E_r = A_r \cap D_r$, care înseamnă

$$\max_{k,j=\overline{1,r}} \left\{ |\bar{F}_n(x_{r,k}) - F(x_{r,k})|, |\bar{F}_n(x_{r,j} - 0) - F(x_{r,j} - 0)| \right\} \longrightarrow 0,$$

atunci $P(E_r) = 1$. În acest fel, dacă $E = \bigcap_{n=1}^{\infty} E_r$, conform Lemei 4.2.35, avem $P(E) = 1$, adică

$$P\left(\sup_{r \in \mathbb{N}} \max_{k,j=\overline{1,r}} \left\{ |\bar{F}_n(x_{r,k}) - F(x_{r,k})|, |\bar{F}_n(x_{r,j} - 0) - F(x_{r,j} - 0)| \right\} \longrightarrow 0\right) = 1.$$

Arătăm în continuare implicația

$$E \subset \left(\lim_{n \rightarrow \infty} \sup_{x \in \mathbb{R}} |\bar{F}_n(x) - F(x)| = 0 \right),$$

din care se va obține

$$1 = P(E) \leq P\left(\lim_{n \rightarrow \infty} \sup_{x \in \mathbb{R}} |\bar{F}_n(x) - F(x)| = 0\right) \leq 1,$$

deci

$$P\left(\lim_{n \rightarrow \infty} \sup_{x \in \mathbb{R}} |\bar{F}_n(x) - F(x)| = 0\right) = 1.$$

Pentru aceasta să considerăm $x \in \mathbb{R}$ oarecare. Pentru orice $r \in \mathbb{N}$ va exista $k \in \mathbb{N}$ astfel încât $x_{r,k} < x \leq x_{r,k+1}$. Având în vedere că \bar{F}_n și F sunt funcții nedescrescătoare, putem să scriem următoarele inegalități:

$$\bar{F}_n(x_{r,k}) \leq \bar{F}_n(x) \leq \bar{F}_n(x_{r,k+1} - 0)$$

și

$$F(x_{r,k}) \leq F(x) \leq F(x_{r,k+1} - 0),$$

de unde se obține

$$\bar{F}_n(x_{r,k}) - F(x_{r,k+1} - 0) \leq \bar{F}_n(x) - F(x) \leq \bar{F}_n(x_{r,k+1} - 0) - F(x_{r,k}).$$

Pe de altă parte, din modul de definire a numerelor $x_{r,k}$, $k = \overline{0, r}$, avem că

$$0 \leq F(x_{r,k+1} - 0) - F(x_{r,k}) \leq \frac{1}{r},$$

de unde

$$F(x_{r,k}) \geq F(x_{r,k+1} - 0) - \frac{1}{r} \quad \text{și} \quad F(x_{r,k+1} - 0) \leq F(x_{r,k}) + \frac{1}{r}.$$

Utilizate în dubla inegalitate de mai înainte, conduc la

$$\bar{F}_n(x_{r,k}) - F(x_{r,k}) - \frac{1}{r} \leq \bar{F}_n(x) - F(x) \leq \bar{F}_n(x_{r,k+1} - 0) - F(x_{r,k+1} - 0) + \frac{1}{r}.$$

Din prima inegalitate se obține

$$F(x) - \bar{F}_n(x) \leq F(x_{r,k}) - \bar{F}_n(x_{r,k}) + \frac{1}{r},$$

de unde

$$\begin{aligned} |\bar{F}_n(x) - F(x)| &\leq |\bar{F}_n(x_{r,k}) - F(x_{r,k})| + \frac{1}{r} \\ &\leq \max_{j=\overline{1,r}} \{|\bar{F}_n(x_{r,j}) - F(x_{r,j})|\} + \frac{1}{r}. \end{aligned}$$

Din a doua inegalitate se obține

$$\begin{aligned} |\bar{F}_n(x) - F(x)| &\leq |\bar{F}_n(x_{r,k+1} - 0) - F(x_{r,k+1} - 0)| + \frac{1}{r} \\ &\leq \max_{k=\overline{1,r}} \{|\bar{F}_n(x_{r,k} - 0) - F(x_{r,k} - 0)|\} + \frac{1}{r}. \end{aligned}$$

Folosind cele două inegalități obținute rezultă că $|\bar{F}_n(x) - F(x)| \leq M_{n,r}$, unde

$$M_{n,r} = \max_{k,j=\overline{1,r}} \{|\bar{F}_n(x_{r,k}) - F(x_{r,k})|, |\bar{F}_n(x_{r,j} - 0) - F(x_{r,j} - 0)|\} + \frac{1}{r},$$

de unde $|\bar{F}_n(x) - F(x)| \leq \sup_{r \in \mathbb{N}} M_{n,r}$, pentru orice $x \in \mathbb{R}$.

În final avem că $\sup_{x \in \mathbb{R}} |\bar{F}_n(x) - F(x)| \leq \sup_{r \in \mathbb{N}} M_{n,r}$, care conduce la implicația care trebuia demonstrată. \square

Teorema 4.2.38 (Kolmogorov). Fie caracteristica X care are funcția de repartiție teoretică F continuă și fie o selecție repetată de volum n relativă la caracteristica X

cu variabilele de selecție X_1, X_2, \dots, X_n și corespunzător funcția de repartiție de selecție \bar{F}_n , atunci

$$\lim_{n \rightarrow \infty} P(\sqrt{n}D_n \leq x) = K(x), \quad x > 0,$$

unde

$$D_n = \sup_{x \in \mathbb{R}} |\bar{F}_n(x) - F(x)|, \quad \text{iar} \quad K(x) = \sum_{k=-\infty}^{+\infty} (-1)^k e^{-2k^2 x^2}, \quad x > 0,$$

este funcția lui Kolmogorov.

Observația 4.2.39. Teorema lui Kolmogorov reprezintă comportarea asimptotică a distanței dintre funcția de repartiție de selecție și funcția de repartiție teoretică. Demonstrația poate fi găsită în [56].

Observația 4.2.40. Funcția lui Kolmogorov este tabelată, pentru anumite valori, în Anexa V.

Observația 4.2.41. Pentru calculul valorilor aproximative ale funcției lui Kolmogorov se pot utiliza următoarele formule de aproximare

$$K(x) \approx \begin{cases} 0, & \text{dacă } x \leq 0.27 \\ \frac{\sqrt{2\pi}}{x} \sum_{i=1}^3 e^{-(2i-1)^2 \pi^2 / (8x^2)}, & \text{dacă } 0.27 < x < 1, \\ 1 - 2 \sum_{i=1}^4 (-1)^{i-1} e^{-2i^2 x^2}, & \text{dacă } 1 \leq x < 3.1, \\ 1, & \text{dacă } x \geq 3.1. \end{cases}$$

4.2.5 Funcția `cdfplot`

Statistics toolbox dispune de funcția `cdfplot`, care reprezintă grafic funcția de repartiție de selecție. Apelul funcției se poate face prin:

```
cdfplot(x)
h=cdfplot(x)
[h,stats]=cdfplot(x)
```

Prima formă reprezintă grafic funcția de repartiție de selecție corespunzătoare datelor conținute de vectorul x , caracteristici cum ar fi culoarea, tipul liniei etc. sunt considerate implicit de sistemul Matlab. Pentru modificarea unor astfel de caracteristici, adică alegerea acestora de către utilizator, se recomandă a doua formă, prin care parametrul h poate fi folosit în comanda `set`.

Parametrul `stats` conține, după executarea funcției, o parte din parametrii statistici pentru `x`. Aceștia sunt: valoarea minimă (`stats.min`), valoarea maximă (`stats.max`), valoarea mediei de selecție (`stats.mean`), valoarea medianei (`stats.median`) și abaterea standard (`stats.std`).

Funcția 4.2.42. Funcția Matlab care urmează generează `n` numere aleatoare, care urmează o lege de probabilitate precizată, după care reprezintă grafic, pe aceeași figură, funcția de repartiție teoretică prin linie-punct, împreună cu funcția de repartiție de selecție.

```
function compar(n,lege)
% Functie de repartitie de selectie
% lege - legea de probabilitate
% n - volumul datelor
clf
switch lege
case 'unid'
    N=input('N=');
    X=unidrnd(N,n,1); [med,var]=unidstat(N);
    x=med-3*sqrt(var):0.01:med+3*sqrt(var);
    f=cdf(lege,x,N);
case 'bino'
    m=input('n='); p=input('p=');
    X=binornd(m,p,n,1); [med,var]=binostat(m,p);
    x=med-3*sqrt(var):0.01:med+3*sqrt(var);
    f=cdf(lege,x,m,p);
case 'hyge'
    M=input('M='); K=input('K='); m=input('n=');
    X=hygernd(M,K,m,n,1); [med,var]=hygestat(M,K,m);
    x=med-3*sqrt(var):0.01:med+3*sqrt(var);
    f=cdf(lege,x,M,K,m);
case 'poiss'
    la=input('lambda=');
    X=poissrnd(la,n,1); [med,var]=poisstat(la);
    x=med-3*sqrt(var):0.01:med+3*sqrt(var);
    f=cdf(lege,x,la);
case 'nbin'
    r=input('r='); p=input('p=');
    X=nbinrnd(r,p,n,1); [med,var]=nbinstat(r,p);
    x=med-3*sqrt(var):0.01:med+3*sqrt(var);
    f=cdf(lege,x,r,p);
case 'geo'
    p=input('p=');
    X=geornd(p,n,1); [med,var]=geostat(p);
    x=med-3*sqrt(var):0.01:med+3*sqrt(var);
    f=cdf(lege,x,p);
case 'unif'
    a=input('a:'); b=input('b(a<b):');
    X=unifrnd(a,b,n,1); [med,var]=unifstat(a,b);
```

```

x=med-3*sqrt(var):0.01:med+3*sqrt(var);
f=cdf(lege,x,a,b);
case 'norm'
mu=input('mu='); s=input('sigma=');
X=normrnd(mu,s,n,1); [med,var]=normstat(mu,s);
x=med-3*sqrt(var):0.01:med+3*sqrt(var);
f=cdf(lege,x,mu,s);
case 'logn'
mu=input('mu='); s=input('sigma=');
X=lognrnd(mu,s,n,1); [med,var]=lognstat(mu,s);
x=med-3*sqrt(var):0.01:med+3*sqrt(var);
f=cdf(lege,x,mu,s);
case 'gam'
a=input('a='); b=input('b=');
X=gamrnd(a,b,n,1); [med,var]=gamstat(a,b);
x=med-3*sqrt(var):0.01:med+3*sqrt(var);
f=cdf(lege,x,a,b);
case 'exp'
mu=input('mu=');
X=exprnd(mu,n,1); [med,var]=expstat(mu);
x=med-3*sqrt(var):0.01:med+3*sqrt(var);
f=cdf(lege,x,mu);
case 'beta'
a=input('a='); b=input('b=');
X=betarnd(a,b,n,1); [med,var]=betastat(a,b);
x=med-3*sqrt(var):0.01:med+3*sqrt(var);
f=cdf(lege,x,a,b);
case 'weib'
a=input('a='); b=input('b=');
X=weibrnd(a,b,n,1); [med,var]=weibstat(a,b);
x=med-3*sqrt(var):0.01:med+3*sqrt(var);
f=cdf(lege,x,a,b);
case 'rayl'
b=input('b=');
X=raylrnd(b,n,1); [med,var]=raylstat(b);
x=med-3*sqrt(var):0.01:med+3*sqrt(var);
f=cdf(lege,x,b);
case 't'
m=input('n=');
X=trnd(m,n,1); [med,var]=tstat(m);
x=med-3*sqrt(var):0.01:med+3*sqrt(var);
f=cdf(lege,x,m);
case 'chi2'
m=input('n=');
X=chi2rnd(m,n,1); [med,var]=chi2stat(m);
x=med-3*sqrt(var):0.01:med+3*sqrt(var);
f=cdf(lege,x,m);
case 'f'
m1=input('m='); m2=input('n=');

```

```

X=frnd(m1,m2,n,1); [med,var]=fstat(m1,m2);
x=med-3*sqrt(var):0.01:med+3*sqrt(var);
f=cdf(lege,x,m1,m2);
otherwise
    error('Lege necunoscuta')
end
plot(x,f,'k-.'), hold on, cdfplot(X), grid off
title('Functia de repartitie de selectie')

```

Apelul funcției compar se face prin

```
>>compar(n,'lege')
```

unde valorile pentru n și $lege$, fie că sunt precizate în acest apel, fie sunt precizate înainte. De exemplu, comanda

```
>>compar(25,'bino')
```

are ca efect apelul funcției pentru legea binomială, iar pe ecran se va cere introducerea parametrilor n și p , după care pe ecran va fi reprezentat graficul din Figura 4.6, în cazul în care $n=7$ și $p=0.4$.

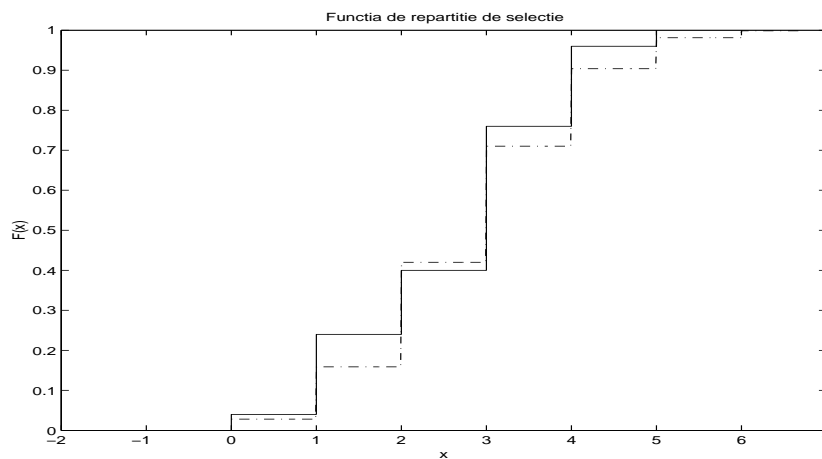
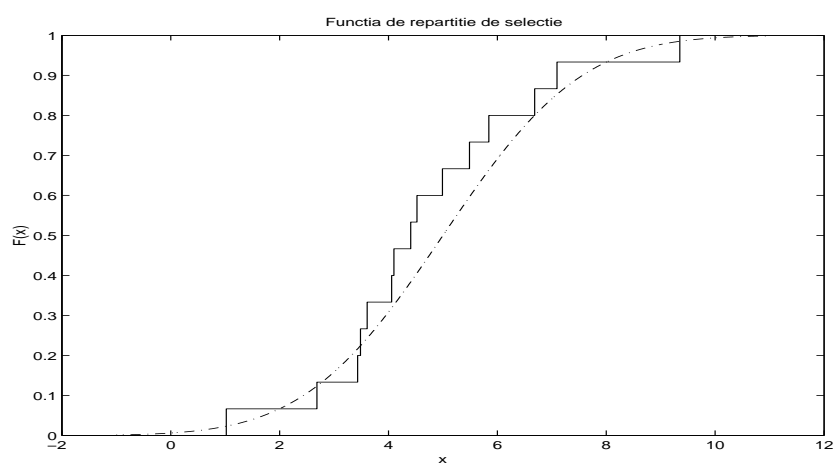


Figura 4.6: Legea $\mathcal{B}(7, 0.4)$

Comanda

```
>>compar(15,'norm')
```

are ca efect apelul funcției pentru legea normală, iar pe ecran se va cere introducerea parametrilor μ și σ , după care pe ecran va fi reprezentat graficul din Figura 4.7, în cazul în care $\mu=5$ și $\sigma=2$.

Figura 4.7: Legea $\mathcal{N}(5, 2)$

Capitolul 5

Teoria estimației

Relativ la colectivitatea \mathcal{C} este cercetată caracteristica X , care urmează legea de probabilitate dată prin funcția de probabilitate $f(x; \theta)$, ce reprezintă funcția de frecvență dacă X este de tip discret, respectiv densitatea de probabilitate în cazul continuu, iar θ este un parametru real necunoscut. Se consideră o selecție repetată de volum n având corespunzător variabilele de selecție X_1, X_2, \dots, X_n .

5.1 Funcție de verosimilitate

Definiția 5.1.1. Numim funcție de verosimilitate, *funcția de selecție*

$$g(X_1, X_2, \dots, X_n; \theta) = \prod_{k=1}^n f(X_k; \theta).$$

Observația 5.1.2. Dacă se consideră funcția de n variabile reale

$$g(x_1, x_2, \dots, x_n; \theta) = \prod_{k=1}^n f(x_k; \theta),$$

aceasta reprezintă funcția frecvență (în cazul discret), respectiv densitatea de probabilitate (în cazul continuu) a vectorului aleator (X_1, X_2, \dots, X_n) .

Definiția 5.1.3. Statistica $S = S(X_1, X_2, \dots, X_n)$ este statistică suficientă (exhaustivă) pentru parametrul θ , dacă există funcția măsurabilă $\varphi: \mathbb{R}^n \rightarrow \mathbb{R}$ nenegativă și funcția măsurabilă $h_\theta: \mathbb{R} \rightarrow \mathbb{R}$ nenegativă, astfel încât

$$g(x_1, x_2, \dots, x_n; \theta) = \varphi(x_1, x_2, \dots, x_n) h_\theta(s) = \varphi(x_1, x_2, \dots, x_n) h(s; \theta),$$

unde $s = S(x_1, x_2, \dots, x_n)$.

Observația 5.1.4. Echivalentă cu condiția precedentă pentru suficiența statisticii S este condiția ca funcția de frecvență (în cazul discret), respectiv densitatea de probabilitate (în cazul continuu) $f(x_1, \dots, x_n; \theta | s)$ a vectorului aleator (X_1, \dots, X_n) condiționată de evenimentul $S(X_1, \dots, X_n) = s$, să nu depindă de parametrul θ .

Exemplul 5.1.5. Fie caracteristica X ce urmează legea lui Poisson cu parametrul $\lambda > 0$ necunoscut, adică are funcția de frecvență

$$f(x; \lambda) = \frac{\lambda^x}{x!} e^{-\lambda}, \quad x = 0, 1, 2, \dots$$

Considerăm statistica $S = \sum_{k=1}^n X_k$. Deoarece variabilele aleatoare X_1, X_2, \dots, X_n sunt independente și identic repartizate cu X avem că S urmează legea lui Poisson de parametru $n\lambda$, adică

$$P(S = s) = \frac{(n\lambda)^s}{s!} e^{-n\lambda}, \quad s = 0, 1, 2, \dots$$

Pe de o parte avem:

$$\begin{aligned} g(x_1, x_2, \dots, x_n; \lambda) &= \prod_{k=1}^n f(x_k; \theta) = \prod_{k=1}^n \left(\frac{\lambda^{x_k}}{x_k!} e^{-\lambda} \right) = \frac{\lambda^{\sum_{k=1}^n x_k}}{\prod_{k=1}^n x_k!} e^{-n\lambda} \\ &= \frac{\lambda^s}{x_1! x_2! \dots x_n!} e^{-n\lambda}. \end{aligned}$$

Se poate considera

$$\varphi(x_1, x_2, \dots, x_n) = \frac{1}{x_1! x_2! \dots x_n!}, \quad h(s; \lambda) = \lambda^s e^{-n\lambda}.$$

Având în vedere că putem scrie și

$$g(x_1, x_2, \dots, x_n; \lambda) = \frac{(n\lambda)^s}{s!} e^{-n\lambda} \cdot \frac{s!}{n^s x_1! x_2! \dots x_n!},$$

se poate considera de asemenea

$$\varphi(x_1, x_2, \dots, x_n) = \frac{s!}{n^s x_1! x_2! \dots x_n!}, \quad h(s; \lambda) = \frac{(n\lambda)^s}{s!} e^{-n\lambda}.$$

Prin urmare, unicitatea acestei scrieri nu se impune.

Pe de altă parte, având în vedere că S urmează legea $\mathcal{P}(n\lambda)$, avem succesiv

$$f(x_1, x_2, \dots, x_n; \lambda | s) = P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n | S = s)$$

$$\begin{aligned}
&= \frac{P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n, S = s)}{P(S = s)} \\
&= \frac{P(X_1 = x_1, \dots, X_{n-1} = x_{n-1}, X_n = s - \sum_{k=1}^{n-1} x_k)}{P(S = s)} \\
&= \frac{P(X_n = s - \sum_{k=1}^{n-1} x_k) \prod_{i=1}^{n-1} P(X_i = x_i)}{P(S = s)} \\
&= \frac{\frac{\lambda^{s - \sum_{k=1}^{n-1} x_k}}{(s - \sum_{k=1}^{n-1} x_k)!} e^{-\lambda} \prod_{i=1}^{n-1} \left(\frac{\lambda^{x_i}}{x_i!} e^{-\lambda} \right)}{\frac{(n\lambda)^s}{s!} e^{-n\lambda}} = \frac{s!}{n^s x_1! x_2! \dots x_n!}.
\end{aligned}$$

Așadar $f(x_1, x_2, \dots, x_n; \lambda | s)$ nu depinde de parametrul λ , decât prin intermediul valorii statisticii S .

Exemplul 5.1.6 (Familia exponențială). Se consideră caracteristica X , care are funcția de probabilitate de forma

$$f(x; \theta) = \exp\{a(x)\alpha(\theta) + b(x) + \beta(\theta)\}.$$

Vom arăta că, în acest caz, statistica

$$S = S(X_1, X_2, \dots, X_n) = \sum_{k=1}^n a(X_k)$$

este suficientă pentru parametrul θ . Într-adevăr avem că

$$\begin{aligned}
g(x_1, x_2, \dots, x_n; \theta) &= \prod_{k=1}^n f(x_k; \theta) \\
&= \exp\left\{\alpha(\theta) \sum_{k=1}^n a(x_k) + n\beta(\theta) + \sum_{k=1}^n b(x_k)\right\} \\
&= \exp\left\{\alpha(\theta) \sum_{k=1}^n a(x_k) + n\beta(\theta)\right\} \exp\left\{\sum_{k=1}^n b(x_k)\right\}.
\end{aligned}$$

Dacă se consideră

$$\varphi(x_1, x_2, \dots, x_n) = \exp\left\{\sum_{k=1}^n b(x_k)\right\} \quad \text{și} \quad h(s; \theta) = \exp\left\{\alpha(\theta)s + n\beta(\theta)\right\},$$

atunci

$$g(x_1, x_2, \dots, x_n; \theta) = \varphi(x_1, x_2, \dots, x_n) h(s; \theta),$$

deci condiția suficienței este îndeplinită.

Definiția 5.1.7. Numim cantitate de informație (Fisher) a unei selecții de volum n relativă la parametrul $\theta \in \mathbb{R}$ necunoscut, valoarea medie

$$I_n(\theta) = E \left[\left(\frac{\partial \ln g(X_1, X_2, \dots, X_n; \theta)}{\partial \theta} \right)^2 \right],$$

când funcția de verosimilitate este derivabilă în raport cu θ .

Observația 5.1.8. Dacă parametrul θ este p -dimensional, matricea

$$\mathbf{F} = \left(\text{Cov} \left(\frac{\partial \ln g(X_1, \dots, X_n; \theta)}{\partial \theta_i}, \frac{\partial \ln g(X_1, \dots, X_n; \theta)}{\partial \theta_j} \right) \right)_{i,j=\overline{1,p}}$$

se numește matricea informației lui Fisher.

Teorema 5.1.9. Dacă domeniul valorilor caracteristicii X nu depinde de parametrul θ , iar funcția de verosimilitate este derivabilă de două ori în raport cu θ , atunci

$$I_n(\theta) = -E \left(\frac{\partial^2 \ln g(X_1, X_2, \dots, X_n; \theta)}{\partial \theta^2} \right).$$

Demonstrație. Se pornește de la relația cunoscută

$$\int \dots \int_{\mathbb{R}^n} g(x_1, x_2, \dots, x_n; \theta) dx_1 dx_2 \dots dx_n = 1,$$

pe care o satisface densitatea de probabilitate. Se derivează această relație în raport cu θ și se ține seama de faptul că

$$\frac{\partial g(x_1, x_2, \dots, x_n; \theta)}{\partial \theta} = g(x_1, x_2, \dots, x_n) \frac{\partial \ln g(x_1, x_2, \dots, x_n; \theta)}{\partial \theta},$$

obținându-se

$$\int \dots \int_{\mathbb{R}^n} \frac{\partial \ln g(x_1, x_2, \dots, x_n; \theta)}{\partial \theta} g(x_1, x_2, \dots, x_n; \theta) dx_1 dx_2 \dots dx_n = 0.$$

Derivând încă odată în raport cu θ rezultă

$$\begin{aligned} & \int \dots \int_{\mathbb{R}^n} \frac{\partial^2 \ln g(x_1, x_2, \dots, x_n; \theta)}{\partial \theta^2} g(x_1, x_2, \dots, x_n; \theta) dx_1 dx_2 \dots dx_n \\ & + \int \dots \int_{\mathbb{R}^n} \frac{\partial \ln g(x_1, x_2, \dots, x_n; \theta)}{\partial \theta} \frac{\partial g(x_1, x_2, \dots, x_n; \theta)}{\partial \theta} dx_1 dx_2 \dots dx_n = 0, \end{aligned}$$

sau, având în vedere relația de mai înainte, rezultă că

$$\int \dots \int_{\mathbb{R}^n} \frac{\partial^2 \ln g(x_1, x_2, \dots, x_n; \theta)}{\partial \theta^2} g(x_1, x_2, \dots, x_n; \theta) dx_1 dx_2 \dots dx_n \\ + \int \dots \int_{\mathbb{R}^n} \left[\frac{\partial \ln g(x_1, x_2, \dots, x_n; \theta)}{\partial \theta} \right]^2 g(x_1, x_2, \dots, x_n; \theta) dx_1 dx_2 \dots dx_n = 0.$$

Așadar am obținut că

$$E \left[\frac{\partial^2 \ln g(X_1, X_2, \dots, X_n; \theta)}{\partial \theta^2} \right] + E \left[\left(\frac{\partial \ln g(X_1, X_2, \dots, X_n; \theta)}{\partial \theta} \right)^2 \right] = 0,$$

de unde avem relația dorită. \square

Observația 5.1.10. În demonstrația teoremei precedente am considerat cazul unei caracteristici X de tip continuu. În mod analog se procedează și în cazul discret, integrala multiplă este înlocuită cu o sumă multiplă.

Corolarul 5.1.11. Când domeniul de definiție al caracteristicii X nu depinde de parametrul θ , atunci $I_n(\theta) = nI_1(\theta)$.

Demonstrație. Deoarece selecția este repetată, avem că

$$\frac{\partial^2 \ln g(x_1, x_2, \dots, x_n; \theta)}{\partial \theta^2} = \sum_{k=1}^n \frac{\partial^2 \ln f(x_k; \theta)}{\partial \theta^2}.$$

Folosind Teorema 5.1.9, se obține

$$I_n(\theta) = - \sum_{k=1}^n E \left[\frac{\partial^2 \ln f(X_k; \theta)}{\partial \theta^2} \right] = \sum_{k=1}^n I_1(\theta) = nI_1(\theta),$$

deoarece

$$I_1(\theta) = -E \left[\frac{\partial^2 \ln f(X; \theta)}{\partial \theta^2} \right].$$

\square

Observația 5.1.12. Din demonstrația Teoremei 5.1.9 avem de asemenea că

$$I_n(\theta) = Var \left[\frac{\partial \ln g(X_1, X_2, \dots, X_n; \theta)}{\partial \theta} \right],$$

deoarece

$$E \left[\frac{\partial \ln g(X_1, X_2, \dots, X_n; \theta)}{\partial \theta} \right] = 0.$$

Exemplul 5.1.13. Să considerăm caracteristica X , ce urmează legea normală $\mathcal{N}(\mu, \sigma)$, unde $\mu \in \mathbb{R}$ este necunoscut, iar $\sigma > 0$ este cunoscut.

Deoarece

$$f(x; \mu) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \quad x \in \mathbb{R},$$

avem că

$$I_1(\mu) = E\left[\left(\frac{\partial \ln f(X; \mu)}{\partial \mu}\right)^2\right] = E\left[\frac{(X - \mu)^2}{\sigma^4}\right] = \frac{1}{\sigma^2}.$$

Prin urmare, cantitatea de informație conținută (adusă) de o observație, relativ la parametrul μ , este cu atât mai mare cu cât dispersia este mai mică.

Teorema 5.1.14. Fie caracteristica X , cu funcția de probabilitate $f(x; \theta)$ derivabilă de două ori în raport cu θ și statistica $S = S(X_1, X_2, \dots, X_n)$ relativă la caracteristica X , cu funcția de probabilitate $h(s; \theta)$, atunci cantitatea de informație

$$I_S(\theta) = -E\left[\frac{\partial^2 \ln h(S; \theta)}{\partial \theta^2}\right],$$

relativă la parametrul θ conținută în statistica S nu depășește cantitatea de informație $I_n(\theta)$ conținută de selecția considerată, adică $I_S(\theta) \leq I_n(\theta)$.

Demonstrație. Dacă $h(s; \theta)$ este funcția de frecvență (în cazul discret), respectiv densitatea de probabilitate (în cazul continuu) pentru statistica S , atunci

$$g(x_1, x_2, \dots, x_n; \theta) = h(s; \theta) f(x_1, x_2, \dots, x_n; \theta | s),$$

unde $f(x_1, x_2, \dots, x_n; \theta | s)$ este funcția de frecvență (în cazul discret), respectiv densitatea de probabilitate (în cazul continuu) pentru vectorul aleator (X_1, \dots, X_n) condiționată de evenimentul $(S = s)$.

Așadar avem că

$$(5.1.1) \quad \frac{\partial^2 \ln g(x_1, x_2, \dots, x_n; \theta)}{\partial \theta^2} = \frac{\partial^2 \ln h(s; \theta)}{\partial \theta^2} + \frac{\partial^2 \ln f(x_1, x_2, \dots, x_n; \theta | s)}{\partial \theta^2},$$

deci

$$I_n(\theta) = I_S(\theta) - E\left[\frac{\partial^2 \ln f(X_1, X_2, \dots, X_n; \theta | s)}{\partial \theta^2}\right] \geq I_S(\theta),$$

ceea ce trebuia arătat. □

Exemplul 5.1.15. Se consideră caracteristica X ce urmează legea normală $\mathcal{N}(\mu, \sigma)$ și o selecție repetată de volum n . Vom compara informația dispersiei de selecție

$$\bar{\sigma}^2 = \frac{1}{n-1} \sum_{k=1}^n (X_k - \bar{X})^2,$$

relativă la dispersia teoretică $Var(X) = \sigma^2$, cu informația selecției relativă la același parametru.

Calculăm prima dată informația selecției de volum n relativă la parametrul notat cu $\theta = \sigma^2 = Var(X)$. Având în vedere că densitatea de probabilitate a caracteristicii X este

$$f(x; \theta) = \frac{1}{\sqrt{2\pi\theta}} e^{-\frac{(x-\mu)^2}{2\theta}}, \quad x \in \mathbb{R},$$

se obține

$$\frac{\partial \ln f(x; \theta)}{\partial \theta} = -\frac{1}{2\theta} + \frac{1}{2\theta^2} (x - \mu)^2 \quad \text{și} \quad \frac{\partial^2 \ln f(x; \theta)}{\partial \theta^2} = \frac{1}{2\theta^2} - \frac{1}{\theta^3} (x - \mu)^2.$$

Rezultă astfel că

$$I_1(\sigma^2) = I_1(\theta) = -\frac{1}{2\sigma^4} + \frac{1}{\sigma^4} E\left[\left(\frac{X - \mu}{\sigma}\right)^2\right] = \frac{1}{2\sigma^4},$$

deci

$$I_n(\sigma^2) = \frac{n}{2\sigma^4}.$$

Pe de altă parte, se știe că statistica

$$\chi^2 = \frac{n-1}{\sigma^2} \bar{\sigma}^2 = \frac{1}{\sigma^2} \sum_{k=1}^n (X_k - \bar{X})^2$$

urmează legea χ^2 cu $n-1$ grade de libertate. Având în vedere legătura dintre statisticile χ^2 și $\bar{\sigma}^2$ rezultă că densitatea de probabilitate pentru $\bar{\sigma}^2$ este dată prin

$$h(s; \theta) = \frac{1}{2^{\frac{n-1}{2}} \Gamma\left(\frac{n-1}{2}\right)} \left(\frac{n-1}{\theta}\right)^{\frac{n-1}{2}} s^{\frac{n-1}{2}-1} e^{-\frac{(n-1)s}{2\theta}}, \quad s > 0.$$

Astfel se obține că

$$\ln h(s; \theta) = C + \frac{n-1}{2} \ln \frac{n-1}{\theta} + \frac{n-3}{2} \ln s - \frac{(n-1)s}{2\theta},$$

de unde

$$\frac{\partial \ln h(s; \theta)}{\partial \theta} = -\frac{n-1}{2\theta} + \frac{(n-1)s}{2\theta^2} \quad \text{și} \quad \frac{\partial^2 \ln h(s; \theta)}{\partial \theta^2} = \frac{n-1}{2\theta^2} - \frac{(n-1)s}{\theta^3}.$$

Calculăm acum informația dispersiei de selecție relativă la σ^2 . Pentru aceasta ținem seama de faptul că

$$E(\bar{\sigma}^2) = \frac{\sigma^2}{n-1} E(\chi^2) = \frac{\sigma^2}{n-1} (n-1) = \sigma^2.$$

Astfel obținem:

$$I_{\bar{\sigma}^2}(\sigma^2) = -\frac{n-1}{2\sigma^4} + \frac{n-1}{\sigma^6} E(\bar{\sigma}^2) = -\frac{n-1}{2\sigma^4} + \frac{n-1}{\sigma^4},$$

deci

$$I_{\bar{\sigma}^2}(\sigma^2) = \frac{n-1}{2\sigma^4} < \frac{n}{2\sigma^4} = I_n(\sigma^2).$$

Observația 5.1.16. Deoarece, în cazul unei statistici S suficiente, avem că funcția $f(x_1, x_2, \dots, x_n; \theta | s)$ nu depinde de θ rezultă că $I_n(\theta) = I_S(\theta)$ în acest caz. Acest lucru se obține imediat, dacă se are în vedere relația (5.1.1).

5.2 Funcții de estimăție

Definiția 5.2.1. Fie caracteristica X cu funcția de probabilitate $f(x; \theta)$, parametrul $\theta \in A$ necunoscut și o selecție repetată de volum n . Numim funcție de estimăție (estimator) pentru parametrul θ , funcția de selecție

$$\hat{\Theta} = \hat{\theta}(X_1, X_2, \dots, X_n),$$

care ia valori în domeniul A , iar valoarea numerică $\hat{\theta} = \hat{\theta}(x_1, x_2, \dots, x_n)$ se numește estimăția lui θ .

Definiția 5.2.2. Estimatorul $\hat{\Theta} = \hat{\theta}(X_1, X_2, \dots, X_n)$ este estimator (funcție de estimăție) nedeplasat pentru parametrul necunoscut θ , dacă $E(\hat{\Theta}) = \theta$, iar valoarea numerică $\hat{\theta} = \hat{\theta}(x_1, x_2, \dots, x_n)$ se numește estimăție nedeplasată pentru parametrul θ .

Definiția 5.2.3. Spunem că estimatorul $\hat{\Theta} = \hat{\theta}(X_1, X_2, \dots, X_n)$ este estimator consistent, pentru parametrul necunoscut θ , dacă $\hat{\Theta} \xrightarrow{P} \theta$, adică

$$\lim_{n \rightarrow \infty} P(|\hat{\Theta} - \theta| < \varepsilon) = 1,$$

pentru orice $\varepsilon > 0$, iar valoarea numerică $\hat{\theta} = \hat{\theta}(x_1, x_2, \dots, x_n)$ se numește estimăție consistentă pentru parametrul θ .

5.2.1 Funcții de estimare absolut corecte

Definiția 5.2.4. Numim funcție de estimare (estimator) absolut corectă pentru parametrul θ , funcția de selecție $\hat{\Theta} = \hat{\theta}(X_1, X_2, \dots, X_n)$, care satisface condițiile

$$(i) \quad E(\hat{\Theta}) = \theta,$$

$$(ii) \quad \lim_{n \rightarrow \infty} Var(\hat{\Theta}) = 0,$$

iar valoarea numerică $\hat{\theta} = \hat{\theta}(x_1, x_2, \dots, x_n)$ se numește estimare absolut corectă pentru parametrul θ .

Proprietatea 5.2.5. Un estimator absolut corect este un estimator consistent.

Demonstrație. Fie estimatorul $\hat{\Theta} = \hat{\theta}(X_1, X_2, \dots, X_n)$, un estimator absolut corect pentru parametrul θ . Din inegalitatea lui Cebîșev avem

$$1 \geq P(|\hat{\Theta} - \theta| < \varepsilon) \geq 1 - \frac{Var(\hat{\Theta})}{\varepsilon^2},$$

pentru orice $\varepsilon > 0$. Făcând pe $n \rightarrow \infty$, se obține

$$\lim_{n \rightarrow \infty} P(|\hat{\Theta} - \theta| < \varepsilon) = 1,$$

pentru orice $\varepsilon > 0$, ceea ce trebuie demonstrat. \square

Exemplul 5.2.6. Fie caracteristica X , ce urmează legea normală $\mathcal{N}(\mu, \sigma)$, unde parametrul $\mu \in \mathbb{R}$ este cunoscut, iar $\sigma > 0$ este necunoscut. Considerând o selecție repetată de volum n , vom determina constanta κ_n astfel încât statistica

$$\bar{S} = \kappa_n \sum_{k=1}^n |X_k - \mu|$$

să fie estimator absolut corect pentru abaterea standard σ .

Având în vedere că variabilele de selecție X_1, X_2, \dots, X_n sunt variabile aleatoare identic repartizate cu X , putem scrie succesiv

$$\begin{aligned} (5.2.1) \quad E(\bar{S}) &= \kappa_n \sum_{k=1}^n E(|X_k - \mu|) = \kappa_n \sum_{k=1}^n E(|X - \mu|) \\ &= n\kappa_n E(|X - \mu|). \end{aligned}$$

Ținem seama de faptul că X urmează legea normală $\mathcal{N}(\mu, \sigma)$ și obținem

$$E(|X - \mu|) = \int_{-\infty}^{+\infty} |x - \mu| f(x; \sigma) dx = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{+\infty} |x - \mu| e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx.$$

Facem schimbarea de variabilă $x - \mu = t$ și ținem seama de faptul că funcția obținută este pară, iar intervalul de integrare este simetric față de origine:

$$\begin{aligned} E(|X - \mu|) &= \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{+\infty} |t| e^{-\frac{t^2}{2\sigma^2}} dt = \frac{1}{\sigma} \sqrt{\frac{2}{\pi}} \int_0^{+\infty} t e^{-\frac{t^2}{2\sigma^2}} dt \\ &= \frac{1}{\sigma} \sqrt{\frac{2}{\pi}} \left(-\sigma^2 e^{-\frac{t^2}{2\sigma^2}} \right) \Big|_0^{+\infty} = \sigma \sqrt{\frac{2}{\pi}}. \end{aligned}$$

Din condiția de nedeplasare și din formula (5.2.1) avem

$$\sigma = E(\bar{S}) = n\kappa_n \sigma \sqrt{\frac{2}{\pi}},$$

adică $\kappa_n = \frac{1}{n} \sqrt{\frac{\pi}{2}}$.

Până aici am arătat că statistica

$$\bar{S} = \frac{1}{n} \sqrt{\frac{\pi}{2}} \sum_{k=1}^n |X_k - \mu|$$

este un estimator nedeplasat pentru parametrul σ .

Deoarece selecția considerată este repetată rezultă că variabilele de selecție sunt independente, prin urmare se poate scrie

$$Var(\bar{S}) = \frac{\pi}{2n^2} \sum_{k=1}^n Var(|X_k - \mu|) = \frac{\pi}{2n^2} \sum_{k=1}^n Var(|X - \mu|) = \frac{\pi}{2n} Var(|X - \mu|).$$

Astfel este îndeplinită și condiția (ii) din Definiția 5.2.4, anume

$$Var(\bar{S}) = \frac{\pi}{2n} Var(|X - \mu|) \rightarrow 0, \quad \text{când } n \rightarrow \infty.$$

Remarcăm că

$$Var(|X - \mu|) = E[(X - \mu)^2] - [E(|X - \mu|)]^2 = \sigma^2 - \sigma^2 \frac{2}{\pi} = \frac{\pi - 2}{\pi} \sigma^2.$$

Proprietatea 5.2.7. Fie caracteristica X pentru care există momentul teoretic de ordinul $2k$, $\nu_{2k} = E(X^{2k})$, și fie o selecție repetată de volum n , atunci momentul de selecție de ordin k

$$\bar{\nu}_k = \frac{1}{n} \sum_{i=1}^n X_i^k$$

este funcție de estimăție absolut corectă pentru parametrul ν_k .

Demonstrație. Din Proprietatea 4.2.11 avem că $E(\bar{\nu}_k) = \nu_k$ și de asemenea

$$\lim_{n \rightarrow \infty} \text{Var}(\bar{\nu}_k) = \lim_{n \rightarrow \infty} \frac{\nu_{2k} - \nu_k^2}{n} = 0.$$

Condițiile Definiției 5.2.4 sunt satisfăcute, deci proprietatea este demonstrată. \square

Observația 5.2.8. Media de selecție \bar{X} ($= \bar{\nu}_1$) este funcție de estimare absolut corectă pentru media teoretică $E(X)$ ($= \nu_1$).

5.2.2 Funcții de estimare corecte

Definiția 5.2.9. Numim funcție de estimare (estimator) corectă, pentru parametrul necunoscut θ , funcția de selecție $\hat{\Theta} = \hat{\theta}(X_1, X_2, \dots, X_n)$, care satisface condițiile

- (i) $\lim_{n \rightarrow \infty} E(\hat{\Theta}) = \theta$,
- (ii) $\lim_{n \rightarrow \infty} \text{Var}(\hat{\Theta}) = 0$,

iar valoarea numerică $\hat{\theta} = \hat{\theta}(x_1, x_2, \dots, x_n)$ se numește estimare corectă pentru parametrul θ .

Proprietatea 5.2.10. Un estimator corect este un estimator consistent.

Demonstrație. Fie estimatorul $\hat{\Theta} = \hat{\theta}(X_1, X_2, \dots, X_n)$ corect pentru parametrul θ , atunci din condițiile (i) și (ii) avem că pentru orice $\varepsilon > 0$ și $\delta > 0$ există numărul natural $N = N(\varepsilon, \delta)$ astfel încât

$$\left| E(\hat{\Theta}) - \theta \right| < \frac{\varepsilon}{2} \quad \text{și} \quad \text{Var}(\hat{\Theta}) < \frac{\varepsilon^2 \delta}{4},$$

pentru $n > N$. Așadar, putem scrie

$$\left| \hat{\Theta} - \theta \right| \leq \left| \hat{\Theta} - E(\hat{\Theta}) \right| + \left| E(\hat{\Theta}) - \theta \right| < \left| \hat{\Theta} - E(\hat{\Theta}) \right| + \frac{\varepsilon}{2},$$

pentru $n > N$, de unde dacă $\left| \hat{\Theta} - E(\hat{\Theta}) \right| < \frac{\varepsilon}{2}$, vom avea că $\left| \hat{\Theta} - \theta \right| < \varepsilon$, pentru $n > N$. Prin urmare avem

$$\left(\left| \hat{\Theta} - E(\hat{\Theta}) \right| < \frac{\varepsilon}{2} \right) \subset \left(\left| \hat{\Theta} - \theta \right| < \varepsilon \right), \quad n > N,$$

care conduce la inegalitatea

$$P\left(\left| \hat{\Theta} - E(\hat{\Theta}) \right| < \frac{\varepsilon}{2}\right) \leq P\left(\left| \hat{\Theta} - \theta \right| < \varepsilon\right), \quad n > N.$$

Pe de altă parte, folosind inegalitatea lui Cebîșev,

$$P\left(\left|\hat{\theta} - E(\hat{\theta})\right| < \frac{\varepsilon}{2}\right) \geq 1 - \frac{4Var(\hat{\theta})}{\varepsilon^2}.$$

Deoarece $Var(\hat{\theta}) < \frac{\varepsilon^2 \delta}{4}$, pentru $n > N$, rezultă că

$$P\left(\left|\hat{\theta} - E(\hat{\theta})\right| < \frac{\varepsilon}{2}\right) \geq 1 - \delta, \quad n > N.$$

Prin urmare se ajunge la

$$P\left(\left|\hat{\theta} - \theta\right| < \varepsilon\right) \geq P\left(\left|\hat{\theta} - E(\hat{\theta})\right| < \frac{\varepsilon}{2}\right) \geq 1 - \delta, \quad n > N.$$

de unde $\hat{\theta} \xrightarrow{p} \theta$, ceea ce trebuie arătat. \square

Proprietatea 5.2.11. Fie caracteristica X pentru care există momentul teoretic de ordin $2k$, $\nu_{2k} = E(X^{2k})$, și fie o selecție repetată de volum n , atunci momentul centrat de selecție de ordin k

$$\bar{\mu}_k = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^k$$

este funcție de estimăție corectă pentru momentul centrat teoretic de ordin k , adică pentru $\mu_k = E[(X - E(X))^k]$.

Demonstrație. Conform Observației 4.2.18 avem

$$\lim_{n \rightarrow \infty} E(\bar{\mu}_k) = \lim_{n \rightarrow \infty} \left[\mu_k + O\left(\frac{1}{n}\right) \right] = \mu_k.$$

De asemenea, avem

$$\lim_{n \rightarrow \infty} Var(\bar{\mu}_k) = \lim_{n \rightarrow \infty} \left[\frac{\mu_{2k} - 2k\mu_{k-1}\mu_{k+1} - \mu_k^2 + k^2\mu_k\mu_{k-1}^2}{n} + O\left(\frac{1}{n^2}\right) \right] = 0.$$

Așadar, condițiile Definiției 5.2.9 sunt satisfăcute. \square

Observația 5.2.12. Momentul centrat de selecție de ordinul 2, $\bar{\mu}_2$, este funcție de estimăție corectă pentru dispersia teoretică $Var(X)$ ($= \mu_2$).

Având în vedere Observația 4.2.17 (formula (4.2.2)) avem că dispersia de selecție

$$\bar{\sigma}^2 = \frac{1}{n-1} \sum_{k=1}^n (X_k - \bar{X})^2$$

este funcție de estimăție absolut corectă pentru dispersia teoretică $Var(X)$ ($= \mu_2$).

5.2.3 Funcții de estimare eficiente

Teorema 5.2.13 (Inegalitatea Rao-Cramer). *Se consideră caracteristica X având funcția de probabilitate $f(x; \theta)$, $\theta \in (a, b)$, pentru care există derivata parțială de ordinul întâi în raport cu θ și fie o statistică $\hat{\Theta} = \hat{\theta}(X_1, \dots, X_n)$ absolut corectă pentru parametrul θ , atunci*

$$\text{Var}(\hat{\Theta}) \geq \frac{1}{I_n(\theta)}.$$

Demonstrație. Deoarece estimatorul $\hat{\Theta}$ este nedeplasat, avem că $E(\hat{\Theta}) = \theta$, adică

$$\int \dots \int_{\mathbb{R}^n} \hat{\theta}(x_1, \dots, x_n) g(x_1, \dots, x_n; \theta) dx_1 \dots dx_n = \theta,$$

unde $g(x_1, \dots, x_n; \theta) = \prod_{k=1}^n f(x_k; \theta)$ este funcția de verosimilitate.

Prin derivarea în raport cu θ a acestei relații se obține că

$$\int \dots \int_{\mathbb{R}^n} \hat{\theta}(x_1, \dots, x_n) \sum_{k=1}^n f(x_1; \theta) \dots \frac{\partial f(x_k; \theta)}{\partial \theta} \dots f(x_n; \theta) dx_1 \dots dx_n = 1$$

sau

$$\int \dots \int_{\mathbb{R}^n} \hat{\theta}(x_1, \dots, x_n) \left[\sum_{k=1}^n \frac{\partial \ln f(x_k; \theta)}{\partial \theta} \right] \left[\prod_{i=1}^n f(x_i; \theta) \right] dx_1 \dots dx_n = 1.$$

Pe de altă parte, deoarece

$$\int_{\mathbb{R}} \frac{\partial \ln f(x; \theta)}{\partial \theta} f(x; \theta) dx = 0,$$

avem

$$\theta \sum_{k=1}^n \int \dots \int_{\mathbb{R}^n} \frac{\partial \ln f(x_k; \theta)}{\partial \theta} \left[\prod_{i=1}^n f(x_i; \theta) \right] dx_1 \dots dx_n = 0,$$

care scăzută din egalitatea obținută înainte conduce la

$$\int \dots \int_{\mathbb{R}^n} \left[\hat{\theta}(x_1, \dots, x_n) - \theta \right] \left[\sum_{k=1}^n \frac{\partial \ln f(x_k; \theta)}{\partial \theta} \right] g(x_1, \dots, x_n; \theta) dx_1 \dots dx_n = 1,$$

adică

$$E \left[\left(\hat{\theta}(X_1, \dots, X_n) - \theta \right) \left(\sum_{k=1}^n \frac{\partial \ln f(X_k; \theta)}{\partial \theta} \right) \right] = 1.$$

Se aplică acum inegalitatea lui Schwarz și se obține

$$\begin{aligned} 1 &= \left[E \left[\left(\hat{\theta}(X_1, \dots, X_n) - \theta \right) \left(\sum_{k=1}^n \frac{\partial \ln f(X_k; \theta)}{\partial \theta} \right) \right] \right]^2 \\ &\leq E \left[\left(\hat{\theta}(X_1, \dots, X_n) - \theta \right)^2 \right] E \left[\left(\sum_{k=1}^n \frac{\partial \ln f(X_k; \theta)}{\partial \theta} \right)^2 \right] \\ &= \text{Var}(\hat{\theta}) \text{Var} \left[\sum_{k=1}^n \frac{\partial \ln f(X_k; \theta)}{\partial \theta} \right] = \text{Var}(\hat{\theta}) \cdot n \text{Var} \left[\frac{\partial \ln f(X; \theta)}{\partial \theta} \right], \end{aligned}$$

adică

$$\text{Var}(\hat{\theta}) \geq \frac{1}{n \text{Var} \left[\frac{\partial \ln f(X; \theta)}{\partial \theta} \right]}.$$

Dar avem că

$$I_n(\theta) = n I_1(\theta) = n E \left[\left(\frac{\partial \ln f(X; \theta)}{\partial \theta} \right)^2 \right] = n \text{Var} \left[\frac{\partial \ln f(X; \theta)}{\partial \theta} \right],$$

care înlocuită în inegalitatea dinainte conduce la inegalitatea Rao–Cramer. \square

Definiția 5.2.14. Se numește eficiență a unei funcții de estimare absolut corecte, $\hat{\theta} = \hat{\theta}(X_1, X_2, \dots, X_n)$, pentru parametrul θ , raportul

$$e(\hat{\theta}) = \frac{I_n^{-1}(\theta)}{\text{Var}(\hat{\theta})}.$$

Definiția 5.2.15. Spunem că funcția de estimare absolut corectă pentru parametrul θ , $\hat{\theta} = \hat{\theta}(X_1, X_2, \dots, X_n)$, este eficientă dacă inegalitatea Rao–Cramer este verificată prin egalitate, adică $e(\hat{\theta}) = 1$.

Exemplul 5.2.16. Se consideră caracteristica X ce urmează legea normală $\mathcal{N}(\mu, \sigma)$ și o selecție repetată de volum n relativă la această caracteristică. Vrem să determinăm eficiența estimatorului lui $\sigma = \sqrt{\text{Var}(X)}$, construit în Exemplul 5.2.6, dat prin

$$\bar{S} = \frac{1}{n} \sqrt{\frac{\pi}{2}} \sum_{i=1}^n |X_i - \mu|.$$

Știm că $E(\bar{S}) = \sigma$ și $\text{Var}(\bar{S}) = \frac{\pi-2}{2n} \sigma^2$, deci \bar{S} este estimator absolut corect pentru σ .

Pe de altă parte avem

$$f(x; \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}},$$

de unde

$$\frac{\partial^2 \ln f(x; \sigma)}{\partial \sigma^2} = \frac{1}{\sigma^2} - \frac{3}{\sigma^4} (x - \mu)^2.$$

Se calculează ușor

$$I_n(\sigma) = nI_1(\sigma) = -nE\left[\frac{1}{\sigma^2} - \frac{3(X - \mu)^2}{\sigma^4}\right] = \frac{2n}{\sigma^2}.$$

Avem astfel eficiența estimatorului \bar{S}

$$e(\bar{S}) = \frac{I_n^{-1}(\sigma)}{\text{Var}(\bar{S})} = \frac{\sigma^2}{2n} \cdot \frac{2n}{(\pi - 2)\sigma^2} = \frac{1}{\pi - 2} < 1,$$

adică estimatorul \bar{S} nu este estimator eficient pentru parametrul σ .

Teorema 5.2.17 (Rao-Cramer). *Se consideră caracteristica X cu funcția de probabilitate $f(x; \theta)$, $\theta \in (a, b)$, care satisface condițiile Teoremei 5.2.13 și fie funcția de estimare absolut corectă $\hat{\Theta} = \hat{\theta}(X_1, X_2, \dots, X_n)$ pentru parametrul θ . Condiția necesară și suficientă ca $\hat{\Theta}$ să fie funcție de estimare eficientă pentru parametrul θ este ca să aibă loc reprezentarea*

$$\ln f(x; \theta) = A'(\theta)[L(x) - \theta] + A(\theta) + N(x),$$

în plus are loc formula

$$\hat{\Theta} = \hat{\theta}(X_1, X_2, \dots, X_n) = \frac{1}{n} \sum_{k=1}^n L(X_k).$$

Demonstrație. Pentru necesitate, din demonstrația inegalității Rao-Cramer, avem că egalitatea în această inegalitate are loc dacă inegalitatea lui Schwarz este verificată prin egalitate. Aceasta are loc dacă și numai dacă variabilele aleatoare considerate depind în mod liniar, adică

$$K(\theta) [\hat{\theta}(X_1, X_2, \dots, X_n) - \theta] = \sum_{k=1}^n \frac{\partial \ln f(X_k; \theta)}{\partial \theta}, \quad K \neq 0.$$

Considerând $X_1 = x$ (oarecare), $X_2 = \dots = X_n = \text{const.}$, (=0 de exemplu) avem

$$K(\theta) [\hat{\theta}(x, 0, \dots, 0) - \theta] = \frac{\partial \ln f(x; \theta)}{\partial \theta} + (n-1) \frac{\partial \ln f(0; \theta)}{\partial \theta},$$

de unde

$$\frac{\partial \ln f(x; \theta)}{\partial \theta} = U(\theta) Q(x) + V(\theta).$$

Astfel s-a obținut că

$$\sum_{k=1}^n \frac{\partial \ln f(x_k; \theta)}{\partial \theta} = U(\theta) \sum_{k=1}^n Q(x_k) + nV(\theta)$$

și prin urmare

$$K(\theta) [\hat{\theta}(x_1, x_2, \dots, x_n) - \theta] = U(\theta) \sum_{k=1}^n Q(x_k) + nV(\theta),$$

de unde

$$\hat{\theta}(x_1, x_2, \dots, x_n) = \frac{U(\theta)}{K(\theta)} \sum_{k=1}^n Q(x_k) + \frac{nV(\theta)}{K(\theta)} + \theta,$$

pentru orice x_1, x_2, \dots, x_n . Rezultă de aici că

$$h = \frac{U(\theta)}{K(\theta)}, \quad g = \frac{nV(\theta)}{K(\theta)} + \theta$$

sunt constante (nu depind de θ), deoarece $\hat{\theta}(x_1, x_2, \dots, x_n)$ nu depinde de θ . Prin urmare

$$\hat{\theta}(x_1, x_2, \dots, x_n) = h \sum_{k=1}^n Q(x_k) + g,$$

și dacă se notează $L(x) = nhQ(x) + g$, atunci

$$\hat{\theta}(x_1, x_2, \dots, x_n) = h \sum_{k=1}^n \frac{L(x_k) - g}{nh} + g = \frac{1}{n} \sum_{k=1}^n L(x_k).$$

Așadar s-a ajuns la

$$\hat{\Theta} = \hat{\theta}(X_1, X_2, \dots, X_n) = \frac{1}{n} \sum_{k=1}^n L(X_k).$$

Pe de altă parte, deoarece

$$\frac{\partial \ln f(x; \theta)}{\partial \theta} = U(\theta) Q(x) + V(\theta),$$

obținem

$$\frac{\partial \ln f(x; \theta)}{\partial \theta} = hK(\theta) \frac{L(x) - g}{nh} + \frac{K(\theta)(g - \theta)}{n},$$

de unde

$$\frac{\partial \ln f(x; \theta)}{\partial \theta} = \frac{K(\theta)}{n} [L(x) - \theta].$$

Integrăm, după ce am notat $A''(\theta) = \frac{K(\theta)}{n}$, ultima relație în raport cu θ și obținem

$$\begin{aligned} \int \frac{\partial \ln f(x; \theta)}{\partial \theta} d\theta &= \int A''(\theta) [L(x) - \theta] d\theta = A'(\theta) [L(x) - \theta] + \int A'(\theta) d\theta \\ &= A'(\theta) [L(x) - \theta] + A(\theta) + N(x), \end{aligned}$$

de unde

$$\ln f(x; \theta) = A'(\theta) [L(x) - \theta] + A(\theta) + N(x).$$

Pentru suficiență se pornește de la relația cunoscută

$$\int_{\mathbb{R}} f(x; \theta) dx = 1,$$

adică

$$\int_{\mathbb{R}} \exp [A'(\theta) [L(x) - \theta] + A(\theta) + N(x)] dx = 1,$$

pe care o derivăm în raport cu θ . Obținem astfel că

$$\int_{\mathbb{R}} A''(\theta) [L(x) - \theta] f(x; \theta) dx = 0,$$

de unde, deoarece $A'' \neq 0$, rezultă că

$$\int_{\mathbb{R}} [L(x) - \theta] f(x; \theta) dx = 0.$$

Derivăm din nou în raport cu θ și se ajunge la

$$-\int_{\mathbb{R}} f(x; \theta) dx + \int_{\mathbb{R}} [L(x) - \theta] \frac{\partial f(x; \theta)}{\partial \theta} dx = 0,$$

de unde

$$\int_{\mathbb{R}} [L(x) - \theta]^2 A''(\theta) f(x; \theta) dx = 1.$$

Deoarece $E[L(X)] = \theta$, am obținut astfel că

$$\text{Var}[L(X)] = \frac{1}{A''(\theta)},$$

deci

$$\text{Var}(\hat{\Theta}) = \text{Var}\left[\frac{1}{n} \sum_{k=1}^n L(X_k)\right] = \frac{1}{nA''(\theta)}.$$

Pe de altă parte

$$\begin{aligned} I_n(\theta) &= nI_1(\theta) = n \int_{\mathbb{R}} \left[\frac{\partial \ln f(x; \theta)}{\partial \theta} \right]^2 f(x; \theta) dx \\ &= n \int_{\mathbb{R}} [A''(\theta)]^2 [L(x) - \theta]^2 f(x; \theta) dx = n[A''(\theta)]^2 \frac{1}{A''(\theta)} = nA''(\theta). \end{aligned}$$

Prin urmare putem scrie

$$e(\hat{\Theta}) = \frac{I_n^{-1}(\theta)}{\text{Var}(\hat{\Theta})} = \frac{\frac{1}{nA''(\theta)}}{\frac{1}{nA''(\theta)}} = 1,$$

deci $\hat{\Theta} = \hat{\theta}(X_1, X_2, \dots, X_n)$ este estimator eficient pentru parametrul θ . \square

Exemplul 5.2.18. Se consideră caracteristica X ce urmează legea $\mathcal{B}(m, p)$, cu m cunoscut și p necunoscut. Vrem să arătăm că media de selecție este estimator eficient pentru parametrul necunoscut $\theta = E(X) = mp$. Pentru aceasta vom considera o selecție repetată de volum n relativă la această caracteristică.

Funcția de frecvență a caracteristicii X este

$$f(x; \theta) = \binom{m}{x} \left(\frac{\theta}{m}\right)^x \left(1 - \frac{\theta}{m}\right)^{m-x},$$

de unde

$$\begin{aligned} \ln f(x; \theta) &= \ln \binom{m}{x} + x \ln \frac{\theta}{m} + (m-x) \ln \left(1 - \frac{\theta}{m}\right) \\ &= (x - \theta) \ln \frac{\theta}{m - \theta} + \theta \ln \frac{\theta}{m - \theta} + m \ln(m - \theta) + \ln \binom{m}{x} - m \ln m. \end{aligned}$$

Considerând

$$A(\theta) = \theta \ln \theta + (m - \theta) \ln(m - \theta), \quad \text{avem} \quad A'(\theta) = \ln \frac{\theta}{m - \theta},$$

deci

$$f(x; \theta) = [L(x) - \theta] A'(\theta) + A(\theta) + N(x),$$

unde $L(x) = x$ și $N(x) = \ln \binom{m}{x} - m \ln m$. Pe baza teoremei Rao–Cramer se obține că

$$\bar{\Theta} = \frac{1}{n} \sum_{k=1}^n L(X_k) = \frac{1}{n} \sum_{k=1}^n X_k$$

este estimator eficient pentru parametrul $\theta = mp$.

5.2.4 Estimatori optimali

Definiția 5.2.19. *Estimatorul nedeplasat $\hat{\Theta} = \hat{\theta}(X_1, X_2, \dots, X_n)$, pentru parametrul θ , este optimal dacă are dispersia minimă dintre toți estimatorii nedeplasați ai lui θ .*

Observația 5.2.20. Un estimator eficient este optimal, invers proprietatea nu are loc.

Proprietatea 5.2.21. *Estimatorul optimal al unui parametru este unic.*

Demonstrație. Fie $\hat{\Theta}_1 = \hat{\theta}_1(X_1, X_2, \dots, X_n)$ și $\hat{\Theta}_2 = \hat{\theta}_2(X_1, X_2, \dots, X_n)$ doi estimatori optimali distincți pentru parametrul θ . Dacă se consideră funcția de selecție $\hat{\Theta} = \frac{1}{2}(\hat{\Theta}_1 + \hat{\Theta}_2)$, atunci $\hat{\Theta}$ este estimator nedeplasat pentru θ , deoarece

$$E(\hat{\Theta}) = \frac{1}{2}E(\hat{\Theta}_1) + \frac{1}{2}E(\hat{\Theta}_2) = \frac{1}{2}(\theta + \theta) = \theta.$$

Pe de altă parte

$$\begin{aligned} Var(\hat{\Theta}) &= \frac{1}{4} \left[Var(\hat{\Theta}_1) + Var(\hat{\Theta}_2) + 2Cov(\hat{\Theta}_1, \hat{\Theta}_2) \right] \\ &= \frac{1}{2} \left[\sigma^2 + Cov(\hat{\Theta}_1, \hat{\Theta}_2) \right], \end{aligned}$$

unde $\sigma^2 = Var(\hat{\Theta}_1) = Var(\hat{\Theta}_2)$.

Din inegalitatea lui Schwarz se obține

$$Cov^2(\hat{\Theta}_1, \hat{\Theta}_2) = E^2[(\hat{\Theta}_1 - \theta)(\hat{\Theta}_2 - \theta)] \leq E[(\hat{\Theta}_1 - \theta)^2] E[(\hat{\Theta}_2 - \theta)^2] = \sigma^4,$$

deci $Cov(\hat{\Theta}_1, \hat{\Theta}_2) \leq \sigma^2$. Prin urmare, putem scrie că $Var(\hat{\Theta}) \leq \sigma^2$, ceea ce contrazice faptul că $\hat{\Theta}_1, \hat{\Theta}_2$ sunt estimatori optimali, afară de cazul în care $Var(\hat{\Theta}) = \sigma^2$, adică $Cov(\hat{\Theta}_1, \hat{\Theta}_2) = \sigma^2$. Având în vedere aceste rezultate, putem scrie:

$$Var(\hat{\Theta}_1 - \hat{\Theta}_2) = Var(\hat{\Theta}_1) + Var(\hat{\Theta}_2) - 2Cov(\hat{\Theta}_1, \hat{\Theta}_2) = \sigma^2 + \sigma^2 - 2\sigma^2 = 0,$$

adică $Var(\hat{\Theta}_1 - \hat{\Theta}_2) = 0$. Dar avem că $E(\hat{\Theta}_1) = E(\hat{\Theta}_2)$, ceea ce implică faptul că $E[(\hat{\Theta}_1 - \hat{\Theta}_2)^2] = 0$. Prin urmare se ajunge la $\hat{\Theta}_1 = \hat{\Theta}_2$, în contradicție cu afirmația că cei doi estimatori sunt distincți. \square

Teorema 5.2.22 (Rao-Blackwell). Fie caracteristica X cu funcția de probabilitate $f(x; \theta)$ și fie $\hat{\Theta} = \hat{\theta}(X_1, \dots, X_n)$ un estimator nedeplasat pentru parametrul θ . Dacă statistica $S = S(X_1, \dots, X_n)$ este o statistică suficientă pentru parametrul θ , atunci estimatorul $\bar{\Theta} = \bar{\theta}(X_1, \dots, X_n) = E(\hat{\Theta} | S)$ este un estimator nedeplasat pentru θ și are loc relația $Var(\bar{\Theta}) \leq Var(\hat{\Theta})$.

Demonstrație. Deoarece statistica S este suficientă rezultă că funcția de probabilitate $f(x_1, \dots, x_n; \theta | s)$ condiționată a vectorului aleator (X_1, \dots, X_n) de evenimentul $(S = s)$ nu depinde de θ , deci

$$E(\hat{\Theta} | S = s) = \int \dots \int_{\mathbb{R}^n} \hat{\theta}(x_1, \dots, x_n) f(x_1, \dots, x_n; \theta | s) dx_1 \dots dx_n$$

nu depinde de θ , adică este o funcție de selecție ce nu depinde de θ . Pe de altă parte, folosind proprietăți ale mediei condiționate avem

$$E(\bar{\Theta}) = E[E(\hat{\Theta} | S)] = E(\hat{\Theta}) = \theta,$$

deci nedeplasarea are loc.

Pentru a stabili inegalitatea dintre dispersiile celor doi estimatori avem că

$$\begin{aligned} Var(\hat{\Theta}) &= E[(\hat{\Theta} - \theta)^2] = E[(\hat{\Theta} - E(\hat{\Theta} | S) + E(\hat{\Theta} | S) - \theta)^2] \\ &= E[(\hat{\Theta} - E(\hat{\Theta} | S))^2] + E[(E(\hat{\Theta} | S) - \theta)^2] \\ &\quad + 2E[(\hat{\Theta} - E(\hat{\Theta} | S))(E(\hat{\Theta} | S) - \theta)]. \end{aligned}$$

Dar avem că

$$E[Var(\hat{\Theta} | S)] = E[E((\hat{\Theta} - E(\hat{\Theta} | S))^2 | S)] = E[(\hat{\Theta} - E(\hat{\Theta} | S))^2],$$

adică primul termen din expresia lui $Var(\hat{\Theta})$. Al doilea termen din aceeași expresie este

$$E[(E(\hat{\Theta} | S) - \theta)^2] = E[(\bar{\Theta} - \theta)^2] = Var(\bar{\Theta}),$$

iar ultimul termen este nul deoarece pentru $S = s$ fixat, avem $E(\hat{\Theta} | s) - \theta = \text{const.}$, deci

$$\begin{aligned} E[(\hat{\Theta} - E(\hat{\Theta} | s))(E(\hat{\Theta} | s) - \theta)] &= [E(\hat{\Theta} | s) - \theta] E[\hat{\Theta} - E(\hat{\Theta} | s)] \\ &= [E(\hat{\Theta} | s) - \theta] [E(\hat{\Theta}) - E(\hat{\Theta})] = 0. \end{aligned}$$

Am stabilit astfel relația $Var(\hat{\Theta}) = E[Var(\hat{\Theta} | S)] + Var(\bar{\Theta})$, de unde $Var(\bar{\Theta}) \leq Var(\hat{\Theta})$. \square

Definiția 5.2.23. Statistica $S = S(X_1, X_2, \dots, X_n)$ este completă pentru familia de legi de probabilitate $f(x; \theta)$, $\theta \in A$, dacă $E[\varphi(S)] = 0$, pentru orice $\theta \in A$, implică faptul că $\varphi = 0$ a.s.

Exemplul 5.2.24. Statistica $S = \sum_{k=1}^n X_k$ este completă pentru familia de legi Poisson cu parametrul $\lambda > 0$.

Prima dată avem că S urmează legea lui Poisson de parametru $n\lambda$, deci

$$E[\varphi(S)] = \sum_{s=0}^{\infty} \varphi(s) \frac{(n\lambda)^s}{s!} e^{-n\lambda} = e^{-n\lambda} \sum_{s=0}^{\infty} \frac{\varphi(s) n^s}{s!} \lambda^s.$$

Pe de altă parte $E[\varphi(S)] = 0$, pentru orice $\lambda > 0$, conduce la

$$\sum_{s=0}^{\infty} \frac{\varphi(s) n^s}{s!} \lambda^s = 0,$$

pentru orice $\lambda > 0$, ceea ce are loc numai dacă fiecare coeficient este nul, adică $\varphi(s) = 0$, când $s \in \mathbb{N}$.

Teorema 5.2.25 (Lehmann–Scheffé). În condițiile Teoremei Rao–Blackwell, dacă statistica $S = S(X_1, X_2, \dots, X_n)$ este completă, atunci estimatorul

$$\bar{\Theta} = \bar{\theta}(X_1, X_2, \dots, X_n) = E(\hat{\Theta} | S)$$

este estimator optimal.

Demonstrație. Fie $\tilde{\Theta} = \tilde{\theta}(X_1, X_2, \dots, X_n)$ un estimator nedeplasat pentru parametrul θ . Folosind Teorema Rao–Blackwell, avem că estimatorul

$$\Theta_1 = \theta_1(X_1, X_2, \dots, X_n) = E(\tilde{\Theta} | S)$$

este nedeplasat pentru parametrul θ , adică $E(\Theta_1) = \theta$, și $Var(\Theta_1) \leq Var(\tilde{\Theta})$. Așadar avem că $E(\Theta_1) = E(\bar{\Theta}) = \theta$ sau $E[E(\hat{\Theta} | S)] = E[E(\tilde{\Theta} | S)] = \theta$, de unde $E[E(\hat{\Theta} | S) - E(\tilde{\Theta} | S)] = 0$. Având în vedere că statistica S este completă rezultă că $E(\hat{\Theta} | s) = E(\tilde{\Theta} | s)$ a.s., deci $\bar{\Theta} = \Theta_1$ a.s., de unde $Var(\bar{\Theta}) = Var(\Theta_1)$.

În final $Var(\bar{\Theta}) = Var(\Theta_1) \leq Var(\tilde{\Theta})$, adică $Var(\bar{\Theta}) \leq Var(\tilde{\Theta})$, pentru orice estimator nedeplasat $\tilde{\Theta}$, ceea ce este chiar condiția din definiția optimalității. \square

Exemplul 5.2.26. Fie caracteristica X , care urmează legea lui Poisson de parametru $\lambda > 0$. Vrem să determinăm un estimator optimal pentru parametrul $\theta = e^{-\lambda}$.

Funcția de frecvență a caracteristicii X este

$$f(x; \theta) = \frac{\lambda^x}{x!} e^{-\lambda} = \theta \frac{1}{x!} \left(\ln \frac{1}{\theta} \right)^x.$$

În Exemplul 5.2.24 am văzut că statistica suficientă $S = \sum_{k=1}^n X_k$ este completă pentru familia de legi Poisson și urmează legea lui Poisson de parametru $n\lambda = n \ln \frac{1}{\theta}$.

Dacă se consideră funcția de selecție

$$\hat{\Theta} = \hat{\theta}(X_1, X_2, \dots, X_n) = \frac{1}{n} \text{card} \{ X_i \mid X_i = 0, i = \overline{1, n} \},$$

atunci $\hat{\Theta}$ are distribuția $\hat{\Theta} \left(\binom{\frac{k}{n}}{\binom{n}{k} \theta^k (1-\theta)^{n-k}} \right)_{k=\overline{0, n}}$, deci $E(\hat{\Theta}) = \frac{1}{n} n\theta = \theta$,

adică $\hat{\Theta}$ este estimator nedeplasat pentru parametrul θ .

Dacă se introduc variabilele aleatoare Y_1, Y_2, \dots, Y_n cu distribuțiile date prin

$$Y_i = \begin{cases} 1, & \text{dacă } X_i = 0, \\ 0, & \text{dacă } X_i \neq 0, \end{cases}$$

deci $P(Y_i = 1) = \theta$, $P(Y_i = 0) = 1 - \theta$, atunci avem că $\hat{\Theta} = \frac{1}{n} \sum_{k=1}^n Y_k$.

Aplicăm Teorema Rao–Blackwell și obținem

$$\bar{\Theta} = E(\hat{\Theta} \mid S) = \frac{1}{n} E\left(\sum_{k=1}^n Y_k \mid S\right) = E(Y_1 \mid S).$$

Pentru $S = s$, avem

$$E(Y_1 \mid S = s) = P(Y_1 = 1 \mid S = s) \cdot 1 = P(X_1 = 0 \mid S = s).$$

Folosind formula lui Bayes avem

$$\begin{aligned} P(X_1 = 0 \mid S = s) &= \frac{P(X_1 = 0) P(S = s \mid X_1 = 0)}{P(S = s)} \\ &= \frac{P(X_1 = 0) P(X_2 + \dots + X_n = s)}{P(S = s)} \\ &= \frac{\exp(-\lambda) \frac{((n-1)\lambda)^s}{s!} \exp[-(n-1)\lambda]}{\frac{(n\lambda)^s}{s!} \exp(-n\lambda)} = \left(\frac{n-1}{n}\right)^s, \end{aligned}$$

de unde se obține că

$$\bar{\Theta} = E\left(\hat{\Theta} \mid S\right) = E(Y_1 \mid S) = \left(1 - \frac{1}{n}\right)^S.$$

Am ajuns astfel la estimatorul optimal pentru parametrul $\theta = e^{-\lambda}$, care este dat prin formula

$$\bar{\Theta} = \bar{\theta}(X_1, X_2, \dots, X_n) = \left(1 - \frac{1}{n}\right)^{n\bar{X}}.$$

5.3 Metode pentru estimarea parametrilor

5.3.1 Metoda momentelor

Se consideră caracteristica X , care are funcția de probabilitate $f(x; \theta)$ cu parametrul necunoscut $\theta = (\theta_1, \theta_2, \dots, \theta_p) \in \mathcal{A} \subset \mathbb{R}^p$ și o selecție repetată de volum n .

Definiția 5.3.1. Numim estimator pentru parametrul θ obținut prin metoda momentelor soluția $\bar{\Theta} = (\bar{\Theta}_1, \bar{\Theta}_2, \dots, \bar{\Theta}_p)$ a sistemului

$$\nu_k = \bar{\nu}_k, \quad k = \overline{1, p},$$

unde ν_k este momentul teoretic de ordin k , $\nu_k = E(X^k)$, iar $\bar{\nu}_k$ este momentul de selecție de ordinul k , adică

$$\bar{\nu}_k = \frac{1}{n} \sum_{i=1}^n X_i^k.$$

Observația 5.3.2. Metoda momentelor este una din cele mai vechi metode de estimare a parametrilor, folosită inițial în mod empiric. Este fundamentată teoretic pe faptul că momentele de selecție sunt estimatori absolut corecți pentru momentele teoretice corespunzătoare. Este cunoscut de asemenea că $\bar{\Theta} \xrightarrow{\text{a.s.}} \theta$, când $n \rightarrow \infty$.

Exemplul 5.3.3. Se consideră caracteristica X , care urmează legea gamma cu parametrii $a, b > 0$ necunoscuți. Densitatea de probabilitate pentru X este

$$f(x; a, b) = \frac{1}{\Gamma(a) b^a} x^{a-1} e^{-\frac{x}{b}}, \quad x > 0.$$

Vrem să estimăm parametrii a și b prin metoda momentelor.

Se știe că

$$\begin{aligned} \nu_1 = E(X) &= \int_{-\infty}^{+\infty} x f(x; a, b) = \frac{1}{\Gamma(a) b^a} \int_0^{+\infty} x^a e^{-\frac{x}{b}} dx = ab \quad \text{și} \\ \nu_2 = E(X^2) &= \int_{-\infty}^{+\infty} x^2 f(x; a, b) = \frac{1}{\Gamma(a) b^a} \int_0^{+\infty} x^{a+1} e^{-\frac{x}{b}} dx = ab^2 (a+1). \end{aligned}$$

Prin urmare, avem sistemul de ecuații

$$\begin{cases} ab = \bar{\nu}_1 = \bar{X} \\ ab^2 (a+1) = \bar{\nu}_2 = \bar{X}^2 + \bar{\mu}_2, \end{cases}$$

care are soluția

$$\bar{a} = \frac{\bar{X}^2}{\bar{\mu}_2}, \quad \bar{b} = \frac{\bar{\mu}_2}{\bar{X}}.$$

5.3.2 Metoda verosimilității maxime

Se consideră caracteristica X cu funcția de probabilitate $f(x; \theta)$, unde parametrul necunoscut $\theta \in \mathcal{A} \subset \mathbb{R}^p$. Relativ la caracteristica X se consideră o selecție repetată de volum n .

Definiția 5.3.4. Numim estimator de verosimilitate maximă, pentru parametrul θ , statistica

$$\hat{\theta} = \hat{\theta}(X_1, X_2, \dots, X_n)$$

pentru care se obține maximul funcției de verosimilitate

$$g(X_1, X_2, \dots, X_n; \theta) = \prod_{k=1}^n f(X_k; \theta),$$

iar $\hat{\theta} = \hat{\theta}(x_1, x_2, \dots, x_n)$ se numește estimație de verosimilitate maximă, pentru parametrul θ .

Observația 5.3.5. În definiția estimatorului de verosimilitate maximă $\hat{\theta}$ nu este necesar ca $f(x; \theta)$ să fie diferențiabilă în raport cu θ . De asemenea, nu este neapărat unic și nedeplasat.

Observația 5.3.6. Dacă funcția de verosimilitate este diferențiabilă de două ori în raport cu θ , atunci estimatorul de verosimilitate maximă se obține ca soluție a sistemului de ecuații

$$\frac{\partial g(X_1, X_2, \dots, X_n; \theta)}{\partial \theta_k} = 0, \quad k = \overline{1, p}.$$

Sistemul de ecuații este echivalent cu

$$\frac{\partial \ln g(X_1, X_2, \dots, X_n; \theta)}{\partial \theta_k} = \sum_{i=1}^n \frac{\partial \ln f(X_i; \theta)}{\partial \theta_k} = 0, \quad k = \overline{1, p},$$

numit *sistemul ecuațiilor de verosimilitate maximă*.

Exemplul 5.3.7. Să determinăm estimatorii de verosimilitate maximă pentru valoarea medie și abaterea standard, dacă se consideră caracteristica X , care urmează legea normală $\mathcal{N}(\mu, \sigma)$.

Se știe că $E(X) = \mu$ și $\sigma(X) = \sigma$, iar densitatea de probabilitate a lui X este

$$f(x; \mu, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}.$$

Pentru a scrie sistemul de verosimilitate maximă, avem că

$$\ln f(x; \mu, \sigma) = -\ln \sqrt{2\pi} - \ln \sigma - \frac{(x-\mu)^2}{2\sigma^2},$$

de unde

$$\frac{\partial \ln f(x; \mu, \sigma)}{\partial \mu} = \frac{x-\mu}{\sigma^2}, \quad \frac{\partial \ln f(x; \mu, \sigma)}{\partial \sigma} = -\frac{1}{\sigma} + \frac{(x-\mu)^2}{\sigma^3}.$$

În acest mod se obține

$$\begin{cases} \frac{\partial \ln g}{\partial \mu} = \sum_{k=1}^n \frac{\partial \ln f(X_k; \mu, \sigma)}{\partial \mu} = \sum_{k=1}^n \frac{X_k - \mu}{\sigma^2} = 0 \\ \frac{\partial \ln g}{\partial \sigma} = \sum_{k=1}^n \frac{\partial \ln f(X_k; \mu, \sigma)}{\partial \sigma} = \sum_{k=1}^n \left[-\frac{1}{\sigma} + \frac{(X_k - \mu)^2}{\sigma^3} \right] = 0, \end{cases}$$

sau

$$\begin{cases} \sum_{k=1}^n (X_k - \mu) = 0 \\ \sum_{k=1}^n \left[-\sigma^2 + (X_k - \mu)^2 \right] = 0, \end{cases}$$

de unde se rezultă estimatorii de verosimilitate maximă

$$\mu^* = \frac{1}{n} \sum_{k=1}^n X_k = \bar{X}, \quad \sigma^* = \sqrt{\frac{1}{n} \sum_{k=1}^n (X_k - \bar{X})^2} = \sqrt{\mu_2},$$

pentru parametrii μ și σ .

Exemplul 5.3.8. Caracteristica X urmează legea uniformă pe intervalul $(0, \theta)$, unde parametrul $\theta > 0$ este necunoscut. Relativ la caracteristica X se consideră o selecție repetată de volum n . Vrem să se determinăm estimatorul de verosimilitate maximă $\hat{\theta}$, pentru parametrul necunoscut θ .

Estimatorul $\hat{\theta}$ de verosimilitate maximă pentru θ se determină astfel încât funcția de verosimilitate

$$g(X_1, X_2, \dots, X_n; \theta) = \prod_{k=1}^n f(x_k; \theta) = \frac{1}{\theta^n}$$

să fie maximă pentru $\hat{\theta} = \hat{\theta}(X_1, X_2, \dots, X_n)$. Deoarece

$$\max \left\{ \frac{1}{\theta^n} \mid X_1 \leq \theta, \dots, X_n \leq \theta \right\} = \left(\frac{1}{\max\{X_i, i = \overline{1, n}\}} \right)^n,$$

se obține că $\hat{\theta} = \hat{\theta}(X_1, X_2, \dots, X_n) = \max\{X_i, i = \overline{1, n}\}$.

Se observă că în acest exemplu nu s-a putut folosi ecuația de verosimilitate maximă, deoarece domeniul valorilor caracteristicii X , care este $(0, \theta)$, depinde de parametrul estimat.

Vom arăta, în continuare, că estimatorul astfel construit este estimator corect pentru parametrul θ . Apoi vom folosi acest estimator pentru obținerea unui estimator absolut corect pentru θ .

Funcția de repartiție a statisticii $\hat{\theta}$ este

$$F_{\hat{\theta}}(x; \theta) = P(\hat{\theta} \leq x) = \prod_{i=1}^n P(X_i \leq x) = [F_X(x; \theta)]^n,$$

deci $\hat{\theta}$ are densitatea de probabilitate

$$f_{\hat{\theta}}(x; \theta) = \frac{\partial F_{\hat{\theta}}(x; \theta)}{\partial x} = n f(x; \theta) [F_X(x; \theta)]^{n-1} = \frac{nx^{n-1}}{\theta^n}, \quad \text{când } x \in (0, \theta).$$

Se calculează ușor

$$E(\hat{\theta}) = \frac{n}{\theta^n} \int_0^\theta x \cdot x^{n-1} dx = \frac{n}{n+1} \theta, \quad E(\hat{\theta}^2) = \frac{n}{\theta^n} \int_0^\theta x^{n+1} dx = \frac{n}{n+2} \theta^2.$$

Astfel se obține că

$$\lim_{n \rightarrow \infty} E(\hat{\Theta}) = \theta$$

$$\lim_{n \rightarrow \infty} Var(\hat{\Theta}) = \lim_{n \rightarrow \infty} \left[\frac{n}{n+2} \theta^2 - \frac{n^2}{(n+1)^2} \theta^2 \right] = \lim_{n \rightarrow \infty} \frac{n}{(n+1)^2 (n+2)} \theta^2 = 0.$$

Prin urmare, $\hat{\Theta}$ este estimator corect pentru θ .

Punând condiția $E(\bar{\Theta}) = \theta$, rezultă că

$$\theta = E(\bar{\Theta}) = \kappa_n E(\hat{\Theta}) = \kappa_n \frac{n}{n+1} \theta,$$

de unde se obține $\kappa_n = \frac{n+1}{n}$ și în final

$$\bar{\Theta} = \frac{n+1}{n} \hat{\Theta} = \frac{n+1}{n} \max \{ X_i, i = \overline{1, n} \}.$$

Deoarece

$$Var(\bar{\Theta}) = \left(\frac{n+1}{n} \right)^2 Var(\hat{\Theta}) = \frac{n(n+1)^2}{n^2(n+1)^2(n+2)} \theta = \frac{1}{n(n+2)} \theta \rightarrow 0,$$

când $n \rightarrow \infty$, rezultă că $\bar{\Theta}$ este estimator absolut corect pentru parametrul θ .

Proprietatea 5.3.9. Dacă $S = S(X_1, X_2, \dots, X_n)$ este statistică suficientă pentru parametrul θ , iar $\hat{\Theta}$ este estimator de verosimilitate maximă pentru θ , atunci $\hat{\Theta}$ este funcție de S .

Demonstrație. Deoarece statistica S este suficientă rezultă că

$$g(x_1, x_2, \dots, x_n; \theta) = \varphi(x_1, x_2, \dots, x_n) h(s, \theta),$$

deci maximul lui g , după θ , se obține atunci și numai atunci când se obține maximul lui h după θ . Astfel că $\hat{\Theta}$ se exprimă în funcție de S . \square

Teorema 5.3.10. Dacă $\hat{\Theta} = \hat{\theta}(X_1, X_2, \dots, X_n)$ este funcție de estimare eficientă pentru parametrul θ , atunci $\hat{\Theta}$ este estimator de verosimilitate maximă pentru θ .

Demonstrație. Deoarece $\hat{\Theta}$ este estimator eficient pentru θ , din inegalitatea Rao–Cramer (cu egalitate), avem că

$$\frac{\partial \ln g(X_1, X_2, \dots, X_n; \theta)}{\partial \theta} = K(\theta) (\hat{\Theta} - \theta).$$

Astfel că

$$\frac{\partial \ln g(X_1, X_2, \dots, X_n; \hat{\Theta})}{\partial \theta} = K(\hat{\Theta}) (\hat{\Theta} - \hat{\Theta}) = 0,$$

deci $\hat{\Theta}$ verifică ecuația verosimilității maxime. \square

Observația 5.3.11. Dacă valoarea adevărată a parametrului θ este θ_0 , se arată că funcția de verosimilitate maximă $\hat{\Theta}$ pentru parametrul θ are următoarele comportări asimptotice: $\hat{\Theta} \xrightarrow{\text{a.s.}} \theta_0$, când $n \rightarrow \infty$, iar $\sqrt{n}(\hat{\Theta} - \theta_0)$ converge în repartiție la legea normală $\mathcal{N}\left(0, \frac{1}{\sqrt{I(\theta_0)}}\right)$.

5.3.3 Metoda minimului χ^2

Se consideră colectivitatea \mathcal{C} și caracteristica X cercetată, cu legea de probabilitate dată prin funcția $f(x; \theta)$, unde $\theta = (\theta_1, \dots, \theta_s) \in \mathbf{A} \subset \mathbb{R}^s$. Domeniul valorilor lui X îl considerăm compus din clasele C_i , $i = \overline{1, k}$. Vom introduce notațiile următoare $p_i = p_i(\theta) = P(X \in C_i)$, adică probabilitatea ca un individ luat la întâmplare din colectivitatea \mathcal{C} să aparțină clasei C_i .

Dacă se consideră o selecție repetată de volum n cu datele de selecție x_1, \dots, x_n , respectiv variabilele de selecție X_1, \dots, X_n , notăm prin n_i frecvența absolută a datelor de selecție din clasa C_i . Fie N_i variabila aleatoare (de selecție) corespunzătoare lui n_i , atunci vectorul aleator $N = (N_1, \dots, N_k)$ urmează legea multinomială de parametri $p_i = p_i(\theta)$, $i = \overline{1, k}$.

Definiția 5.3.12. Estimatorul cu χ^2 minim pentru parametrul θ este estimatorul (funcția de selecție) $\bar{\Theta} = \bar{\theta}(X_1, X_2, \dots, X_n)$, care realizează valoarea minimă a expresiei

$$\chi^2 = \sum_{i=1}^k \frac{[N_i - np_i(\theta)]^2}{np_i(\theta)}$$

în raport cu parametrul θ , iar $\bar{\theta} = \bar{\theta}(x_1, x_2, \dots, x_n)$ se numește estimatie cu χ^2 minim pentru parametrul θ .

Observația 5.3.13. Dacă notăm $\hat{p}_i = \frac{n_i}{n}$, atunci valoarea lui χ^2 se poate scrie în felul următor:

$$\chi^2 = \sum_{i=1}^k \frac{(n_i - np_i)^2}{np_i} = \sum_{i=1}^k \frac{(n\hat{p}_i - np_i)^2}{np_i} = n \left(\sum_{i=1}^k \frac{\hat{p}_i^2}{p_i} - 1 \right),$$

prin urmare, minimul lui χ^2 se obține când este minimă expresia $\sum_{i=1}^k \frac{\hat{p}_i^2}{p_i}$.

Observația 5.3.14. Dacă $p_i(\theta)$, $i = \overline{1, s}$, sunt diferențiabile de două ori în raport cu θ , atunci $\bar{\Theta}$ se obține ca soluție a sistemului

$$\frac{\partial}{\partial \theta_j} \left[\sum_{i=1}^k \frac{(N_i - np_i(\theta))^2}{np_i(\theta)} \right] = 0, \quad j = \overline{1, s}.$$

Observația 5.3.15. În locul expresiei lui χ^2 de mai sus, numit *indicatorul lui Pearson*, se pot utiliza alte expresii ca:

$$D^2(\theta) = \sum_{i=1}^k \frac{[N_i - np_i(\theta)]^2}{N_i}, \quad \text{indicatorul lui Neyman},$$

$$D^2(\theta) = 2 \sum_{i=1}^k N_i \ln \left[\frac{N_i}{np_i(\theta)} \right], \quad \text{indicatorul de verosimilitate},$$

$$D^2(\theta) = 2n \sum_{i=1}^k p_i(\theta) \ln \left[\frac{np_i(\theta)}{N_i} \right], \quad \text{indicatorul lui Kullbach},$$

$$D^2(\theta) = \sum_{i=1}^k n\hat{p}_i(1 - \hat{p}_i) \left[\ln \frac{\hat{p}_i}{1 - \hat{p}_i} - \ln \frac{p_i(\theta)}{1 - p_i(\theta)} \right],$$

indicatorul lui Berkson.

5.3.4 Metoda intervalelor de încredere

Se consideră caracteristica X cu legea de probabilitate $f(x; \theta)$, unde parametrul necunoscut $\theta \in A \subset \mathbb{R}$. Fie o selecție repetată de volum n și numărul $\alpha \in (0, 1)$, numit *probabilitate de risc*, $1 - \alpha$ numindu-se *probabilitate de încredere*.

Definiția 5.3.16. Numim interval de încredere pentru parametrul θ , intervalul aleator

$$(\bar{\theta}_1, \bar{\theta}_2) = (\bar{\theta}_1(X_1, X_2, \dots, X_n), \bar{\theta}_2(X_1, X_2, \dots, X_n)),$$

unde statisticile $\bar{\theta}_1$ și $\bar{\theta}_2$ sunt astfel încât $P(\bar{\theta}_1 < \theta < \bar{\theta}_2) = 1 - \alpha$, iar intervalul numeric $(\bar{\theta}_1, \bar{\theta}_2) = (\bar{\theta}_1(x_1, x_2, \dots, x_n), \bar{\theta}_2(x_1, x_2, \dots, x_n))$, se numește valoarea intervalului de încredere pentru parametrul θ .

Observația 5.3.17. Pentru determinarea intervalului de încredere se consideră o statistică $S = S(X_1, X_2, \dots, X_n)$, care urmează o lege de probabilitate cunoscută, dar în expresia căreia apare parametrul θ . Se determină apoi intervalul numeric (s_1, s_2) astfel încât $P(S \in (s_1, s_2)) = 1 - \alpha$, relație care se va scrie, în mod echivalent, $P(\bar{\theta}_1 < \theta < \bar{\theta}_2) = 1 - \alpha$.

Observația 5.3.18. Cu cât α este mai mic și intervalul de încredere are lungimea mai mică cu atât estimăția parametrului necunoscut este mai bună.

Observația 5.3.19. Din relația $P(S \in (s_1, s_2)) = 1 - \alpha$, intervalul (s_1, s_2) nu este determinat în mod unic.

De la problemă la problemă se mai adaugă o condiție suplimentară, de exemplu fixarea valorii fie a lui s_1 , fie a lui s_2 . De asemenea, se poate considera o legătură între s_1 și s_2 dată de anumite ipoteze de lucru.

Interval de încredere pentru media teoretică – dispersie cunoscută

Se consideră caracteristica X ce urmează legea normală $\mathcal{N}(m, \sigma)$, unde $m \in \mathbb{R}$ este necunoscut, iar $\sigma > 0$ este cunoscut.

Pentru construirea unui interval de încredere pentru media teoretică m necunoscută efectuăm o selecție repetată de volum n și considerăm probabilitatea de încredere $1 - \alpha$ dată, $\alpha \in (0, 1)$.

Se construiește statistica

$$Z = \frac{\bar{X} - m}{\frac{\sigma}{\sqrt{n}}},$$

care urmează legea normală $\mathcal{N}(0, 1)$ (a se vedea Observația 4.2.6). Prin urmare, pentru α dat determinăm intervalul numeric (z_1, z_2) astfel încât

$$P(Z \in (z_1, z_2)) = \Phi(z_2) - \Phi(z_1) = 1 - \alpha,$$

unde

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

este funcția lui Laplace. Funcția lui Laplace este tabelată în Anexa I pentru valori pozitive ale argumentului, având în vedere că $\Phi(-x) = -\Phi(x)$.

Deoarece dubla inegalitate

$$z_1 < \frac{\bar{X} - m}{\frac{\sigma}{\sqrt{n}}} < z_2$$

este echivalentă cu

$$\bar{m}_1 = \bar{X} - z_2 \frac{\sigma}{\sqrt{n}} < m < \bar{X} - z_1 \frac{\sigma}{\sqrt{n}} = \bar{m}_2$$

rezultă că $P(\bar{m}_1 < m < \bar{m}_2) = 1 - \alpha$, adică (\bar{m}_1, \bar{m}_2) este un interval de încredere pentru media teoretică m .

Desigur că intervalul numeric (z_1, z_2) nu este în mod unic determinat. Dacă nu există nici o informație suplimentară relativă la valoarea medie, atunci se va considera intervalul de încredere de lungime minimă, pentru α fixat, și se obține când $z_1 = -z_2$.

În acest caz $z_2 = z_{1-\frac{\alpha}{2}}$, va fi dat prin relația $\Phi(z_{1-\frac{\alpha}{2}}) - \Phi(-z_{1-\frac{\alpha}{2}}) = 1 - \alpha$, ceea ce este echivalent cu $\Phi(z_{1-\frac{\alpha}{2}}) = \frac{1-\alpha}{2}$.

Când se folosește funcția lui Laplace definită prin

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{t^2}{2}} dt,$$

atunci $z_{1-\frac{\alpha}{2}}$ se determină din relația $\Phi\left(z_{1-\frac{\alpha}{2}}\right) = 1 - \frac{\alpha}{2}$ și reprezintă de fapt cuantila de ordin $1 - \frac{\alpha}{2}$.

Intervalul de încredere pentru parametrul m are extremitățile

$$(5.3.1) \quad \begin{aligned} \bar{m}_1 &= \bar{m}_1(X_1, X_2, \dots, X_n) = \bar{X} - z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}, \\ \bar{m}_2 &= \bar{m}_2(X_1, X_2, \dots, X_n) = \bar{X} + z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}. \end{aligned}$$

Observația 5.3.20. Geometric, dacă se consideră graficele funcției de repartiție și densității de probabilitate a legii normale $\mathcal{N}(0, 1)$, modul de obținere a cuantilei z_γ este prezentat în Figura 5.1. Adică, fie intersectând graficul funcției de repartiție a legii normale $\mathcal{N}(0, 1)$ cu dreapta de ecuație $y = \gamma$, fie alegându-l pe z_γ astfel ca aria umbră de sub graficul densității de probabilitate să fie γ .

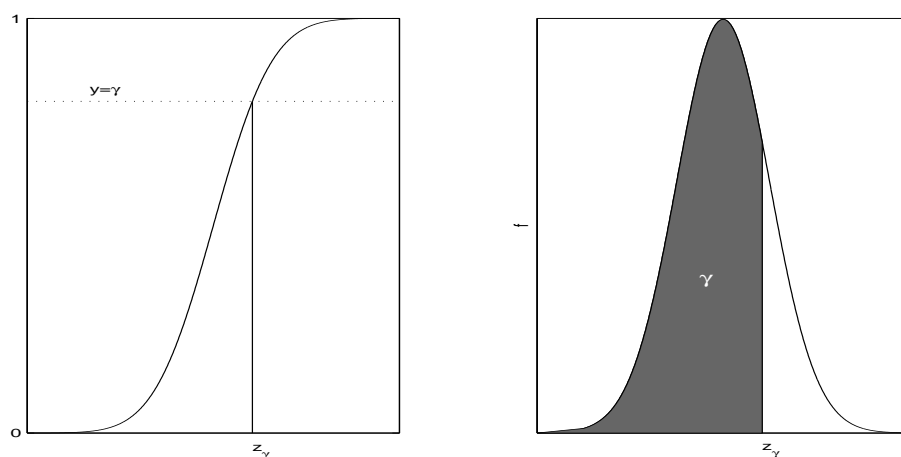


Figura 5.1: Determinarea cuantilei z_γ

Observația 5.3.21. Dacă există o informație relativă la valoarea medie de forma că aceasta nu este limitată superior, atunci intervalul numeric $(z_1, z_2) = (-\infty, z_{1-\alpha})$, care conduce la intervalul de încredere

$$(\bar{m}_1, \bar{m}_2) = \left(\bar{X} - \frac{\sigma}{\sqrt{n}} z_{1-\alpha}, \infty \right).$$

În mod analog, dacă se cunoaște că valoarea medie nu este limitată inferior, adică are tendința de a avea valori mici în raport cu ceea ce ne așteptăm, atunci intervalul

numeric $(z_1, z_2) = (z_\alpha, \infty)$, care conduce la intervalul de încredere

$$(\bar{m}_1, \bar{m}_2) = \left(-\infty, \bar{X} - \frac{\sigma}{\sqrt{n}} z_\alpha \right).$$

Pe baza teoremei limită centrală avem că rezultatul obținut se menține când X urmează o lege de probabilitate oarecare, pentru $n > 30$ (a se vedea Proprietatea 4.2.7).

Exemplul 5.3.22. Relativ la populația \mathcal{C} se cercetează caracteristica X privind media teoretică $E(X) = m$. Știind că dispersia teoretică a caracteristicii X este $Var(X) = 0.35$, să se stabilească un interval de încredere pentru media teoretică m cu probabilitatea de încredere $1 - \alpha = 0.95$, utilizând distribuția empirică de selecție

$$X \begin{pmatrix} 22.7 & 22.8 & 22.9 & 23.0 & 23.1 & 23.2 & 23.3 & 23.4 \\ 1 & 3 & 7 & 4 & 6 & 7 & 5 & 2 \end{pmatrix}.$$

Deoarece volumul selecției este $n = 35 > 30$, putem considera că statistica

$$Z = \frac{\bar{X} - m}{\frac{\sigma}{\sqrt{n}}}, \quad \text{unde } \sigma = \sqrt{Var(X)},$$

urmează legea normală $\mathcal{N}(0, 1)$. Așadar, extremitățile intervalului de încredere pentru media teoretică m sunt date prin (5.3.1). Calculăm valorile acestor extremități pe baza datelor de selecție.

Valoarea mediei de selecție \bar{X} este

$$\bar{x} = \frac{1}{35}(1 \cdot 22.7 + 3 \cdot 22.8 + 7 \cdot 22.9 + 4 \cdot 23 + 6 \cdot 23.1 + 7 \cdot 23.2 + 5 \cdot 23.3 + 2 \cdot 23.4) = 23.077,$$

iar din Anexa I, pentru $\frac{1-\alpha}{2} = 0.475$, se găsește $z_{1-\frac{\alpha}{2}} = 1.96$.

De asemenea, avem că

$$\frac{\sigma}{\sqrt{n}} = \sqrt{\frac{Var(X)}{n}} = \sqrt{\frac{0.35}{35}} = \sqrt{0.01} = 0.1.$$

Obținem, în acest fel, intervalul de încredere pentru media teoretică $m = E(X)$

$$\begin{aligned} \left(\bar{x} - \frac{\sigma}{\sqrt{n}} z_{1-\frac{\alpha}{2}}; \bar{x} + \frac{\sigma}{\sqrt{n}} z_{1-\frac{\alpha}{2}} \right) &= (23.077 - 0.196; 23.077 + 0.196) \\ &= (22.881; 23.273). \end{aligned}$$

Programul 5.3.23. Calculele pot fi făcute cu următorul program Matlab:

```

x=[22.7,22.8*ones(1,3),22.9*ones(1,7),23*ones(1,4),...
  23.1*ones(1,6),23.2*ones(1,7),23.3*ones(1,5),23.4*ones(1,2)];
s=sqrt(0.35)*norminv(0.975)/sqrt(35);ma=mean(x);
m1=ma-s; m2=ma+s;
fprintf(' (m1,m2)=(%6.3f,%6.3f)',m1,m2)

```

și care, în urma executării, afișează rezultatul

```
(m1,m2)=(22.881,23.273)
```

Interval de încredere pentru media teoretică – dispersia necunoscută

În condițiile secțiunii precedente, dar cu $\sigma > 0$ necunoscut, se va considera statistica

$$T = \frac{\bar{X} - m}{\frac{\bar{\sigma}}{\sqrt{n}}} = \frac{\bar{X} - m}{\sqrt{\frac{\bar{\mu}_2}{n-1}}},$$

care, conform Proprietății 4.2.26, urmează legea Student cu $n - 1$ grade de libertate.

Se determină intervalul numeric (t_1, t_2) astfel încât

$$P(T \in (t_1, t_2)) = F_{n-1}(t_2) - F_{n-1}(t_1) = 1 - \alpha,$$

unde

$$F_m(x) = \frac{\Gamma\left(\frac{m+1}{2}\right)}{\sqrt{m\pi}\Gamma\left(\frac{m}{2}\right)} \int_{-\infty}^x \left(1 + \frac{t^2}{m}\right)^{-\frac{m+1}{2}} dt, \quad x \in \mathbb{R},$$

este funcția de repartiție a legii Student cu m grade de libertate și care este tabelată, pentru anumite valori, în *Anexa II*.

Luând $t_2 = t_{n-1, 1-\frac{\alpha}{2}}$, $t_1 = -t_2$, avem că $F_{n-1}\left(t_{n-1, 1-\frac{\alpha}{2}}\right) = 1 - \frac{\alpha}{2}$, iar $P(\bar{m}_1 < m < \bar{m}_2) = 1 - \alpha$, deci, intervalul de încredere pentru media teoretică m are extremitățile date prin

$$(5.3.2) \quad \begin{aligned} \bar{m}_1 &= \bar{X} - t_{n-1, 1-\frac{\alpha}{2}} \frac{\bar{\sigma}}{\sqrt{n}}, \\ \bar{m}_2 &= \bar{X} + t_{n-1, 1-\frac{\alpha}{2}} \frac{\bar{\sigma}}{\sqrt{n}}. \end{aligned}$$

Observația 5.3.24. Dacă există o informație relativă la valoarea medie de forma că aceasta nu este limitată superior, adică are tendința de a avea valori mai mari decât cele așteptate, atunci intervalul numeric $(t_1, t_2) = (-\infty, t_{n-1, 1-\alpha})$, care conduce la intervalul de încredere pentru m :

$$(\bar{m}_1, \bar{m}_2) = \left(\bar{X} - \frac{\bar{\sigma}}{\sqrt{n}} t_{n-1, 1-\alpha}, \infty \right).$$

În mod analog, dacă se cunoaște că valoarea medie nu este limitată inferior, atunci intervalul numeric $(t_1, t_2) = (t_{n-1, \alpha}, \infty)$, care conduce la intervalul de încredere

$$(\bar{m}_1, \bar{m}_2) = \left(-\infty, \bar{X} - \frac{\bar{\sigma}}{\sqrt{n}} t_{n-1, \alpha} \right).$$

Observația 5.3.25. Geometric, dacă se consideră graficele funcției de repartiție și densității de probabilitate a legii Student cu $n-1$ grade de libertate, modul de obținere a cuantilei $t_{n-1, \gamma}$ este prezentat în Figura 5.2.

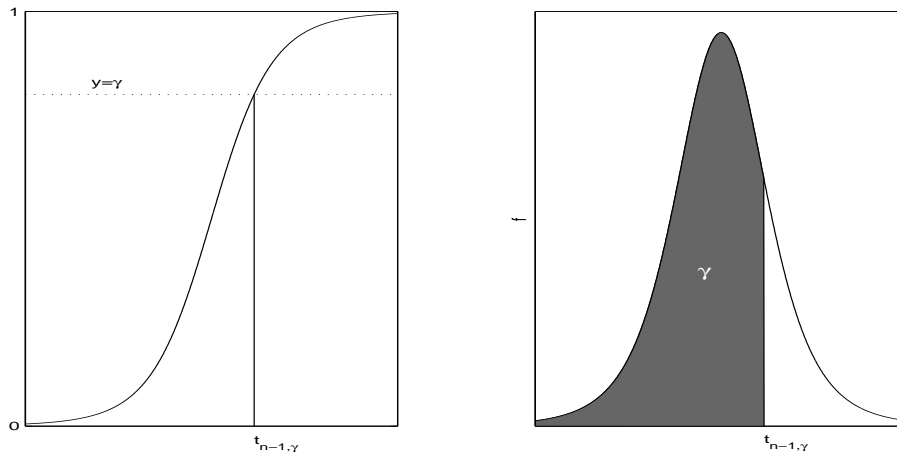


Figura 5.2: Determinarea cuantilei $t_{n-1, \gamma}$

Din teorema limită centrală avem că rezultatele pot fi aplicate pentru o caracteristică X ce urmează o lege de probabilitate oarecare, pentru $n > 30$.

Exemplul 5.3.26. Pentru recepționarea unei mărfi ambalată în cutii, se efectuează un control, prin sondaj, privind greutatea X a cutiilor. Pentru 22 de cutii cântărite s-a obținut distribuția empirică de selecție, relativ la caracteristica X :

$$X \begin{pmatrix} 2.7 & 2.8 & 2.9 & 3.0 & 3.1 & 3.2 & 3.3 \\ 1 & 2 & 5 & 3 & 5 & 4 & 2 \end{pmatrix}.$$

Folosind probabilitatea de încredere 0.98, să determinăm un interval de încredere pentru valoarea medie a greutății cutiilor, presupunând că X urmează legea normală $\mathcal{N}(m, \sigma)$.

Deoarece abaterea standard $\sigma = \sqrt{\text{Var}(X)}$ este necunoscută, se consideră statistica

$$T = \frac{\bar{X} - m}{\frac{\bar{\sigma}}{\sqrt{n}}},$$

care urmează legea Student cu $n - 1$ grade de libertate.

Extremitățile intervalului de încredere pentru valoarea medie teoretică $m = E(X)$ sunt date prin (5.3.2).

Pentru $n - 1 = 21$ și $1 - \alpha = 0.98$ ($\alpha = 0.02$), din *Anexa II* se determină $t_{n-1, 1-\frac{\alpha}{2}} = 2.518$.

De asemenea, folosind datele de selecție, obținem valoarea \bar{x} a mediei de selecție \bar{X} , anume

$$\bar{x} = \frac{1}{22} (1 \cdot 2.7 + 2 \cdot 2.8 + 5 \cdot 2.9 + 3 \cdot 3 + 5 \cdot 3.1 + 4 \cdot 3.2 + 2 \cdot 3.3) = 3.032$$

și valoarea abaterii standard de selecție

$$\bar{\sigma} = \sqrt{\frac{1}{21} \sum_{k=1}^7 f_k (x_k - \bar{x})^2} = \sqrt{\frac{0.587728}{21}} = 0.167.$$

Putem scrie atunci intervalul (numeric) de încredere

$$\begin{aligned} & \left(\bar{x} - t_{n-1, 1-\frac{\alpha}{2}} \frac{\bar{\sigma}}{\sqrt{n}}; \bar{x} + t_{n-1, 1-\frac{\alpha}{2}} \frac{\bar{\sigma}}{\sqrt{n}} \right) \\ &= \left(3.032 - \frac{2.518 \cdot 0.167}{\sqrt{22}}; 3.032 + \frac{2.518 \cdot 0.167}{\sqrt{22}} \right) = (2.942; 3.122). \end{aligned}$$

Programul 5.3.27. Dacă se execută programul Matlab

```
x=[2.7,2.8*ones(1,2),2.9*ones(1,5),3*ones(1,3),...
    3.1*ones(1,5),3.2*ones(1,4),3.3*ones(1,2)];
ma=mean(x); va=var(x);
s=tinv(0.99,21)*sqrt(va)/sqrt(22);
m1=ma-s; m2=ma+s;
fprintf(' (m1,m2)=(%6.3f,%6.3f)',m1,m2)
```

relativ la problema precedentă, se obține intervalul de încredere:

```
(m1,m2)=( 2.942, 3.122)
```

Interval de încredere pentru dispersie

Considerăm caracteristica X ce urmează legea normală $\mathcal{N}(m, \sigma)$ cu $m \in \mathbb{R}$ necunoscut și $\sigma > 0$ necunoscut. Determinăm un interval de încredere pentru dispersia teoretică σ^2 a caracteristicii X .

Statistica ce se consideră în acest caz este

$$\chi^2 = \frac{(n-1)\bar{\sigma}^2}{\sigma^2} = \frac{1}{\sigma^2} \sum_{k=1}^n (X_k - \bar{X})^2,$$

care, conform Proprietății 4.2.25, urmează legea χ^2 cu $n-1$ grade de libertate.

Se determină intervalul numeric (χ_1^2, χ_2^2) astfel încât

$$P(\chi^2 \in (\chi_1^2, \chi_2^2)) = F_{n-1}(\chi_2^2) - F_{n-1}(\chi_1^2) = 1 - \alpha,$$

unde

$$F_m(x) = \frac{1}{2^{\frac{m}{2}} \Gamma(\frac{m}{2})} \int_0^x t^{\frac{m}{2}-1} e^{-\frac{t}{2}} dt, \quad x > 0,$$

este funcția de repartiție a legii χ^2 cu m grade de libertate și este tabelată, pentru anumite valori, în *Anexa III*.

Dacă se alege $\chi_1^2 = \chi_{n-1, \frac{\alpha}{2}}^2$ și $\chi_2^2 = \chi_{n-1, 1-\frac{\alpha}{2}}^2$, adică astfel încât

$$F_{n-1}(\chi_{n-1, \frac{\alpha}{2}}^2) = \frac{\alpha}{2} \quad \text{și} \quad F_{n-1}(\chi_{n-1, 1-\frac{\alpha}{2}}^2) = 1 - \frac{\alpha}{2},$$

se obține $P(\bar{\sigma}^2 < \sigma^2 < \bar{\sigma}^2) = 1 - \alpha$, unde

$$(5.3.3) \quad \begin{aligned} \bar{\sigma}_1^2 &= \bar{\sigma}_1^2(X_1, X_2, \dots, X_n) = \frac{(n-1)\bar{\sigma}^2}{\chi_{n-1, 1-\frac{\alpha}{2}}^2}, \\ \bar{\sigma}_2^2 &= \bar{\sigma}_2^2(X_1, X_2, \dots, X_n) = \frac{(n-1)\bar{\sigma}^2}{\chi_{n-1, \frac{\alpha}{2}}^2}. \end{aligned}$$

Observația 5.3.28. Modul geometric de obținere a cuantilei $\chi_{n-1, \gamma}^2$ este prezentat în Figura 5.3, dacă se consideră graficele funcției de repartiție și densității de probabilitate a legii χ^2 cu $n-1$ grade de libertate.

Exemplul 5.3.29. Fie X caracteristica ce reprezintă timpul de producere a unei reacții chimice, măsurat în secunde. Dacă X urmează legea normală $\mathcal{N}(m, \sigma)$ și având o selecție repetată de volum $n = 11$, cu datele de selecție 4.21, 4.03, 3.99, 4.05, 3.89, 3.98, 4.01, 3.92, 4.23, 3.85, 4.20, vom determina intervalul de încredere pentru dispersia $\sigma^2 = \text{Var}(X)$ și pentru abaterea standard $\sigma = \sqrt{\text{Var}(X)}$, cu probabilitatea de încredere 0.95.

Se consideră statistica

$$\chi^2 = \frac{(n-1)\bar{\sigma}^2}{\sigma^2}, \quad \text{unde} \quad \bar{\sigma}^2 = \frac{1}{n-1} \sum_{k=1}^n (X_k - \bar{X})^2, \quad \bar{X} = \frac{1}{n} \sum_{k=1}^n X_k,$$

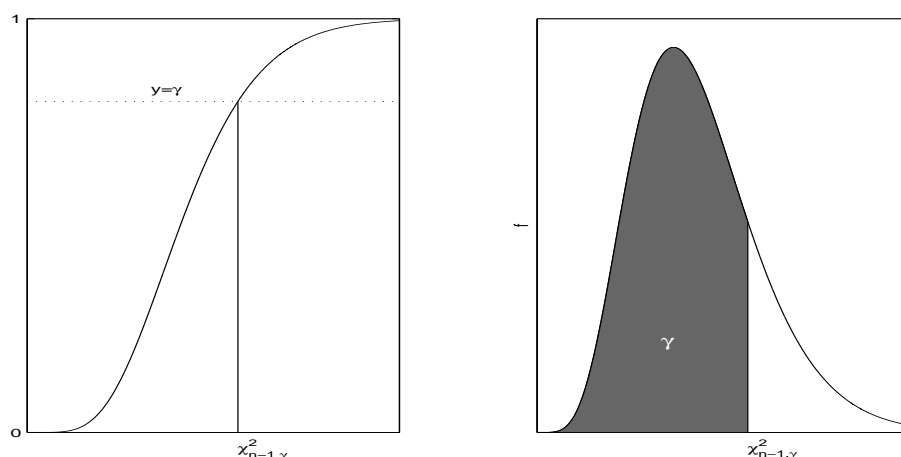


Figura 5.3: Determinarea cuantilei $\chi^2_{n-1, \gamma}$

care urmează legea χ^2 cu $n - 1$ grade de libertate. Extremitățile intervalului de încredere pentru σ^2 vor fi date prin (5.3.3).

Pentru determinarea valorilor numerice ale acestor intervale de încredere, calculăm:

$$\bar{x} = \frac{1}{11} \sum_{k=1}^{11} x_k = \frac{1}{11} (4.21 + 4.03 + \dots + 4.20) = 4.033;$$

$$\bar{\sigma}^2 = \frac{1}{10} \sum_{k=1}^{11} (x_k - \bar{x})^2 = 0.017; \quad \chi^2_{10;0.975} = 20.5 \quad \text{și} \quad \chi^2_{10;0.025} = 3.25.$$

Așadar, intervalele de încredere pentru σ^2 și σ sunt

$$\left(\frac{10 \cdot 0.017}{20.5}; \frac{10 \cdot 0.017}{3.25} \right) = (0.008; 0.052),$$

respectiv

$$\left(\sqrt{\frac{10 \cdot 0.017}{20.5}}; \sqrt{\frac{10 \cdot 0.017}{3.25}} \right) = (0.091; 0.229).$$

Programul 5.3.30. Un program Matlab, care să rezolve această problemă, ar putea fi următorul:

```
x=[4.21,4.03,3.99,4.05,3.89,3.98,4.01,3.92,4.23,3.85,4.20];
va=var(x); c1=chi2inv(0.025,10); c2=chi2inv(0.975,10);
```

```

v1=10*va/c2; v2=10*va/c1; s1=sqrt(v1); s2=sqrt(v2);
fprintf(' (v1,v2)=(%.3f,%.3f)\n',v1,v2)
fprintf(' (s1,s2)=(%.3f,%.3f)',s1,s2)

```

În urma executării, se obțin intervalele de încredere pentru dispersie și respectiv pentru abaterea standard:

```

(v1,v2)=( 0.008, 0.052)
(s1,s2)=( 0.091, 0.229)

```

Interval de încredere pentru raportul dispersiilor

Se consideră caracteristicile independente X' și X'' fiecare urmând legea normală, respectiv $\mathcal{N}(m', \sigma')$ și $\mathcal{N}(m'', \sigma'')$. Relativ la cele două caracteristici se consideră câte o selecție repetată, respectiv de volume n' și n'' . Vom determina un interval de încredere pentru raportul σ'^2 / σ''^2 corespunzător probabilității de încredere $1 - \alpha$ dată.

Pentru aceasta se consideră statistica

$$F = \frac{\bar{\sigma}''^2}{\sigma''^2} \bigg/ \frac{\bar{\sigma}'^2}{\sigma'^2},$$

care, conform Proprietății 4.2.30, urmează legea Fisher–Snedecor cu $m = n'' - 1$ și $n = n' - 1$ grade de libertate.

Se determină intervalul numeric (f_1, f_2) astfel încât

$$P(F \in (f_1, f_2)) = F_{m,n}(f_2) - F_{m,n}(f_1) = 1 - \alpha,$$

unde

$$F_{m,n}(x) = \left(\frac{m}{n}\right)^{\frac{m}{2}} \frac{\Gamma\left(\frac{m+n}{2}\right)}{\Gamma\left(\frac{m}{2}\right)\Gamma\left(\frac{n}{2}\right)} \int_0^x t^{\frac{m}{2}-1} \left(1 + \frac{m}{n}t\right)^{-\frac{m+n}{2}} dt, \quad x > 0,$$

este funcția de repartiție a legii Fisher–Snedecor cu m și n grade de libertate și care este tabelată pentru anumite valori în *Anexa IV*.

Dacă se alege $f_1 = f_{m,n;\frac{\alpha}{2}}$ și $f_2 = f_{m,n;1-\frac{\alpha}{2}}$, adică astfel încât

$$F_{m,n}\left(f_{m,n;\frac{\alpha}{2}}\right) = \frac{\alpha}{2} \quad \text{și} \quad F_{m,n}\left(f_{m,n;1-\frac{\alpha}{2}}\right) = 1 - \frac{\alpha}{2},$$

atunci

$$(5.3.4) \quad P\left(f_{m,n;\frac{\alpha}{2}} \frac{\bar{\sigma}'^2}{\bar{\sigma}''^2} < \frac{\sigma'^2}{\sigma''^2} < f_{m,n;1-\frac{\alpha}{2}} \frac{\bar{\sigma}'^2}{\bar{\sigma}''^2}\right) = 1 - \alpha.$$

Această relație pune în evidență intervalul de încredere pentru raportul celor două dispersii.

Observația 5.3.31. Modul geometric de obținere a cuantilei $f_{m,n;\gamma}$ este prezentat în Figura 5.4, dacă se consideră graficele funcției de repartiție și densității de probabilitate a legii Fisher–Snedecor cu m și n grade de libertate.

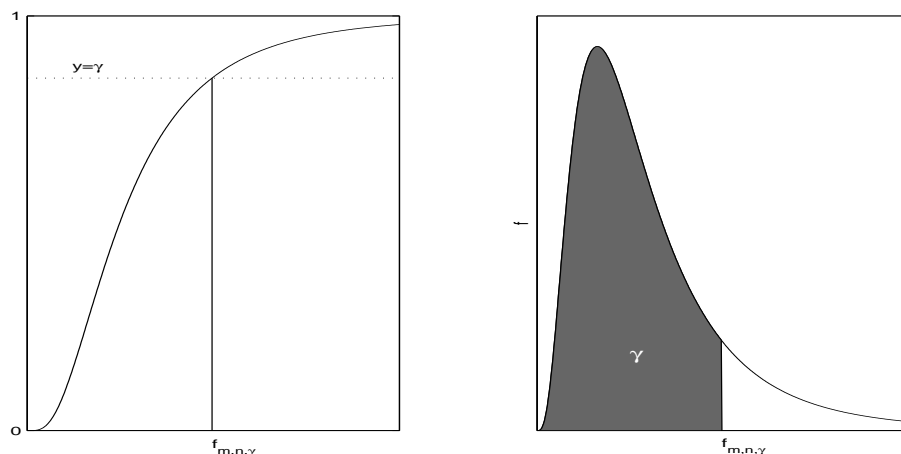


Figura 5.4: Determinarea cuantilei $f_{m,n;\gamma}$

Exemplul 5.3.32. Se cercetează precizia cu care două mașini produc conserve de același tip. Pentru aceasta, se consideră câte un eșantion din conservele produse de cele două mașini și se măsoară greutatea acestora. Fie X' greutatea în grame a unei conserve produsă de prima mașină, respectiv X'' pentru a doua mașină. Măsurătorile obținute sunt:

X' :	1021	980	988	1017	1005	998	1014	985	995	1004	1030
		1015	995	1023	1008	1013					
X'' :	1003	988	993	1013	1006	1002	1014	997	1002	1010	975

Să se calculeze un interval de încredere pentru raportul abaterilor standard, adică pentru $\frac{\sigma'}{\sigma''} = \frac{\sqrt{Var(X')}}{\sqrt{Var(X'')}}$, folosind probabilitatea de risc $\alpha = 0.05$.

Vom considera că cele două caracteristici X' și X'' sunt independente și că urmează legile normale $\mathcal{N}(m', \sigma')$ și respectiv $\mathcal{N}(m'', \sigma'')$.

Statistica ce se utilizează pentru stabilirea intervalului de încredere pentru raportul σ' / σ'' este

$$F = \frac{\bar{\sigma}''^2}{\sigma''^2} \bigg/ \frac{\bar{\sigma}'^2}{\sigma'^2},$$

care urmează legea Fisher–Snedecor cu $m = n'' - 1 = 10$ și $n = n' - 1 = 15$ grade de libertate.

Intervalul de încredere este dat prin (5.3.4).

Pentru aceasta avem că

$$\bar{x}' = \frac{1}{n'} \sum_{k=1}^{n'} x'_k = \frac{1}{16} (1021 + 980 + \dots + 1013) = 1005.7,$$

$$\bar{x}'' = \frac{1}{n''} \sum_{k=1}^{n''} x''_k = \frac{1}{11} (1003 + 988 + \dots + 975) = 1000.3,$$

$$\begin{aligned} \bar{\sigma}'^2 &= \frac{1}{n' - 1} \sum_{k=1}^{n'} (x'_k - \bar{x}')^2 = \frac{1}{15} [(1021 - 1005.7)^2 + \dots + (1013 - 1005.7)^2] \\ &= 210.63, \end{aligned}$$

$$\begin{aligned} \bar{\sigma}''^2 &= \frac{1}{n'' - 1} \sum_{k=1}^{n''} (x''_k - \bar{x}'')^2 = \frac{1}{10} [(1003 - 1000.3)^2 + \dots + (975 - 1000.3)^2] \\ &= 134.42. \end{aligned}$$

Pe de altă parte, se obțin cuantilele

$$f_{m,n;\frac{\alpha}{2}} = f_{10,15;0.025} = 0.28, \quad f_{m,n;1-\frac{\alpha}{2}} = f_{10,15;0.975} = 3.06,$$

pe baza *Anexei IV*. Facem observația că în acest calcul s-a avut în vedere că

$$f_{m,n;\frac{\alpha}{2}} = \frac{1}{f_{n,m;1-\frac{\alpha}{2}}}, \quad \text{adică} \quad f_{10,15;0.025} = \frac{1}{f_{15,10;0.975}} = \frac{1}{3.52} = 0.28.$$

Se ajunge astfel la faptul că

$$\frac{\sigma'}{\sigma''} \in \left(\sqrt{0.28 \cdot \frac{210.63}{134.42}}, \sqrt{3.06 \cdot \frac{210.63}{134.42}} \right) = (0.66, 2.19).$$

La același rezultat se ajunge și dacă se consideră statistica $\frac{1}{F}$.

Programul 5.3.33. Programul Matlab, care urmează, calculează intervalele de încredere pentru raportul dispersiilor, respectiv al abaterilor standard, pentru datele considerate în problema precedentă.

```
x1=[1021,980,988,1017,1005,998,1014,985,995,1004,...
    1030,1015,995,1023,1008,1013];
x2=[1003,988,993,1013,1006,1002,1014,997,1002,1010,975];
v1=var(x1); v2=var(x2); r=v1/v2;
```

```

f1=finv(0.025,10,15); f2=finv(0.975,10,15);
r1=f1*r; r2=f2*r; s1=sqrt(r1); s2=sqrt(r2);
fprintf(' (r1,r2)=(%6.3f,%6.3f)\n',r1,r2)
fprintf(' (s1,s2)=(%6.3f,%6.3f)',s1,s2)

```

În urma executării programului, s-a obținut:

```

(r1,r2)=( 0.445, 4.795)
(s1,s2)=( 0.667, 2.190)

```

Interval de încredere pentru diferența mediilor

Caracteristicile independente X' și X'' urmează respectiv legile normale $\mathcal{N}(m'; \sigma')$ și $\mathcal{N}(m''; \sigma'')$. Folosind câte o selecție repetată de volum n' și n'' pentru cele două caracteristici, vrem să determinăm un interval de încredere pentru diferența $m' - m''$.

Distingem următoarele trei situații:

- A. abaterile standard ale celor două caracteristici sunt cunoscute,
- B. abaterile standard sunt necunoscute, dar se știe că sunt egale,
- C. abaterile standard sunt necunoscute și diferite,

pe care le tratăm pe rând în continuare.

A. *Abaterile standard σ' și σ'' sunt cunoscute.*

Se consideră statistica

$$(5.3.5) \quad Z = \frac{(\bar{X}' - \bar{X}'') - (m' - m'')}{\sqrt{\frac{\sigma'^2}{n'} + \frac{\sigma''^2}{n''}}},$$

care urmează legea normală $\mathcal{N}(0, 1)$ (a se vedea Observația 4.2.29). Astfel, pentru probabilitatea de risc $\alpha \in (0, 1)$ dată, se poate determina intervalul numeric $(z_1, z_2) = \left(-z_{1-\frac{\alpha}{2}}, z_{1-\frac{\alpha}{2}}\right)$ astfel încât $P(z_1 < Z < z_2) = 1 - \alpha$. Anume, $z_{1-\frac{\alpha}{2}}$ se calculează din relația $\Phi\left(z_{1-\frac{\alpha}{2}}\right) = \frac{1-\alpha}{2}$, unde

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

este funcția lui Laplace, tabelată în *Anexa I*. Se ajunge astfel la relația

$$(5.3.6) \quad P(\bar{X}' - \bar{X}'' - \kappa_\alpha < m' - m'' < \bar{X}' - \bar{X}'' + \kappa_\alpha) = 1 - \alpha,$$

unde

$$\kappa_{\alpha} = z_{1-\frac{\alpha}{2}} \sqrt{\frac{\sigma'^2}{n'} + \frac{\sigma''^2}{n''}},$$

care pune în evidență extremitățile intervalului de încredere pentru diferența celor două medii.

B. Abaterile standard σ' și σ'' sunt egale cu σ , dar necunoscut.

Se consideră statistica

$$(5.3.7) \quad T = \frac{(\bar{X}' - \bar{X}'') - (m' - m'')}{\sqrt{(n' - 1) \bar{\sigma}'^2 + (n'' - 1) \bar{\sigma}''^2}} \sqrt{\frac{n' + n'' - 2}{\frac{1}{n'} + \frac{1}{n''}}},$$

unde

$$\bar{\sigma}'^2 = \frac{1}{n' - 1} \sum_{k=1}^{n'} (X'_k - \bar{X}'), \quad \bar{\sigma}''^2 = \frac{1}{n'' - 1} \sum_{k=1}^{n''} (X''_k - \bar{X}'')$$

și care urmează legea Student cu $m = n' + n'' - 2$ grade de libertate (a se vedea Proprietatea 4.2.28).

Ca la punctul precedent se obțin extremitățile intervalului de încredere

$$(5.3.8) \quad \bar{m}_{1,2} = \bar{X}' - \bar{X}'' \pm t_{m, 1-\frac{\alpha}{2}} \sqrt{\frac{\frac{1}{n'} + \frac{1}{n''}}{n' + n'' - 2}} \sqrt{(n' - 1) \bar{\sigma}'^2 + (n'' - 1) \bar{\sigma}''^2},$$

unde $t_{m, 1-\frac{\alpha}{2}}$ este cuantila de ordin $1 - \frac{\alpha}{2}$ pentru legea Student cu m grade de libertate.

C. Abaterile standard σ' și σ'' sunt diferite și necunoscute.

Se consideră statistica

$$(5.3.9) \quad T = \frac{(\bar{X}' - \bar{X}'') - (m' - m'')}{\sqrt{\frac{\bar{\sigma}'^2}{n'} + \frac{\bar{\sigma}''^2}{n''}}},$$

care urmează legea Student cu n grade de libertate și care se calculează prin formula

$$(5.3.10) \quad \frac{1}{n} = \frac{c^2}{n' - 1} + \frac{(1 - c)^2}{n'' - 1}, \quad \text{unde} \quad c = \frac{\bar{\sigma}'^2}{n'} \bigg/ \left(\frac{\bar{\sigma}'^2}{n'} + \frac{\bar{\sigma}''^2}{n''} \right).$$

Ca la punctul precedent se ajunge la extremitățile intervalului de încredere pentru diferența celor două medii:

$$(5.3.11) \quad \bar{m}_{1,2} = \bar{X}' - \bar{X}'' \pm t_{n, 1-\frac{\alpha}{2}} \sqrt{\frac{\bar{\sigma}'^2}{n'} + \frac{\bar{\sigma}''^2}{n''}}.$$

Exemplul 5.3.34. Mașinile M_1 și M_2 ambalează carne în pachete de 1000 grame. Greutatea cutiilor este o caracteristică X' , ce urmează legea normală $\mathcal{N}(m', \sigma')$ și respectiv o caracteristică X'' , ce urmează legea normală $\mathcal{N}(m'', \sigma'')$.

Cântărind 100 de pachete din cele produse de mașina M_1 s-a obținut valoarea mediei de selecție $\bar{x}' = 1007$ grame, iar din cântărirea a 150 de pachete de la mașina M_2 s-a obținut $\bar{x}'' = 1002$ grame.

Folosind probabilitatea de încredere 0.98, să determinăm intervalul de încredere pentru diferența $m' - m''$, dacă se știe că abaterile standard sunt $\sigma' = 3$ și $\sigma'' = 4$.

Se folosește statistica (5.3.5), care urmează legea normală $\mathcal{N}(0, 1)$. Astfel, intervalul de încredere pentru diferența $m' - m''$ este dat prin (5.3.6), unde $z_{1-\frac{\alpha}{2}}$ se determină astfel ca $\Phi\left(z_{1-\frac{\alpha}{2}}\right) = \frac{1-\alpha}{2} = 0.49$. Folosind *Anexa I*, obținem $z_{1-\frac{\alpha}{2}} = 2.33$.

De asemenea, avem că

$$\sqrt{\frac{\sigma'^2}{n'} + \frac{\sigma''^2}{n''}} = \sqrt{\frac{9}{100} + \frac{16}{150}} = \sqrt{\frac{59}{300}} = 0.4435.$$

Astfel, intervalul de încredere pentru diferența $m' - m''$ este

$$\begin{aligned} & \left((\bar{x}' - \bar{x}'') - z_{1-\frac{\alpha}{2}} \sqrt{\frac{\sigma'^2}{n'} + \frac{\sigma''^2}{n''}}; (\bar{x}' - \bar{x}'') + z_{1-\frac{\alpha}{2}} \sqrt{\frac{\sigma'^2}{n'} + \frac{\sigma''^2}{n''}} \right) \\ &= (5 - 2.33 \cdot 0.4435; 5 + 2.33 \cdot 0.4435) = (3.97; 6.03). \end{aligned}$$

Programul 5.3.35. Executarea programului Matlab:

```
z1=norminv(0.01,0,1); z2=norminv(0.99,0,1);
s=sqrt(9/100+16/150); d=1007-1002;
m1=d+z1*s; m2=d+z2*s;
fprintf(' (m1,m2)=(%6.3f,%6.3f)',m1,m2)
```

conduce la rezultatul:

```
(m1,m2)=( 3.968, 6.032)
```

Exemplul 5.3.36. Fie caracteristica X' ce urmează legea normală $\mathcal{N}(m', \sigma)$ și care reprezintă vânzările în milioane lei pe săptămână la magazinele alimentare în orașul \mathcal{A} și X'' vânzările în milioane lei la magazinele alimentare din orașul \mathcal{B} și care urmează legea normală $\mathcal{N}(m'', \sigma)$. S-au efectuat două sondaje, respectiv pentru X' și X'' și s-au obținut următoarele date de selecție:

X' : 226.5, 224.1, 218.6, 220.1, 228.8, 229.6, 222.5;

X'' : 221.5, 230.2, 223.4, 224.3, 230.8, 223.8.

Cu probabilitatea de încredere 0.95, vrem să se construim un interval de încredere pentru diferența $m' - m''$, dacă $\sigma > 0$ este necunoscut.

Folosind statistica (5.3.7), care urmează legea Student cu $n = n' + n'' - 2 = 11$ grade de libertate, se va construi intervalul de încredere pentru $m' - m''$. Anume, acest interval de încredere este precizat prin (5.3.8).

Pentru a determina valoarea numerică a intervalului de încredere, se calculează pe rând

$$\begin{aligned}\bar{x}' &= \frac{1}{7} (226.5 + 224.1 + 218.6 + 220.1 + 228.8 + 229.6 + 222.5) = 224.314, \\ \bar{x}'' &= \frac{1}{6} (221.5 + 230.2 + 223.4 + 224.3 + 230.8 + 223.8) = 225.667, \\ \bar{\sigma}'^2 &= \frac{1}{6} \sum_{k=1}^7 (x'_k - \bar{x}')^2 = 17.765, \quad \bar{\sigma}''^2 = \frac{1}{5} \sum_{k=1}^6 (x''_k - \bar{x}'')^2 = 14.951, \\ s &= \sqrt{\frac{\frac{1}{7} + \frac{1}{6}}{11} (6 \cdot 17.765 + 5 \cdot 14.951)} = 2.259.\end{aligned}$$

De asemenea, din *Anexa II*, pentru $1 - \alpha = 0.95$ și $n = 11$, obținem $t_{n, 1-\frac{\alpha}{2}} = 2.201$, astfel că intervalul de încredere pentru $m' - m''$ va fi

$$\begin{aligned}& \left((\bar{x}' - \bar{x}'') - t_{n, 1-\frac{\alpha}{2}} \cdot s; (\bar{x}' - \bar{x}'') + t_{n, 1-\frac{\alpha}{2}} \cdot s \right) \\ &= (-1.353 - 2.201 \cdot 2.259; -1.353 + 2.201 \cdot 2.259) = (-6.32; 3.62).\end{aligned}$$

Programul 5.3.37. Prin executarea programului Matlab:

```
x1=[226.5,224.1,218.6,220.1,228.8,229.6,222.5];
x2=[221.5,230.2,223.4,224.3,230.8,223.8];
ma1=mean(x1); ma2=mean(x2); md=ma1-ma2;
v1=var(x1); v2=var(x2);
t1=tinv(0.025,11); t2=tinv(0.975,11);
s=sqrt((1/7+1/6)/11)*sqrt(6*v1+5*v2);
m1=md+t1*s; m2=md+t2*s;
fprintf(' (m1,m2)=(%6.3f,%6.3f)\n',m1,m2)
```

regăsim intervalul de încredere pentru diferența valorilor medii:

```
(m1,m2)=(-6.324, 3.619)
```

Exemplul 5.3.38. Se consideră două aparate de îmbuteliat vin în sticle de 750 ml. Fie caracteristica X' ce reprezintă cantitatea de vin (în ml) îmbuteliată de primul aparat într-o sticlă și respectiv X'' aceeași caracteristică pentru al doilea aparat. Pentru compararea modului de îmbuteliere pentru cele două aparate, se consideră câte o selecție din sticlele îmbuteliate de cele două aparate, respectiv

X' : 746 743 748 748 750 745 744 753 750 751 743 747 749 742
 X'' : 749 748 748 753 750 749 747 745 744 751

Vom construi un interval de încredere pentru diferența $m' - m'' = E(X') - E(X'')$, folosind probabilitatea de risc $\alpha = 0.05$. Vom considera că cele două caracteristici X' și X'' sunt independente și că urmează legile normale $\mathcal{N}(m', \sigma')$ și $\mathcal{N}(m'', \sigma'')$.

Distingem două cazuri, unul când se știe că $\sigma' = \sigma''$, caz în care se folosește statistica (5.3.7), care urmează legea Student cu $n = n' + n'' - 2$ grade de libertate, respectiv când se știe că $\sigma' \neq \sigma''$, caz în care se folosește statistica (5.3.9), care urmează legea Student cu n grade de libertate. Numărul n al gradelor de libertate se calculează, în acest caz, cu formula (5.3.10)

În cazul $\sigma' = \sigma''$, extremitățile intervalului de încredere pentru diferența $m' - m''$ sunt date prin formulele

$$\bar{m}_1 = (\bar{x}' - \bar{x}'') - t_{n, 1-\frac{\alpha}{2}} \sqrt{\frac{\frac{1}{n'} + \frac{1}{n''}}{n' + n'' - 2}} \cdot \sqrt{(n' - 1) \bar{\sigma}'^2 + (n'' - 1) \bar{\sigma}''^2},$$

$$\bar{m}_2 = (\bar{x}' - \bar{x}'') + t_{n, 1-\frac{\alpha}{2}} \sqrt{\frac{\frac{1}{n'} + \frac{1}{n''}}{n' + n'' - 2}} \cdot \sqrt{(n' - 1) \bar{\sigma}'^2 + (n'' - 1) \bar{\sigma}''^2}.$$

Cuantila $t_{n, 1-\frac{\alpha}{2}} = t_{22, 0.975} = 2.07$, s-a determinat, conform *Anexei II*, din relația $F_{22}(t_{22, 0.975}) = 1 - \frac{\alpha}{2} = 0.975$.

Pe de altă parte

$$\bar{x}' = \frac{1}{14} (746 + 743 + \dots + 742) = 747.1,$$

$$\bar{x}'' = \frac{1}{10} (749 + 748 + \dots + 751) = 748.4,$$

$$(n' - 1) \bar{\sigma}'^2 = \sum_{k=1}^{n'} (x'_k - \bar{x}')^2 = (746 - 747.1)^2 + \dots + (742 - 747.1)^2 = 146.94,$$

$$(n'' - 1) \bar{\sigma}''^2 = \sum_{k=1}^{n''} (x''_k - \bar{x}'')^2 = (749 - 748.4)^2 + \dots + (751 - 748.4)^2 = 64.4.$$

Astfel avem că

$$\bar{m}_1 = (747.1 - 748.4) - 2.07 \sqrt{\frac{\frac{1}{14} + \frac{1}{10}}{14 + 10 - 2}} \sqrt{146.94 + 64.4} = -4.0,$$

$$\bar{m}_2 = (747.1 - 748.4) + 2.07 \sqrt{\frac{\frac{1}{14} + \frac{1}{10}}{14 + 10 - 2}} \sqrt{146.94 + 64.4} = 1.3.$$

Pentru cazul $\sigma' \neq \sigma''$, calculăm prima dată numărul n al gradelor de libertate. Astfel avem că

$$c = \frac{\bar{\sigma}'^2}{n'} \bigg/ \left(\frac{\bar{\sigma}'^2}{n'} + \frac{\bar{\sigma}''^2}{n''} \right) = \frac{11.30}{14} \bigg/ \left(\frac{11.3}{14} + \frac{7.16}{10} \right) = \frac{0.807}{0.807 + 0.716} = 0.53,$$

iar apoi

$$\frac{1}{n} = \frac{c^2}{n' - 1} + \frac{(1 - c)^2}{n'' - 1} = \frac{0.281}{13} + \frac{0.221}{9} = 0.046,$$

de unde avem că $n = 22$.

Extremitățile intervalului de încredere sunt date prin

$$\begin{aligned} \bar{m}_1 &= (\bar{x}' - \bar{x}'') - t_{n,1-\frac{\alpha}{2}} \sqrt{\frac{\bar{\sigma}'^2}{n'} + \frac{\bar{\sigma}''^2}{n''}}, \\ \bar{m}_2 &= (\bar{x}' - \bar{x}'') + t_{n,1-\frac{\alpha}{2}} \sqrt{\frac{\bar{\sigma}'^2}{n'} + \frac{\bar{\sigma}''^2}{n''}}. \end{aligned}$$

Având în vedere că $t_{22,0.975} = 2.07$ rezultă că

$$\begin{aligned} \bar{m}_1 &= (747.1 - 748.4) - 2.07\sqrt{0.807 - 0.716} = -3.9, \\ \bar{m}_2 &= (747.1 - 748.4) + 2.07\sqrt{0.807 - 0.716} = 1.2. \end{aligned}$$

Programul 5.3.39. Programul Matlab, care urmează, rezolvă problema precedentă, adică, va determina intervale de încredere pentru diferența mediilor, în cele două cazuri, când se consideră dispersii egale și necunoscute, respectiv dispersii diferite necunoscute.

```
x1=[746,743,748,748,750,745,744,753,750,751,743,747,749,742];
x2=[749,748,748,753,750,749,747,745,744,751];
m1=mean(x1); m2=mean(x2); md=m1-m2;
v1=var(x1); v2=var(x2);
t1=tinv(0.025,22); t2=tinv(0.975,22);
s=sqrt((1/14+1/10)/22)*sqrt(13*v1+9*v2);
m1=md+t1*s; m2=md+t2*s;
fprintf(' (m1,m2)=(%6.3f,%6.3f)\n',m1,m2)
c=(v1/14)/(v1/14+v2/10); n=c^2/13+(1-c)^2/9;
n=ceil(1/n); t1=tinv(0.025,n); t2=tinv(0.975,n);
s=sqrt(v1/14+v2/10); m1=md+t1*s; m2=md+t2*s;
fprintf(' (m1,m2)=(%6.3f,%6.3f)\n',m1,m2)
```

După executarea programului, se obțin cele două intervale de încredere:

```
(m1,m2)=(-3.990, 1.333)
(m1,m2)=(-3.888, 1.231)
```

5.3.5 Metoda intervalelor de încredere pentru selecții mari

Fie caracteristica X cercetată, cu legea de probabilitate $f(x; \theta)$, unde $\theta \in A \subset \mathbb{R}$ este un parametru necunoscut. Considerăm o selecție repetată de volum n relativă la caracteristica X , pentru care avem variabilele de selecție X_1, X_2, \dots, X_n .

Proprietatea 5.3.40. Dacă se consideră variabilele aleatoare Y_1, Y_2, \dots, Y_n definite prin relația

$$Y_k = \frac{\partial \ln f(X_k; \theta)}{\partial \theta},$$

și pentru care există dispersia $\text{Var}(Y_k) = d^2 > 0$, atunci statistica

$$Z = \frac{1}{d\sqrt{n}} \sum_{k=1}^n Y_k = \frac{1}{d\sqrt{n}} \sum_{k=1}^n \frac{\partial \ln f(X_k; \theta)}{\partial \theta},$$

pentru $n \rightarrow \infty$, urmează legea normală $\mathcal{N}(0, 1)$.

Demonstrație. Variabilele aleatoare X_k , $k = \overline{1, n}$, fiind independente și identic repartizate, rezultă că și variabilele aleatoare Y_k , $k = \overline{1, n}$, sunt independente și identic repartizate. Conform teoremei limită centrală, avem că

$$Y_{(n)} = \frac{1}{d\sqrt{n}} \sum_{k=1}^n (Y_k - E(Y_k))$$

converge în repartiție la legea normală $\mathcal{N}(0, 1)$.

Deoarece

$$E(Y_k) = E\left(\frac{\partial \ln f(X_k; \theta)}{\partial \theta}\right) = 0,$$

rezultă că $Y_{(n)} = Z$, ceea ce încheie demonstrația. \square

Pentru probabilitatea de încredere $1 - \alpha$ dată se va determina intervalul numeric $\left(-z_{1-\frac{\alpha}{2}}, z_{1-\frac{\alpha}{2}}\right)$ astfel încât $P\left(Z \in \left(-z_{1-\frac{\alpha}{2}}, z_{1-\frac{\alpha}{2}}\right)\right) = 2\Phi\left(z_{1-\frac{\alpha}{2}}\right) = 1 - \alpha$. Ceea ce revine la determinarea cuantilei $z_{1-\frac{\alpha}{2}}$ astfel încât $\Phi\left(z_{1-\frac{\alpha}{2}}\right) = \frac{1-\alpha}{2}$. Funcția Φ a lui Laplace este aici, cea tabelată în Anexa I.

Prin operații algebrice se înlocuiește inegalitatea $|Z| < z_{1-\frac{\alpha}{2}}$, cu dubla inegalitate echivalentă de forma $\bar{\theta}_1(X_1, X_2, \dots, X_n) < \theta < \bar{\theta}_2(X_1, X_2, \dots, X_n)$, care definește intervalul de încredere pentru parametrul θ .

Aplicația 5.3.41. Fie caracteristica X ce ia numai valorile 1 și 0 cu probabilitățile p și respectiv $1 - p$, adică are funcția de frecvență $f(x; p) = p^x (1 - p)^{1-x}$, $x = 0, 1$, unde $p \in (0, 1)$ este un parametru necunoscut.

Vom folosi metoda intervalului de încredere pentru selecții mari în vederea estimării parametrului p . Considerăm o selecție repetată de volum (mare) n și probabilitatea de încredere $1 - \alpha$. Deoarece $\ln f(x; p) = x \ln p + (1 - x) \ln(1 - p)$, avem că

$$\frac{\partial \ln f(x; p)}{\partial p} = \frac{x}{p} - \frac{1 - x}{1 - p} = \frac{x - p}{p(1 - p)},$$

și prin urmare se obține statistica

$$Z = \frac{1}{d\sqrt{n}} \sum_{k=1}^n \frac{X_k - p}{p(1 - p)} = \frac{\sqrt{n}}{p(1 - p)d} (\bar{X} - p).$$

Știind că $\text{Var}(X) = p(1 - p)$ rezultă

$$\begin{aligned} d^2 &= \text{Var}\left(\frac{\partial \ln f(X_k; p)}{\partial p}\right) = \text{Var}\left(\frac{X_k - p}{p(1 - p)}\right) = \frac{1}{p^2(1 - p)^2} \text{Var}(X_k) \\ &= \frac{1}{p^2(1 - p)^2} p(1 - p) = \frac{1}{p(1 - p)}. \end{aligned}$$

Statistica Z devine așadar

$$Z = \frac{\sqrt{n}}{\sqrt{p(1 - p)}} (\bar{X} - p)$$

și care urmează legea normală $\mathcal{N}(0, 1)$, când $n \rightarrow \infty$.

Pentru α dat, se determină $z_{1-\frac{\alpha}{2}} = z$ astfel încât $P(|Z| < z) = 1 - \alpha$.

Putem scrie că $|Z| < z$ este echivalentă cu

$$\frac{n(\bar{X} - p)^2}{p(1 - p)} < z^2 \quad \text{sau} \quad (n + z^2)p^2 - (2n\bar{X} + z^2)p + n\bar{X}^2 < 0.$$

Discriminantul trinomialului este pozitiv, anume $\Delta = z^2[z^2 + 4n\bar{X}(1 - \bar{X})] > 0$, deci inecuația în p are soluția de forma unui interval (\bar{p}_1, \bar{p}_2) și care va reprezenta intervalul de încredere pentru parametrul p .

Astfel extremitățile intervalului de încredere au următoarele expresii:

$$\begin{aligned} \bar{p}_1 &= \frac{\left(2\bar{X} + \frac{z^2}{n}\right) - \sqrt{\frac{z^4}{n^2} + 4\bar{X}\frac{z^2}{n} - 4\bar{X}^2\frac{z^2}{n}}}{2\left(1 + \frac{z^2}{n}\right)}, \\ \bar{p}_2 &= \frac{\left(2\bar{X} + \frac{z^2}{n}\right) + \sqrt{\frac{z^4}{n^2} + 4\bar{X}\frac{z^2}{n} - 4\bar{X}^2\frac{z^2}{n}}}{2\left(1 + \frac{z^2}{n}\right)}. \end{aligned}$$

Aceste formule de calcul sunt complicate. Având în vedere că aceste formule au fost deduse pentru n mare, să zicem $n > 30$, când se obțin rezultate bune, putem face simplificări ale acestor formule. În primul rând putem folosi următoarea scriere asimptotică

$$\frac{2\bar{X} + \frac{z^2}{n}}{2\left(1 + \frac{z^2}{n}\right)} \cong \bar{X},$$

iar apoi

$$\begin{aligned} \sqrt{\frac{z^4 + 4\bar{X}nz^2 - 4\bar{X}^2nz^2}{4(n+z^2)^2}} &= \sqrt{\frac{z^4 + 4\bar{X}nz^2 - 4\bar{X}^2nz^2}{4n^2 + 8nz^2 + 4z^4}} \\ &\cong \sqrt{\frac{\bar{X}nz^2 - \bar{X}^2nz^2}{n^2}} = z\sqrt{\frac{\bar{X}(1-\bar{X})}{n}}. \end{aligned}$$

S-a ajuns, în acest mod, la intervalul de încredere pentru p

$$(\bar{p}_1, \bar{p}_2) = \left(\bar{X} - z_{1-\frac{\alpha}{2}} \sqrt{\frac{\bar{X}(1-\bar{X})}{n}}, \bar{X} + z_{1-\frac{\alpha}{2}} \sqrt{\frac{\bar{X}(1-\bar{X})}{n}} \right).$$

Observația 5.3.42. Dacă se dorește să se determine parametrul p cu o incertitudine $\pm \Delta p$, pentru o probabilitate de încredere $1 - \alpha$, atunci volumul selecției se determină în felul următor. Având în vedere că

$$z_{1-\frac{\alpha}{2}} \sqrt{\frac{\bar{X}(1-\bar{X})}{n}} \leq \frac{z_{1-\frac{\alpha}{2}}}{2\sqrt{n}},$$

condiția $\frac{z_{1-\frac{\alpha}{2}}}{2\sqrt{n}} \leq \Delta p$, echivalentă cu $n \geq \frac{z_{1-\frac{\alpha}{2}}^2}{4(\Delta p)^2}$, conduce la volumul optim al selecției.

Mai jos prezentăm un tabel cu valorile optime ale volumului selecției pentru diferite valori ale nivelului de încredere și ale lui Δp :

$\Delta p \setminus 1 - \alpha$	0.90	0.95	0.98
0.01	6760	9600	13530
0.02	1700	2400	3380
0.05	270	380	540

5.3.6 Funcții Matlab privind estimația

Sistemul Matlab, prin *Statistics toolbox*, dispune de funcții corespunzătoare teoriei estimației, pe care le prezentăm în continuare.

Funcția `histfit`

Funcția `histfit` are drept rezultat reprezentarea pe aceeași figură a histogramei și a densității legii normale $\mathcal{N}(\bar{x}, \bar{\sigma})$.

Lansarea executării funcției se face prin una din formele

```
histfit(x)
histfit(x,nc)
```

unde x este un vector ce conține datele ce urmează a fi prelucrate, iar nc este numărul claselor. Dacă parametrul nc este absent, atunci numărul claselor se consideră ca fiind radicalul volumului datelor, adică radicalul lungimii vectorului x .

Programul 5.3.43. Programul următor generează n numere aleatoare, ce urmează legea normală $\mathcal{N}(\mu, \sigma)$, după care apelează la funcția `histfit`, iar graficul obținut, pentru $n=100$, $\mu = 10$ și $\sigma = 2$, este prezentat în Figura 5.5.

```
n=input('n='); mu=input('mu=');
s=input('sigma='); x=normrnd(mu,s,1,n);
histfit(x), colormap spring
```

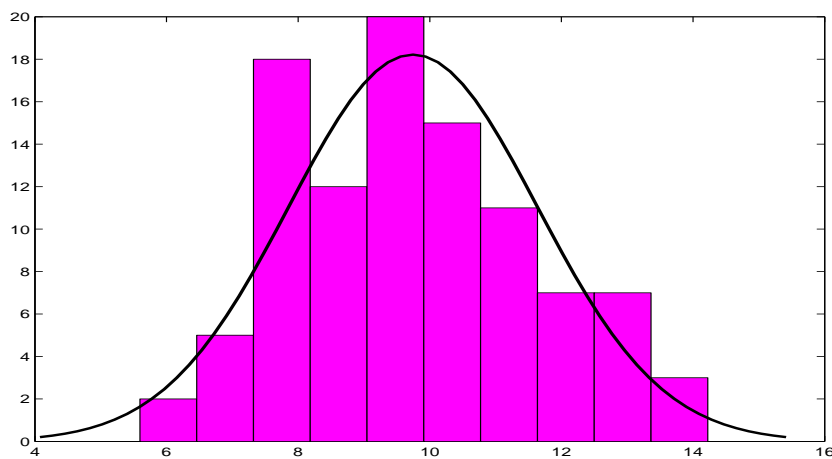


Figura 5.5: Legea $\mathcal{N}(10, 2)$

Funcția `mle`

Pentru estimarea parametrilor, folosind metoda verosimilității maxime, precum și metoda intervalului de încredere, poate fi utilizată funcția `mle`. Apelul funcției se poate face prin una din instrucțiunile:


```

P=mle('lege',x)
[P,int]=mle('lege',x)
[P,int]=mle('lege',x,alpha)
[P,int]=mle('lege',x,alpha,n)

```

unde lege specifică una din legile: bernoulli, unid, bino, poiss, geo, unif, norm, gam, exp, beta, weib și rayl. Parametrul x este un vector, iar în urma executării unei astfel de instrucțiuni, P și int vor conține estimatori de verosimilitate maximă, respectiv intervale de încredere, pentru parametrii legii considerate, obținute pe baza datelor conținute de x .

Parametrul opțional $alpha$, având valoarea implicită $alpha=0.05$, specifică probabilitate de încredere, $1-alpha$, în construirea intervalelor de încredere.

Ultima formă este specifică numai legii binomiale, adică $lege=bino$, parametrul n fiind parametrul cel de la legea $\mathcal{B}(n, p)$. Dacă x și n sunt scalari, se returnează în P raportul acestora, dacă x este vector și n este scalar, se va returna în P fiecare element a lui x împărțit la n , iar dacă x și n sunt vectori de aceleași dimensiuni, P va conține rapoartele pe componente ale lui x și n .

Mai remarcăm că pentru prima dată întâlnim legea lui Bernoulli, care reprezintă cazul particular al legii binomiale, $\mathcal{B}(1, p)$. Dacă $lege=bernoulli$, atunci x , desigur, trebuie să conțină numai valorile 0 și 1.

Funcția 5.3.44. Vom scrie o funcție, care generează n numere aleatoare ce urmează una din legile de probabilitate: uniformă discretă, Poisson, uniformă, normală, Gamma, exponențială, Weibull și Rayleigh. Se va apela această funcție pentru obținerea estimațiilor de verosimilitate maximă ale parametrilor legii de probabilitate considerate. De asemenea, se va reprezenta grafic, pe aceeași figură, funcția de repartiție a legii de probabilitate considerată, împreună cu funcția de repartiție a aceleași legi, dar având parametrii precizați înlocuiți cu estimațiile acestora.

```

function veros(lege,n)
switch lege
case 'unid'
    P1=input('N=');
case 'poiss'
    P1=input('lambda=');
case 'unif'
    P1(1)=input('a='); P1(2)=input('b=');
case 'norm'
    P1(1)=input('mu='); P1(2)=input('sigma=');
case 'gam'
    P1(1)=input('a='); P1(2)=input('b=');
case 'exp'
    P1=input('mu=');
case 'beta'
    P1(1)=input('a='); P1(2)=input('b=');
case 'weib'
    P1(1)=input('a='); P1(2)=input('b=');

```

```

case 'rayl'
    P1=input('b=');
otherwise
    error('Eroare')
end
if length(P1)==1
    x=random(lege,P1,1,n);
else
    x=random(lege,P1(1),P1(2),1,n);
end
P2=mle(lege,x); xx=min(x):0.01:max(x);
if length(P1)==1
    F1=cdf(lege,xx,P1); F2=cdf(lege,xx,P2);
else
    F1=cdf(lege,xx,P1(1),P1(2)); F2=cdf(lege,xx,P2(1),P2(2));
end
plot(xx,F1,'k-.',xx,F2,'k-')
legend('Functia de repartitie teoretica',...
       'Functia de repartitie estimata',2)

```

Apelul funcției veros se face prin

```
>>veros('lege',n)
```

unde valorile pentru n și $lege$, fie că sunt precizate în acest apel, fie sunt precizate înainte. De exemplu, comanda

```
>>veros('poiss',50)
```

are ca efect apelul funcției pentru legea lui Poisson, iar pe ecran se va cere introducerea parametrului λ , după care pe ecran va fi reprezentat graficul din Figura 5.6, în cazul în care $\lambda=7$. Să remarcăm totuși că la un nou apel, cu aceeași parametri, graficul diferă, deoarece sunt generate alte numere aleatoare.

Dacă se consideră comanda

```
>>veros('norm',50)
```

are ca efect apelul funcției pentru legea normală, iar pe ecran se va cere introducerea parametrilor μ și σ , după care pe ecran va fi reprezentat graficul din Figura 5.7, în cazul în care $\mu=10$ și $\sigma=2$.

Funcțiile normlike, gamalike, betalike și weiblike

Funcțiile normlike, gamlike, betalike și weiblike sunt specifice pentru *Statistics toolbox* și se pot apela cu una din comenzile

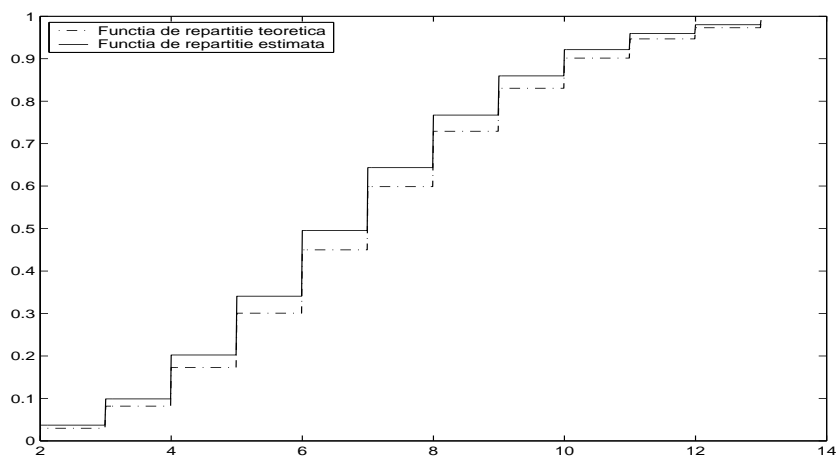
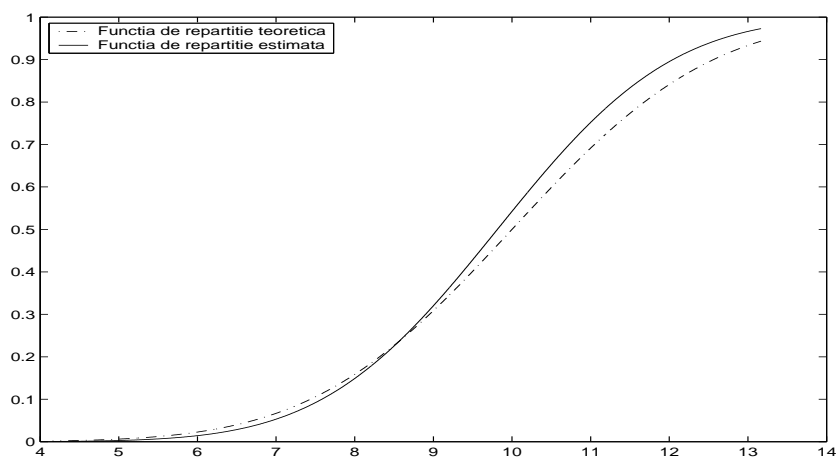
```

L=numef(par,x)
[L,cov]=numef(par,x)

```

unde numef reprezintă numele uneia din cele patru funcții.

Parametrul par conține respectiv parametrii legilor normală, gamma, beta și Weibull, iar x este un vector, care conține datele de selecție.

Figura 5.6: Legea $\mathcal{P}o(7)$ Figura 5.7: Legea $\mathcal{N}(10, 2)$

Prima formă returnează în L valoarea logaritmului funcției de verosimilitate de parametrul precizat prin `par`, înmulțită cu -1 .

Dacă se folosește a doua formă, se mai returnează matricea `var`, care reprezintă valoarea asimptotică a matricei covarianțelor estimatorilor parametrilor legii considerate, când s-au considerat parametri de intrare `par`, ca fiind estimațiile de verosimilitate maximă ale acestora. Remarcăm faptul că `var` reprezintă inversa matricei informației lui Fisher, când volumul selecției tinde la infinit.

Funcțiile `binofit`, `poissfit`, `unifit`, `normfit`, `gamfit`, `expfit`, `betafit`, `raylfit` și `weibfit`

Funcțiile Matlab din această categorie returnează estimații punctuale și intervale de încredere pentru parametrii respectiv ai legilor de probabilitate binomială, Poisson, uniformă, normală, Gamma, exponențială, Beta, Rayleigh și Weibull.

Apelul acestor funcții nu are o sintaxă generală.

Pentru funcția `binofit` se pot folosi comenzile

```
P=binofit(x,n)
[P,int]=binofit(x,n)
[P,int]=binofit(x,n,alpha)
```

parametrul x fiind un vector, iar n precizează parametrul întreg pozitiv al legii binomiale $\mathcal{B}(n, p)$. Pentru funcțiile `poiss` și `exp`, respectiv pentru `gam`, `beta`, `rayl` și `weib`, apelurile sunt de forma:

```
P=numef(x)
[P,int]=numef(x,alpha)
```

unde parametrul `numef` reprezintă numele funcției, iar x este fie un vector, fie o matrice, pentru prima grupă de funcții, iar pentru a doua, numai vector. Dacă x este matrice, atunci funcția operează pentru fiecare coloană în parte.

Pentru legile `unif` și `norm`, apelurile sunt:

```
[par1,par2,int1,int2]=unifit(x)
[par1,par2,int1,int2]=normfit(x,alpha)
```

unde x poate fi vector sau matrice. În al doilea caz, operațiunea se execută pentru fiecare coloană în parte.

De regulă, estimațiile punctuale returnate, sunt obținute prin metoda verosimilității maxime, folosindu-se datele conținute în parametrul x , în parametrii P , `par1` și `par2`. Intervalele de încredere sunt returnate în parametrii `int`, `int1` și `int2`, și sunt determinate folosind probabilitatea de încredere $1-\alpha$. Valoarea implicită pentru α fiind $\alpha=0.05$.

Pentru legile `unif` și `norm`, în comenzile de apel pot lipsi `int2`, `int1` și `par2`, caz în care sunt returnate numai valorile parametrilor rămași.

Mai remarcăm faptul că `int`, `int1` și `int2` au câte două componente, extremitățile intervalului de încredere corespunzător.

Funcția 5.3.45. Vom scrie o funcție, care generează n numere aleatoare, ce urmează una din legile Gamma, Beta, Weibull și Rayleigh. Folosind aceste numere, se vor da estimări punctuale pentru parametrii acestora, precum și intervale de încredere, considerând o probabilitate de încredere $1 - \alpha$ specificată. Folosind estimațiile punctuale obținute pentru parametri, se vor reprezenta grafic densitățile de probabilitate, folosind parametrii estimați și respectiv parametrii considerați la generarea numerelor aleatoare. De asemenea, se vor afișa intervalele de încredere obținute.

```
function estimat(lege,n,alpha)
switch lege
case 'gam'
    a=input('a='); b=input('b=');
    x=gamrnd(a,b,n,1); [P,I]=gamfit(x,alpha);
case 'beta'
    a=input('a='); b=input('b=');
    x=betarnd(a,b,n,1); [P,I]=betafit(x,alpha);
case 'weib'
    a=input('a='); b=input('b=');
    x=weibrnd(a,b,n,1); [P,I]=weibfit(x,alpha);
case 'rayl'
    b=input('b='); x=raylrnd(b,n,1);
    [P,I]=raylfit(x,alpha);
otherwise
    error('Eroare')
end
t=min(x):0.01:max(x);
if length(P)==1
    y=pdf(lege,t,P); yy=pdf(lege,t,b);
    fprintf('  bbarat=%6.3f\n',P)
    fprintf('  (b1,b2)=(%6.3f,%6.3f)\n',I)
else
    y=pdf(lege,t,P(1),P(2)); yy=pdf(lege,t,a,b);
    fprintf('  abarat=%6.3f,  bbarat=%6.3f\n',P)
    fprintf('  (a1,a2)=(%6.3f,%6.3f)\n',I(1,1),I(2,1))
    fprintf('  (b1,b2)=(%6.3f,%6.3f)\n',I(1,2),I(2,2))
end
plot(t,y,'k-',t,yy,'k-.'), plot(t,y,'k-',t,yy,'k-.')
legend('Legea estimata', 'Legea teoretica')
```

Apelul funcției estimat se face prin

```
>>estimat('lege',n,alpha)
```

unde valorile pentru n , lege și α , fie că sunt precizate în acest apel, fie sunt precizate înainte. De exemplu, comanda

```
>>estimat('beta',500,0.01)
```

are ca efect apelul funcției pentru legea Beta, iar pe ecran se va cere introducerea parametrilor a și b , după care pe ecran va fi reprezentat graficul din Figura 5.8, în cazul în care $a=3$ și $b=7$. Să remarcăm totuși că la un nou apel, cu aceiași parametri, graficul diferă, deoarece sunt generate alte numere aleatoare. De asemenea, se vor afișa și estimațiile punctuale, respectiv intervalele de încredere:

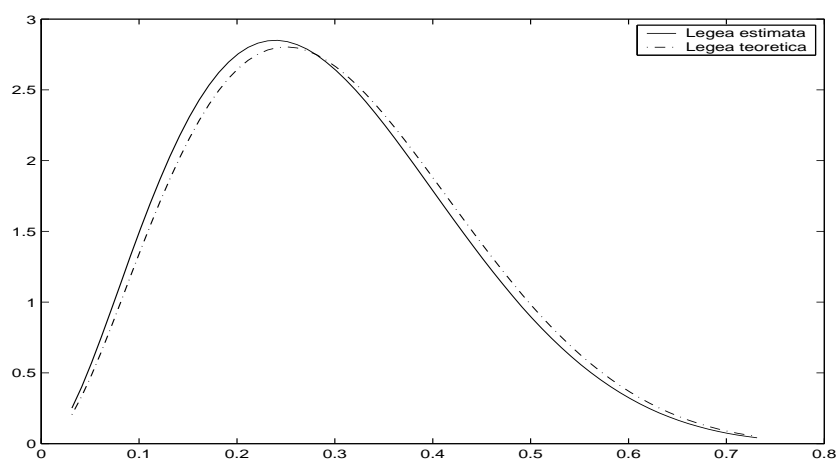


Figura 5.8: Legea $Beta(3, 7)$

```
abarat= 2.925, bbarat= 7.130  
(a1,a2)=( 2.431, 3.418)  
(b1,b2)=( 5.928, 8.333)
```

Capitolul 6

Verificarea ipotezelor statistice

6.1 Concepte de bază

Fie colectivitatea \mathcal{C} cercetată din punct de vedere al caracteristicii X , care are legea de probabilitate dată prin funcția de probabilitate $f(x; \theta)$, ce reprezintă funcția de frecvență în cazul discret, respectiv densitatea de probabilitate în cazul continuu, pentru variabila aleatoare X .

Definiția 6.1.1. Numim ipoteză statistică o presupunere relativă la legea de probabilitate pe care o urmează caracteristica X .

Definiția 6.1.2. Metoda de stabilire a veridicității unei ipoteze statistice, se numește test (criteriu) de verificare a ipotezei statistice.

Definiția 6.1.3. Când ipoteza statistică se referă la parametri de care depinde legea de probabilitate a caracteristicii X se obține un test parametric, iar în caz contrar se obține un test neparametric.

Observația 6.1.4. Pentru testele parametrice vom considera că $\theta \in A = A_0 \cup A_1$, unde $A_0 \cap A_1 = \emptyset$. Ipoteza $H_0 : \theta \in A_0$ o vom numi ipoteză nulă, iar ipoteza $H_1 : \theta \in A_1$ o vom numi ipoteză alternativă.

Definiția 6.1.5. O ipoteză parametrică se numește ipoteză simplă, dacă mulțimea la care se presupune că aparține parametrul necunoscut este formată dintr-un singur element, iar în caz contrar se numește ipoteză compusă.

Observația 6.1.6. Ipoteza nulă este aceea care o intuim a fi cea apropiată de realitate, intuiție pe care o obținem din informațiile pe care le deținem referitoare la caracteristica cercetată X .

Observația 6.1.7. Construirea unui test revine la obținerea unui domeniu $\mathcal{U} \subset \mathbb{R}^n$, numit *regiune critică*, pentru un *nivel de semnificație (probabilitate de risc)* α dat, astfel încât

$$P((X_1, X_2, \dots, X_n) \in \mathcal{U} \mid H_0) = \alpha,$$

unde X_1, X_2, \dots, X_n sunt variabilele de selecție corespunzătoare unei selecții de volum n considerată.

Folosind datele de selecție și regiunea critică \mathcal{U} astfel determinată, avem că ipoteza nulă H_0 va fi admisă (acceptată) dacă $(x_1, x_2, \dots, x_n) \notin \mathcal{U}$, iar în caz contrar va fi respinsă, altfel spus ipoteza alternativă H_1 va fi admisă (acceptată) în al doilea caz.

6.2 Testul Z privind media teoretică

Se consideră caracteristica X care urmează legea normală $\mathcal{N}(m, \sigma)$, unde $m \in \mathbb{R}$ este necunoscut, iar $\sigma > 0$ este cunoscut.

Relativ la media teoretică $m = E(X)$ facem ipoteza nulă $H_0 : m = m_0$, cu una din alternativele:

$$H_1 : m \neq m_0 \quad (\text{testul } Z \text{ bilateral}),$$

$$H_1 : m > m_0 \quad (\text{testul } Z \text{ unilateral dreapta}),$$

$$H_1 : m < m_0 \quad (\text{testul } Z \text{ unilateral stânga}).$$

Pentru verificare ipotezei nule H_0 cu una din alternativele precizate mai înainte, considerăm o selecție repetată de volum n și un nivel de semnificație $\alpha \in (0, 1)$.

Se cunoaște că statistica

$$Z = \frac{\bar{X} - m}{\frac{\sigma}{\sqrt{n}}},$$

urmează legea normală $\mathcal{N}(0, 1)$. Prin urmare, pentru $\alpha \in (0, 1)$, putem determina un interval numeric (z_1, z_2) astfel încât

$$P(z_1 < Z < z_2 \mid H_0) = \Phi(z_2) - \Phi(z_1) = 1 - \alpha.$$

Intervalul (z_1, z_2) nu este determinat în mod unic, dar având în vedere alternativa H_1 considerată, adăugăm condiția suplimentară:

$$z_1 = -z_2, \text{ dacă } H_1 : m \neq m_0, \text{ adică } \Phi\left(z_{1-\frac{\alpha}{2}}\right) = \frac{1-\alpha}{2}, \text{ unde } z_{1-\frac{\alpha}{2}} = z_2;$$

$$z_1 = -\infty, z_2 = z_{1-\alpha}, \text{ unde } \Phi(z_{1-\alpha}) = \frac{1}{2} - \alpha, \text{ dacă } H_1 : m > m_0;$$

$z_1 = z_\alpha$, $z_2 = +\infty$, unde $\Phi(-z_\alpha) = \frac{1}{2} - \alpha$, dacă $H_1 : m < m_0$.

Corespunzător celor trei alternative definim regiunea critică respectiv prin:

$$\begin{aligned}\mathcal{U} &= \left\{ (u_1, u_2, \dots, u_n) \in \mathbb{R}^n \left| \frac{|\bar{u} - m_0|}{\frac{\sigma}{\sqrt{n}}} \geq z_{1-\frac{\alpha}{2}} \right. \right\}, \\ \mathcal{U} &= \left\{ (u_1, u_2, \dots, u_n) \in \mathbb{R}^n \left| \frac{\bar{u} - m_0}{\frac{\sigma}{\sqrt{n}}} \geq z_{1-\alpha} \right. \right\}, \\ \mathcal{U} &= \left\{ (u_1, u_2, \dots, u_n) \in \mathbb{R}^n \left| \frac{\bar{u} - m_0}{\frac{\sigma}{\sqrt{n}}} \leq z_\alpha \right. \right\}.\end{aligned}$$

Aici s-a folosit notația $\bar{u} = \frac{1}{n} \sum_{k=1}^n u_k$.

Într-adevăr, pentru fiecare din cele trei regiuni critice scrise mai înainte, avem că $P((X_1, X_2, \dots, X_n) \in \mathcal{U} \mid H_0) = \alpha$. Modul în care am ales regiunea critică ne conduce respectiv la testul Z bilateral, unilateral dreapta și unilateral stânga.

Mai târziu va fi fundamentată riguros modalitatea de construire a regiunii critice pentru alternativele considerate.

Odată construită regiunea critică \mathcal{U} , ipoteza nulă va fi admisă dacă datele de selecție satisfac condiția $(x_1, x_2, \dots, x_n) \notin \mathcal{U}$, iar în caz contrar va fi respinsă.

Remarcăm de asemenea că regiunea critică \mathcal{U} corespunde mulțimii complementare intervalului (z_1, z_2) .

Etapele aplicării testului Z

1. Se dau: α ; x_1, x_2, \dots, x_n ; m_0 ; σ ;
2. Se calculează intervalul (z_1, z_2) astfel încât $\Phi(z_2) - \Phi(z_1) = 1 - \alpha$ (după cum s-a precizat mai înainte);
3. Se calculează $z = \frac{\bar{x} - m_0}{\frac{\sigma}{\sqrt{n}}}$, unde $\bar{x} = \frac{1}{n} \sum_{k=1}^n x_k$;
4. Concluzia: dacă $z \in (z_1, z_2)$ ipoteza H_0 este admisă, în caz contrar ipoteza este respinsă.

Observația 6.2.1. Testul Z se poate aplica și în cazul unei caracteristici X care nu urmează legea normală, dacă volumul selecției este mare ($n > 30$), considerându-se media teoretică $m = E(X)$ necunoscută și abaterea standard $\sigma = \sqrt{Var(X)}$ cunoscută. Această considerație se bazează pe Teorema 2.7.8.

Exemplul 6.2.2. Caracteristica X reprezintă cheltuielile lunare în mii lei pentru abonamentele la ziare și reviste ale unei familii. Să se verifice, cu nivelul de semnificație $\alpha = 0.01$, dacă media acestor cheltuieli lunare pentru o familie este de 16 mii lei, știind că abaterea standard $\sigma = 3$ mii lei și având o selecție repetată de volum $n = 40$, care ne dă distribuția empirică de selecție

$$X \begin{pmatrix} 11 & 13 & 15 & 17 & 20 \\ 4 & 6 & 12 & 10 & 8 \end{pmatrix}.$$

Deoarece $n = 40 > 30$ și abaterea standard $\sigma = 3$ este cunoscută, vom folosi testul Z pentru verificarea ipotezei nule

$$H_0 : m = E(X) = 16, \quad \text{cu ipoteza alternativă } H_1 : m \neq 16.$$

Pentru $\alpha = 0.01$, folosind *Anexa I*, se determină $z_{1-\frac{\alpha}{2}} = z_{0.995}$, astfel încât

$$\Phi(z_{0.995}) = \frac{1-\alpha}{2} = 0.495.$$

Anume, se obține că $z_{0.995} = 2.58$, care ne dă intervalul numeric $(-2.58; 2.58)$ pentru statistica notată prin:

$$Z = \frac{\bar{X} - m}{\frac{\sigma}{\sqrt{n}}}.$$

Calculăm succesiv

$$\begin{aligned} \bar{x} &= \frac{1}{n} \sum_{k=1}^n x_k = \frac{1}{40} (4 \cdot 11 + 6 \cdot 13 + 12 \cdot 15 + 10 \cdot 17 + 8 \cdot 20) = 15.8; \\ z &= \frac{\bar{x} - m_0}{\frac{\sigma}{\sqrt{n}}} = \frac{15.8 - 16}{\frac{3}{\sqrt{40}}} = \frac{-0.2 \cdot \sqrt{40}}{3} = -0.422. \end{aligned}$$

Deoarece $z = -0.422 \in (-2.58; 2.58)$, rezultă că se acceptă ipoteza că cheltuielile medii lunare ale unei familii pentru abonamentele la ziare și reviste sunt de 16 mii lei, cu nivelul de semnificație 0.01.

Programul 6.2.3. Programul Matlab, care urmează, execută aceste calcule și afișează intervalul de încredere, împreună cu valoarea calculată a statisticii Z .

```
x=[11 13 15 17 20]; f=[4 6 12 10 8];
z=(sum(x.*f)/40-16)/(3/sqrt(40));
z1=norminv(0.005,0,1); z2=norminv(0.995,0,1);
fprintf(' (z1,z2)=(%6.3f,%6.3f)\n',z1,z2)
fprintf(' z=%6.3f,%6.3f)',z)
```

În urma executării programului, s-au obținut rezultatele:

$$\begin{aligned}(z_1, z_2) &= (-2.576, 2.576) \\ z &= -0.422,\end{aligned}$$

Observația 6.2.4. Având în vedere că funcția Φ a lui Laplace este monoton crescătoare și că pentru determinarea intervalului (z_1, z_2) este necesară inversarea funcției lui Laplace, se poate renunța la inversarea acesteia. Pentru aceasta se determină ceea ce se numește *valoare critică* și pe care o notăm prin c . Vom analiza pe rând cele trei alternative.

Dacă se consideră testul bilateral, făcând notația

$$1 - \frac{c}{2} = \Phi(|z|), \quad \text{adică} \quad c = 1 - 2\Phi(|z|),$$

ipoteza nulă va fi respinsă dacă $1 - \frac{c}{2} \geq 1 - \frac{\alpha}{2}$, adică $c \leq \alpha$.

Pentru testul unilateral dreapta, se notează

$$1 - c = \frac{1}{2} + \Phi(z), \quad \text{adică} \quad c = \frac{1}{2} - \Phi(z),$$

iar ipoteza nulă este respinsă, dacă $1 - c \geq 1 - \alpha$, adică, la fel ca mai înainte, dacă $c \leq \alpha$.

Un raționament analog, pentru testul unilateral stânga ne conduce la valoarea critică, calculată după formulă $c = \frac{1}{2} + \Phi(z)$. Drept urmare, ipoteza nulă este respinsă, dacă $c \leq \alpha$.

În rezumat, ipoteza nulă este respinsă, dacă $c \leq \alpha$, unde

$$c = \begin{cases} 1 - 2\Phi(|z|), & \text{dacă } H_1 : m \neq m_0, \\ \frac{1}{2} - \Phi(z), & \text{dacă } H_1 : m > m_0, \\ \frac{1}{2} + \Phi(z), & \text{dacă } H_1 : m < m_0. \end{cases}$$

Pentru datele din Exemplul 6.2.2, se obține

$$c = \begin{cases} 0.673, & \text{pentru testul bilateral,} \\ 0.663, & \text{pentru testul unilateral dreapta,} \\ 0.337, & \text{pentru testul unilateral stânga.} \end{cases}$$

6.2.1 Funcțiile `zscore` și `ztest`

Sistemul Matlab, prin *Statistics toolbox*, dispune de funcțiile `zscore` și `ztest`, cu aplicabilitate la testul Z . Apelarea acestor funcții se face prin:

```
z=zscore(d)
h=ztest(x,m0,sigma)
h=ztest(x,m0,sigma,alpha)
h=ztest(x,m0,sigma,alpha,tail)
[h,c,ci,zc]=ztest(x,m0,sigma,alpha,tail)
```

În urma executării primei instrucțiuni, dacă d este un vector, se calculează abaterile componentelor vectorului de la valoarea medie a componentelor, raportate la abaterea standard a componentelor, adică $x_i = \frac{d_i - \bar{d}}{\bar{\sigma}_d}$. Dacă d este matrice, atunci operația este executată pentru fiecare coloană în parte.

Celelalte instrucțiuni se referă la funcția `ztest` și efectuează testul Z asupra datelor conținute în vectorul x , folosind nivelul de semnificație α , care are valoarea implicită $\alpha=0.05$.

Parametrul `tail` specifică una din cele trei alternative, care conduc la testul bilateral (`tail=0`, implicit), unilateral dreapta (`tail=1`) și unilateral stânga (`tail=-1`). Dacă $h=1$, atunci ipoteza nulă va fi respinsă, respectiv dacă $h=0$, ipoteza nu poate fi respinsă.

Ultima formă de apel permite, de asemenea, obținerea valorii critice c , a valorii zc a statisticii Z , precum și a intervalului de încredere pentru media teoretică, corespunzător probabilității de încredere $1-\alpha$, obținut în vectorul cu două componente ci .

Programul 6.2.5. Programul Matlab, ce urmează, rezolvă problema din Exemplul 6.2.2. Mai mult, se consideră și testele unilateral dreapta și unilateral stânga, precum și obținerea intervalelor de încredere.

```
x=[11*ones(1,4),13*ones(1,6),15*ones(1,12),...
  17*ones(1,10),20*ones(1,8)];
fprintf(' h      c              ci              z\n')
fprintf('_____ \n')
for i=-1:1
    [h,c,ci,zc]=ztest(x,16,3,0.01,i);
    fprintf(' %d  %5.4f  (%4.2f,%4.2f)  %6.4f\n',...
            h,c,ci,zc)
end
```

În urma executării programului, se obțin rezultatele:

h	c	ci	z
0	0.3366	(-Inf,16.90)	-0.4216
0	0.6733	(14.58,17.02)	-0.4216
0	0.6634	(14.70, Inf)	-0.4216

Se observă că pentru $\text{tail} = -1$ și $\text{tail} = 1$, se construiesc intervale nemărginite, respectiv la stânga și la dreapta.

6.3 Puterea unui test

Definiția 6.3.1. Dacă se consideră un test relativ la ipoteza nulă H_0 cu alternativa H_1 , se numește eroare de genul (speța) întâi respingerea unei ipoteze adevărate, iar probabilitatea acestei erori se numește risc de speța întâi (risc al furnizorului) și este dată de nivelul α de semnificație, adică

$$\alpha = P\left((X_1, X_2, \dots, X_n) \in \mathcal{U} \mid H_0\right).$$

Definiția 6.3.2. Se numește eroare de genul (speța) al doilea admiterea unei ipoteze false, iar probabilitatea acestei erori se numește risc de speța a doua (risc al beneficiarului) și este notată β ,

$$\beta = P\left((X_1, X_2, \dots, X_n) \notin \mathcal{U} \mid H_1\right).$$

Observația 6.3.3. Producerea unei erori de genul al doilea este mai gravă decât a unei erori de genul întâi, când verificăm concentrația unui medicament, care peste un anumit grad de concentrație devine dăunător. Pe de altă parte, producerea unei erori de genul întâi este mai gravă decât cea a unei erori de genul doi, când verificăm calitatea unui articol de îmbrăcăminte.

Definiția 6.3.4. Numim puterea unui test probabilitatea respingerii unei ipoteze false, adică

$$\pi(\tilde{\theta}) = \pi(\mathcal{U}; \tilde{\theta}) = P\left((X_1, X_2, \dots, X_n) \in \mathcal{U} \mid \theta = \tilde{\theta}\right),$$

când θ este parametrul asupra căruia se face ipoteza statistică, iar \mathcal{U} este regiunea critică construită sub ipoteza nulă cu nivelul de semnificație $\alpha \in (0, 1)$ fixat.

Observația 6.3.5. Dacă testul considerat se referă la ipoteza nulă $H_0 : \theta = \theta_0$ cu ipoteza alternativă $H_1 : \theta = \theta_1$, atunci $\pi(\theta_0) = \alpha$ și $\pi(\theta_1) = 1 - \beta$.

Aplicația 6.3.6. Calculăm în continuare puterea testului Z , când avem ipoteza alternativă $H_1 : m = m_1 \neq m_0$.

După cum am văzut la testul Z regiunea critică în acest caz este

$$\mathcal{U} = \left\{ (u_1, u_2, \dots, u_n) \in \mathbb{R}^n \mid \frac{|\bar{u} - m_0|}{\frac{\sigma}{\sqrt{n}}} \geq z_{1-\frac{\alpha}{2}} \right\}.$$

Pentru calculul puterii $\pi(m_1)$ a testului avem

$$\pi(m_1) = P((X_1, X_2, \dots, X_n) \in \mathcal{U} \mid H_1) = P\left(\left|\frac{\bar{X} - m_0}{\frac{\sigma}{\sqrt{n}}}\right| \geq z_{1-\frac{\alpha}{2}} \mid H_1\right).$$

Dacă se consideră evenimentul contrar, vom putea scrie succesiv

$$\begin{aligned} \beta = 1 - \pi(m_1) &= P\left(\left|\frac{\bar{X} - m_0}{\frac{\sigma}{\sqrt{n}}}\right| < z_{1-\frac{\alpha}{2}} \mid H_1\right) \\ &= P\left(-z_{1-\frac{\alpha}{2}} < \frac{\bar{X} - m_1}{\frac{\sigma}{\sqrt{n}}} - \frac{m_0 - m_1}{\frac{\sigma}{\sqrt{n}}} < z_{1-\frac{\alpha}{2}} \mid H_1\right) \\ &= P\left(\frac{m_0 - m_1}{\frac{\sigma}{\sqrt{n}}} - z_{1-\frac{\alpha}{2}} < \frac{\bar{X} - m_1}{\frac{\sigma}{\sqrt{n}}} < \frac{m_0 - m_1}{\frac{\sigma}{\sqrt{n}}} + z_{1-\frac{\alpha}{2}} \mid H_1\right), \end{aligned}$$

de unde avem că

$$(6.3.1) \quad \beta = \Phi\left(\frac{m_0 - m_1}{\frac{\sigma}{\sqrt{n}}} + z_{1-\frac{\alpha}{2}}\right) - \Phi\left(\frac{m_0 - m_1}{\frac{\sigma}{\sqrt{n}}} - z_{1-\frac{\alpha}{2}}\right),$$

adică

$$(6.3.2) \quad \pi(m_1) = 1 - \Phi\left(\frac{m_0 - m_1}{\frac{\sigma}{\sqrt{n}}} + z_{1-\frac{\alpha}{2}}\right) + \Phi\left(\frac{m_0 - m_1}{\frac{\sigma}{\sqrt{n}}} - z_{1-\frac{\alpha}{2}}\right).$$

Deoarece $\Phi(+\infty) = 0.5$ și $\Phi(-\infty) = -0.5$, rezultă că pentru $n \rightarrow \infty$, vom avea $\beta \rightarrow 0$, adică $\pi(m_1) \rightarrow 1$. Așadar, putem determina, din relația precedentă valoarea lui n (volumul selecției) astfel încât puterea testului (corespunzător riscul de speța a doua) să fie atinsă pentru acel n .

Exemplul 6.3.7. Relativ la caracteristica X s-a efectuat o selecție de volum $n = 24$, obținându-se datele de selecție

3.59	4.27	2.83	2.28	3.69	2.81	3.16	3.03	3.65	3.84	4.32	3.21
3.20	3.32	3.82	3.27	3.77	3.56	3.29	4.08	2.97	3.29	3.78	3.58.

Știind că X urmează legea normală $\mathcal{N}(m, \sigma)$ cu abaterea standard teoretică cunoscută $\sigma = \sqrt{\text{Var}(X)} = 0.5$, vrem să verificăm ipoteza nulă $H_0 : m = 3.5$, când se consideră ipoteza alternativă $H_1 : m \neq 3.5$, cu nivelul de semnificație $\alpha = 0.01$, iar apoi să calculăm puterea testului, când $m = 3.4, 3.5, 3.6, 3.7$, și să determinăm volumul n optim al selecției când se consideră alternativă $H_1 : m = 3.4$, și riscul de speța a doua $\beta = 0.05$.

Când se consideră alternativa $H_1 : m \neq 3.5$ se folosește testul Z bilateral. Se determină cuantila

$$z_{1-\frac{\alpha}{2}} = 2.58 \quad \text{din relația} \quad \Phi\left(z_{1-\frac{\alpha}{2}}\right) = \frac{1-\alpha}{2} = \frac{0.99}{2} = 0.495.$$

Pe de altă parte avem că

$$\bar{x} = \frac{1}{n} \sum_{k=1}^n x_k = \frac{1}{24} (3.59 + 4.27 + \dots + 3.58) = 3.44,$$

astfel că

$$z = \frac{\bar{x} - m_0}{\frac{\sigma}{\sqrt{n}}} = \frac{3.44 - 3.5}{\frac{0.5}{\sqrt{24}}} = -0.57.$$

Deoarece $|z| = 0.57 < 2.58 = z_{1-\frac{\alpha}{2}}$, ipoteza $H_0 : m = 3.5$ este admisă.

Formula de calcul a puterii testului Z bilateral este (6.3.2). Prin urmare, se obține

$$\pi(3.4) = 1 - \Phi(0.4\sqrt{6} + 2.58) + \Phi(0.4\sqrt{6} - 2.58) = 0.06,$$

$$\pi(3.5) = 1 - \Phi(2.58) + \Phi(-2.58) = 0.01,$$

$$\pi(3.6) = 1 - \Phi(-0.4\sqrt{6} + 2.58) + \Phi(-0.4\sqrt{6} - 2.58) = 0.05,$$

$$\pi(3.7) = 1 - \Phi(-0.8\sqrt{6} + 2.58) + \Phi(-0.8\sqrt{6} - 2.58) = 0.27.$$

Pentru determinarea volumului optim al selecției, când se consideră alternativa $H_1 : m = 3.4$, se folosește formula (6.3.1). Astfel se obține

$$\begin{aligned} 0.05 = \beta &= \Phi\left(\frac{3.5 - 3.4}{\frac{0.5}{\sqrt{n}}} + 2.58\right) - \Phi\left(\frac{3.5 - 3.4}{\frac{0.5}{\sqrt{n}}} - 2.58\right) \\ &= \Phi\left(\frac{1}{5}\sqrt{n} + 2.58\right) - \Phi\left(\frac{1}{5}\sqrt{n} - 2.58\right). \end{aligned}$$

Când n este astfel încât $\frac{1}{5}\sqrt{n} - 2.58 < 0$, se observă că valoarea $\beta = 0.05$ nu poate fi atinsă. Dacă $\frac{1}{5}\sqrt{n} - 2.58 > 0$, atunci $n > 166$, astfel avem $\Phi\left(\frac{1}{5}\sqrt{n} + 2.58\right) \approx 0.5$. Prin urmare, pentru determinarea lui n optim avem relația $\Phi\left(\frac{1}{5}\sqrt{n} - 2.58\right) = 0.45$, și folosind *Anexa I* se obține $\frac{1}{5}\sqrt{n} - 2.58 = 1.65$, adică $\sqrt{n} = 5 \times 4.23 = 21.15$, de unde se obține $n = 448$.

Programul 6.3.8. Programul Matlab, ce urmează, pe lângă faptul că va efectua calculele din exemplul precedent, va reprezenta grafic și funcția putere. Punctele de pe curba puterii, ale căror valori au fost calculate în exemplul de mai sus, sunt marcate prin cerușe în Figura 6.1.

```

clear,clf
x=[3.59,4.27,2.83,2.28,3.69,2.81,3.16,3.03,3.65,...
   3.84,4.32,3.21,3.20,3.32,3.82,3.27,3.77,3.56,...
   3.29,4.08,2.97,3.29,3.78,3.58];
z=(mean(x)-3.5)/(0.5/sqrt(24));
z1=norminv(0.005,0,1); z2=norminv(0.995,0,1);
fprintf(' z=%6.3f\n z1=%6.3f\n z2=%6.3f\n',z,z1,z2)
m=3:0.01:4; m1=3.4:0.1:3.7;
c1=norminv(0.005,0,1); c2=norminv(0.995,0,1);
pib=1-normcdf((3.5-m)/(0.5/sqrt(24))+c2,0,1)...
    +normcdf((3.5-m)/(0.5/sqrt(24))+c1,0,1);
pib1=1-normcdf((3.5-m1)/(0.5/sqrt(24))+c2,0,1)...
    +normcdf((3.5-m1)/(0.5/sqrt(24))+c1,0,1);
plot(m,pib,'k-',m1,pib1,'o'), grid on, len=length(m1);
for i=1:len
    fprintf(' pi(%3.1f)=%5.3f\n',m1(i),pib1(i))
end
za=norminv(0.005,0,1); zb=norminv(0.95,0,1);
n=ceil(0.5^2*(zb-za)^2/(3.5-3.4)^2);
fprintf(' nopt=%3d',n)

```

Rezultatele obținute, în urma executării programului, sunt următoarele:

```

z=-0.567
z1=-2.576
z2= 2.576
pi(3.4)=0.055
pi(3.5)=0.010
pi(3.6)=0.055
pi(3.7)=0.269
nopt=446

```

Lema 6.3.9 (Neyman–Pearson). Fie caracteristica X cu legea de probabilitate dată prin $f(x; \theta)$, unde $\theta \in A \subset \mathbb{R}$ este parametrul necunoscut asupra căruia se face ipoteza nulă simplă $H_0 : \theta = \theta_0$, cu ipoteza alternativă simplă $H_1 : \theta = \theta_1 \neq \theta_0$. Se consideră o selecție repetată de volum n relativă la caracteristica X și nivelul de semnificație dat $\alpha \in (0, 1)$, atunci

$$\begin{aligned} \max \left\{ P \left((X_1, \dots, X_n) \in \mathcal{U} \mid H_1 \right) \mid \mathcal{U} \subset \mathbb{R}^n, P \left((X_1, \dots, X_n) \in \mathcal{U} \mid H_0 \right) = \alpha \right\} \\ = P \left((X_1, \dots, X_n) \in \tilde{\mathcal{U}} \mid H_0 \right), \end{aligned}$$

regiunea critică $\tilde{\mathcal{U}}$, fiind definită prin

$$\tilde{\mathcal{U}} = \left\{ (u_1, \dots, u_n) \in \mathbb{R}^n \mid \frac{g(u_1, \dots, u_n; \theta_1)}{g(u_1, \dots, u_n; \theta_0)} \geq k_\alpha > 0 \right\},$$

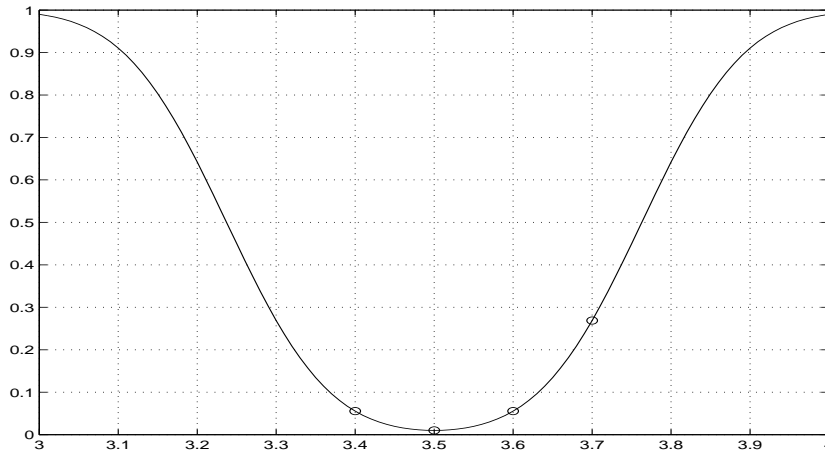


Figura 6.1: Curba funcției putere

unde

$$g(u_1, \dots, u_n; \theta) = \prod_{k=1}^n f(u_k; \theta).$$

Demonstrație. Facem demonstrația în două etape.

În prima etapă, considerăm că se cunoaște constanta $k_\alpha > 0$ și arătăm că

$$1 - \beta = P((X_1, \dots, X_n) \in \mathcal{U} \mid H_1) \leq P((X_1, \dots, X_n) \in \tilde{\mathcal{U}} \mid H_1) = 1 - \tilde{\beta}.$$

Pentru aceasta, dacă notăm $\mathcal{W} = \mathcal{U} \cap \tilde{\mathcal{U}}$, obținem succesiv

$$\begin{aligned} \int \dots \int_{\mathcal{U}} g(u_1, \dots, u_n; \theta_0) du_1 \dots du_n &= P((X_1, \dots, X_n) \in \mathcal{U} \mid H_0) = \alpha \\ &= P((X_1, \dots, X_n) \in \tilde{\mathcal{U}} \mid H_0) = \int \dots \int_{\tilde{\mathcal{U}}} g(u_1, \dots, u_n; \theta_0) du_1 \dots du_n, \end{aligned}$$

de unde avem că

$$\begin{aligned} \int \dots \int_{\mathcal{U} \setminus \mathcal{W}} g(u_1, \dots, u_n; \theta_0) du_1 \dots du_n &+ \int \dots \int_{\mathcal{W}} g(u_1, \dots, u_n; \theta_0) du_1 \dots du_n \\ &= \int \dots \int_{\tilde{\mathcal{U}} \setminus \mathcal{W}} g(u_1, \dots, u_n; \theta_0) du_1 \dots du_n + \\ &\quad + \int \dots \int_{\mathcal{W}} g(u_1, \dots, u_n; \theta_0) du_1 \dots du_n, \end{aligned}$$

adică

$$(6.3.3) \quad \int \dots \int_{\mathcal{U} \setminus \mathcal{W}} g(u_1, \dots, u_n; \theta_0) du_1 \dots du_n \\ = \int \dots \int_{\tilde{\mathcal{U}} \setminus \mathcal{W}} g(u_1, \dots, u_n; \theta_0) du_1 \dots du_n.$$

Putem trece acum la calculul diferenței

$$\begin{aligned} \beta - \tilde{\beta} &= (1 - \tilde{\beta}) - (1 - \beta) \\ &= P((X_1, \dots, X_n) \in \tilde{\mathcal{U}} \mid H_1) - P((X_1, \dots, X_n) \in \mathcal{U} \mid H_1) \\ &= \int \dots \int_{\tilde{\mathcal{U}}} g(u_1, \dots, u_n; \theta_1) du_1 \dots du_n - \int \dots \int_{\mathcal{U}} g(u_1, \dots, u_n; \theta_1) du_1 \dots du_n \\ &= \int \dots \int_{\tilde{\mathcal{U}} \setminus \mathcal{W}} g(u_1, \dots, u_n; \theta_1) du_1 \dots du_n \\ &\quad - \int \dots \int_{\mathcal{U} \setminus \mathcal{W}} g(u_1, \dots, u_n; \theta_1) du_1 \dots du_n. \end{aligned}$$

Scriem această diferență sub forma

$$\begin{aligned} \beta - \tilde{\beta} &= \int \dots \int_{\tilde{\mathcal{U}} \setminus \mathcal{W}} \frac{g(u_1, \dots, u_n; \theta_1)}{g(u_1, \dots, u_n; \theta_0)} g(u_1, \dots, u_n; \theta_0) du_1 \dots du_n \\ &\quad - \int \dots \int_{\mathcal{U} \setminus \mathcal{W}} \frac{g(u_1, \dots, u_n; \theta_1)}{g(u_1, \dots, u_n; \theta_0)} g(u_1, \dots, u_n; \theta_0) du_1 \dots du_n. \end{aligned}$$

Pentru fiecare integrală din membrul drept aplicăm formula de medie. Prin urmare există $(\xi_1, \xi_2, \dots, \xi_n) \in \tilde{\mathcal{U}} \setminus \mathcal{W}$ și $(\eta_1, \eta_2, \dots, \eta_n) \in \mathcal{U} \setminus \mathcal{W}$ astfel încât

$$\begin{aligned} \beta - \tilde{\beta} &= \frac{g(\xi_1, \dots, \xi_n; \theta_1)}{g(\xi_1, \dots, \xi_n; \theta_0)} \int \dots \int_{\tilde{\mathcal{U}} \setminus \mathcal{W}} g(u_1, \dots, u_n; \theta_0) du_1 \dots du_n \\ &\quad - \frac{g(\eta_1, \dots, \eta_n; \theta_1)}{g(\eta_1, \dots, \eta_n; \theta_0)} \int \dots \int_{\mathcal{U} \setminus \mathcal{W}} g(u_1, \dots, u_n; \theta_0) du_1 \dots du_n. \end{aligned}$$

Folosind relația (6.3.3), stabilită mai înainte, avem că

$$\begin{aligned} \beta - \tilde{\beta} &= \left[\frac{g(\xi_1, \dots, \xi_n; \theta_1)}{g(\xi_1, \dots, \xi_n; \theta_0)} - \frac{g(\eta_1, \dots, \eta_n; \theta_1)}{g(\eta_1, \dots, \eta_n; \theta_0)} \right] \times \\ &\quad \times \int \dots \int_{\mathcal{U} \setminus \mathcal{W}} g(u_1, \dots, u_n; \theta_0) du_1 \dots du_n. \end{aligned}$$

Dar integrala din membrul drept este pozitivă, iar din modul cum s-a definit regiunea critică $\tilde{\mathcal{U}}$ avem că

$$\frac{g(\xi_1, \dots, \xi_n; \theta_1)}{g(\xi_1, \dots, \xi_n; \theta_0)} \geq k_\alpha > \frac{g(\eta_1, \dots, \eta_n; \theta_1)}{g(\eta_1, \dots, \eta_n; \theta_0)},$$

deci $\beta - \tilde{\beta} > 0$ și ca urmare $1 - \tilde{\beta} > 1 - \beta$. Ceea ce încheie partea întâi a demonstrației.

În partea a doua arătăm existența constantei $k_\alpha > 0$.

Notăm prin $A(k)$ următoarea regiune

$$A(k) = \left\{ (u_1, \dots, u_n) \in \mathbb{R}^n \mid g(u_1, \dots, u_n; \theta_1) \geq k g(u_1, \dots, u_n; \theta_0) \right\},$$

și prin $s(k)$, următoarea probabilitate

$$s(k) = P((X_1, \dots, X_n) \in A(k) \mid H_0) = \int \dots \int_{A(k)} g(u_1, \dots, u_n; \theta_0) du_1 \dots du_n.$$

Se observă că $s(+\infty) = 0$, iar

$$s(0) = \int \dots \int_{\mathbb{R}^n} g(u_1, \dots, u_n; \theta_0) du_1 \dots du_n = 1.$$

De asemenea, avem că funcția $s(k)$ este monoton descrescătoare. Prin urmare, pentru $\alpha \in (0, 1)$ fixat, există $k_\alpha > 0$ astfel încât $s(k_\alpha) = \alpha$, adică

$$P((X_1, \dots, X_n) \in A(k_\alpha) \mid H_0) = \alpha.$$

Ceea ce încheie demonstrația lemei. □

Observația 6.3.10. Demonstrația lemei Neyman–Pearson a fost făcută în cazul continuu. În mod analog se demonstrează și pentru cazul discret, integralele multiple care apar fiind înlocuite cu sume multiple.

Definiția 6.3.11. *Testul pentru care puterea este maximă îl numim cel mai puternic test.*

Definiția 6.3.12. *Testul pentru care are loc inegalitatea*

$$1 - \beta = P((X_1, X_2, \dots, X_n) \in \mathcal{U} \mid H_1) > P((X_1, X_2, \dots, X_n) \in \mathcal{U} \mid H_0) = \alpha,$$

adică puterea testului este mai mare decât riscul de speța întâi se numește test nedepășat.

Proprietatea 6.3.13. Cel mai puternic test dat de lema Neyman–Pearson este un test nedeplasat, adică $1 - \tilde{\beta} > \alpha$.

Demonstrație. Cu notațiile de la lema Neyman–Pearson, avem că

$$\int \dots \int_{\tilde{\mathcal{U}}} g(u_1, \dots, u_n; \theta_1) du_1 \dots du_n \geq k_\alpha \int \dots \int_{\tilde{\mathcal{U}}} g(u_1, \dots, u_n; \theta_0) du_1 \dots du_n.$$

Distingem două cazuri.

Dacă avem $k_\alpha > 1$, atunci obținem

$$\begin{aligned} 1 - \tilde{\beta} &= \int \dots \int_{\tilde{\mathcal{U}}} g(u_1, \dots, u_n; \theta_1) du_1 \dots du_n \\ &\geq k_\alpha \int \dots \int_{\tilde{\mathcal{U}}} g(u_1, \dots, u_n; \theta_0) du_1 \dots du_n \\ &> \int \dots \int_{\tilde{\mathcal{U}}} g(u_1, \dots, u_n; \theta_0) du_1 \dots du_n = \alpha, \end{aligned}$$

de unde $1 - \tilde{\beta} > \alpha$.

În cazul când $k_\alpha \leq 1$, avem

$$\begin{aligned} \tilde{\beta} &= \int \dots \int_{\mathbb{C}\tilde{\mathcal{U}}} g(u_1, \dots, u_n; \theta_1) du_1 \dots du_n \\ &< k_\alpha \int \dots \int_{\mathbb{C}\tilde{\mathcal{U}}} g(u_1, \dots, u_n; \theta_0) du_1 \dots du_n \\ &\leq \int \dots \int_{\mathbb{C}\tilde{\mathcal{U}}} g(u_1, \dots, u_n; \theta_0) du_1 \dots du_n = 1 - \alpha, \end{aligned}$$

de unde $1 - \tilde{\beta} > \alpha$, ceea ce trebuie demonstrat. \square

Observația 6.3.14. Prin micșorarea riscului de speța întâi, pentru n fixat, crește constanta k_α și prin urmare scade puterea testului, $1 - \tilde{\beta}$, deci crește riscul de speța a doua.

Aplicația 6.3.15. Fie caracteristica X ce urmează legea normală $\mathcal{N}(m, \sigma)$, unde $m \in \mathbb{R}$ este necunoscut și $\sigma > 0$ este cunoscut. Fie ipoteza nulă $H_0 : m = m_0$, cu ipoteza alternativă $H_1 : m = m_1 \neq m_0$. Vrem să determinăm cel mai puternic test pentru verificarea acestei ipoteze nule.

Pentru aceasta aplicăm lema Neyman–Pearson.

Prima dată determinăm regiunea critică $\tilde{\mathcal{U}}$.

Deoarece

$$g(u_1, u_2, \dots, u_n; m) = \left(\frac{1}{\sigma\sqrt{2\pi}} \right)^n \exp \left\{ -\frac{1}{2\sigma^2} \sum_{k=1}^n (u_k - m)^2 \right\},$$

rezultă că

$$\frac{g(u_1, u_2, \dots, u_n; m_1)}{g(u_1, u_2, \dots, u_n; m_0)} = \exp \left\{ \frac{1}{2\sigma^2} \left[\sum_{k=1}^n (u_k - m_0)^2 - \sum_{k=1}^n (u_k - m_1)^2 \right] \right\}.$$

Condiția care definește regiunea critică $\tilde{\mathcal{U}}$

$$\frac{g(u_1, u_2, \dots, u_n; m_1)}{g(u_1, u_2, \dots, u_n; m_0)} \geq k_\alpha,$$

devine prin logaritmare

$$\sum_{k=1}^n (u_k - m_0)^2 - \sum_{k=1}^n (u_k - m_1)^2 \geq 2\sigma^2 \ln k_\alpha,$$

sau

$$2(m_1 - m_0) \sum_{k=1}^n u_k - n(m_1 - m_0)(m_1 + m_0) \geq 2\sigma^2 \ln k_\alpha,$$

adică

$$2(m_1 - m_0) \sum_{k=1}^n u_k \geq 2\sigma^2 \ln k_\alpha + n(m_1 - m_0)(m_1 + m_0).$$

Distingem în continuare două cazuri.

Dacă $m_0 > m_1$, avem

$$\frac{1}{n} \sum_{k=1}^n u_k = \bar{u} \leq \frac{\sigma^2 \ln k_\alpha}{(m_1 - m_0)n} + \frac{m_1 + m_0}{2},$$

și notând constanta din membrul drept prin K_α , avem că $\bar{u} \leq K_\alpha$. Așadar regiunea critică a celui mai puternic test este

$$\tilde{\mathcal{U}} = \{(u_1, u_2, \dots, u_n) \in \mathbb{R}^n \mid \bar{u} \leq K_\alpha\}.$$

Dacă $m_0 < m_1$, procedând analog, se ajunge la regiunea critică

$$\tilde{\mathcal{U}} = \{(u_1, u_2, \dots, u_n) \in \mathbb{R}^n \mid \bar{u} \geq K_\alpha\}.$$

Pentru determinarea constantei K_α (cazul $m_0 > m_1$), se scrie relația

$$\alpha = P\left((X_1, X_2, \dots, X_n) \in \tilde{\mathcal{U}} \mid H_0\right) = P\left(\bar{X} \leq K_\alpha \mid H_0\right)$$

sau

$$P\left(\frac{\bar{X} - m_0}{\frac{\sigma}{\sqrt{n}}} \leq \frac{K_\alpha - m_0}{\frac{\sigma}{\sqrt{n}}} \mid H_0\right) = \alpha.$$

Prin urmare, K_α se va determina astfel încât $\Phi\left(\frac{K_\alpha - m_0}{\frac{\sigma}{\sqrt{n}}}\right) = \alpha$, unde funcția Φ a lui Laplace este cea definită prin

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{t^2}{2}} dt.$$

Așadar avem cuantila z_α , ca fiind $z_\alpha = \frac{K_\alpha - m_0}{\frac{\sigma}{\sqrt{n}}}$, de unde $K_\alpha = m_0 + \frac{\sigma}{\sqrt{n}} z_\alpha$.

În acest fel, regiunea critică se poate scrie

$$\begin{aligned} \tilde{\mathcal{U}} &= \left\{ (u_1, u_2, \dots, u_n) \in \mathbb{R}^n \mid \bar{u} \leq m_0 + \frac{\sigma}{\sqrt{n}} z_\alpha \right\} \\ &= \left\{ (u_1, u_2, \dots, u_n) \in \mathbb{R}^n \mid \frac{\bar{u} - m_0}{\frac{\sigma}{\sqrt{n}}} \leq z_\alpha \right\}. \end{aligned}$$

În cazul $m_0 < m_1$, procedând la fel, se ajunge la

$$\begin{aligned} \tilde{\mathcal{U}} &= \left\{ (u_1, u_2, \dots, u_n) \in \mathbb{R}^n \mid \bar{u} \geq m_0 + \frac{\sigma}{\sqrt{n}} z_{1-\alpha} \right\} \\ &= \left\{ (u_1, u_2, \dots, u_n) \in \mathbb{R}^n \mid \frac{\bar{u} - m_0}{\frac{\sigma}{\sqrt{n}}} \geq z_{1-\alpha} \right\}. \end{aligned}$$

Observația 6.3.16. În aplicația precedentă, regiunea critică $\tilde{\mathcal{U}}$ nu depinde de m_1 numai prin faptul că $m_1 < m_0$, respectiv $m_1 > m_0$. Drept urmare, ipoteza H_1 se poate înlocui, pentru cele două situații prin $H_1 : m > m_0$, ceea ce ne conduce la testul Z unilateral dreapta, respectiv $H_1 : m < m_0$, care ne conduce la testul Z unilateral stânga. Aceasta confirmă justetea alegerii regiunilor critice în cazul testului Z unilateral dreapta și respectiv stânga.

Observația 6.3.17. Când se consideră ipoteza alternativă $H_1 : m \neq m_0$ (testul Z bilateral), utilizând lema Neyman–Pearson, nu putem construi un cel mai puternic test. Altă metodă va fi prezentată pentru construirea acestuia în acest caz.

Aplicația 6.3.18. Vrem să determinăm în continuare puterea testului Z stabilit la Aplicația 6.3.15.

Din nou distingem cele două cazuri.

Dacă $m_0 > m_1$, avem succesiv

$$\begin{aligned} \pi(m_1) &= 1 - \tilde{\beta} = P\left((X_1, X_2, \dots, X_n) \in \tilde{\mathcal{U}} \mid m = m_1\right) \\ &= P\left(\bar{X} \leq m_0 + z_\alpha \frac{\sigma}{\sqrt{n}} \mid m = m_1\right) \\ &= P\left(\frac{\bar{X} - m_1}{\frac{\sigma}{\sqrt{n}}} \leq \frac{m_0 - m_1}{\frac{\sigma}{\sqrt{n}}} + z_\alpha \mid m = m_1\right) = \Phi\left(\frac{m_0 - m_1}{\frac{\sigma}{\sqrt{n}}} + z_\alpha\right), \end{aligned}$$

unde funcția Φ a lui Laplace este cea din Aplicația 6.3.15.

Dacă $m_0 < m_1$, de asemenea, avem

$$\begin{aligned}\pi(m_1) &= 1 - \tilde{\beta} = P\left((X_1, X_2, \dots, X_n) \in \tilde{\mathcal{U}} \mid m = m_1\right) \\ &= P\left(\bar{X} \geq m_0 + z_{1-\alpha} \frac{\sigma}{\sqrt{n}} \mid m = m_1\right) \\ &= P\left(\frac{\bar{X} - m_1}{\frac{\sigma}{\sqrt{n}}} \geq \frac{m_0 - m_1}{\frac{\sigma}{\sqrt{n}}} + z_{1-\alpha} \mid m = m_1\right) = 1 - \Phi\left(\frac{m_0 - m_1}{\frac{\sigma}{\sqrt{n}}} + z_{1-\alpha}\right).\end{aligned}$$

Se vede, încă odată, în ambele cazuri, că $\pi(m_1) \rightarrow 1$, când $n \rightarrow \infty$, adică $\tilde{\beta} \rightarrow 0$.

Să încheiem prin modul de determinare a volumului n al selecției, pentru a se atinge o putere a testului (sau un risc de speța a doua) apriori fixată. Desigur că α este considerat de asemenea dat.

În cazul $m_0 > m_1$, din $P\left((X_1, X_2, \dots, X_n) \in \tilde{\mathcal{U}} \mid H_0\right) = \alpha$, am determinat $K_\alpha = m_0 + \frac{\sigma}{\sqrt{n}} z_\alpha$. Pe de altă parte, $P\left((X_1, X_2, \dots, X_n) \in \tilde{\mathcal{U}} \mid H_1\right) = 1 - \tilde{\beta}$ se rescrie în felul următor

$$P\left(\frac{\bar{X} - m_1}{\frac{\sigma}{\sqrt{n}}} \leq \frac{K_\alpha - m_1}{\frac{\sigma}{\sqrt{n}}} \mid H_1\right) = 1 - \tilde{\beta}, \quad \text{adică} \quad \Phi\left(\frac{K_\alpha - m_1}{\frac{\sigma}{\sqrt{n}}}\right) = 1 - \tilde{\beta},$$

de unde se obține o altă relație ce îl conține pe K_α , anume $\frac{K_\alpha - m_1}{\frac{\sigma}{\sqrt{n}}} = z_{1-\tilde{\beta}}$. Comparând cele două relații ce îl conțin pe K_α , avem că $m_0 + \frac{\sigma}{\sqrt{n}} z_\alpha = m_1 + \frac{\sigma}{\sqrt{n}} z_{1-\tilde{\beta}}$, de unde se obține

$$(6.3.4) \quad n = \frac{\sigma^2 \left(z_{1-\tilde{\beta}} - z_\alpha\right)^2}{(m_0 - m_1)^2}.$$

În cazul $m_0 < m_1$, avem următoarele două relații, care îl conțin pe K_α ,

$$K_\alpha = m_0 + \frac{\sigma}{\sqrt{n}} z_{1-\alpha} \quad \text{și} \quad \frac{K_\alpha - m_1}{\frac{\sigma}{\sqrt{n}}} = z_{\tilde{\beta}},$$

de unde se obține

$$(6.3.5) \quad n = \frac{\sigma^2 \left(z_{\tilde{\beta}} - z_{1-\alpha}\right)^2}{(m_0 - m_1)^2}.$$

Exemplul 6.3.19. Să considerăm aceleași date de selecție de la Exemplul 6.3.7. Vom verifica ipoteza nulă $H_0 : m = 3.5$, când se consideră respectiv ipotezele alternative $H_1 : m < 3.5$, $H_1 : m > 3.5$, iar nivelul de semnificație este $\alpha = 0.01$. Vom calcula puterea pentru fiecare din testele considerate, când $m = 3.4, 3.5, 3.6, 3.7$, iar apoi vom determina volumul n optim al selecției când se consideră alternativa $H_1 : m = 3.4$, respectiv $H_1 : m = 3.7$, cu riscul de speța a doua fixat $\beta = 0.05$.

Când se consideră alternativa $H_1 : m > 3.5$ se folosește testul Z unilateral dreapta. În acest caz, cuantila $z_{1-\alpha} = 2.33$ se determină din $\Phi(z_{1-\alpha}) = \frac{1}{2} - \alpha = 0.49$. Deoarece $z = -0.57 < 2.33 = z_{1-\alpha}$, rezultă că ipoteza nulă este admisă.

Pentru alternativa $H_1 : m < 3.5$ se folosește testul Z unilateral stânga. Se determină $z_\alpha = -2.33$ din relația $\Phi(-z_\alpha) = \frac{1}{2} - \alpha = 0.49$. Deoarece are loc inegalitatea $z = -0.57 > -2.33 = z_\alpha$, se admite ipoteza H_0 și în acest ultim caz.

Pentru testul Z unilateral dreapta, formula de calcul a puterii este

$$\pi(m) = \frac{1}{2} - \Phi\left(\frac{m_0 - m}{\frac{\sigma}{\sqrt{n}}} + z_{1-\alpha}\right),$$

și prin urmare

$$\begin{aligned}\pi(3.4) &= \frac{1}{2} - \Phi(0.4\sqrt{6} + 2.33) = 0, & \pi(3.6) &= \frac{1}{2} - \Phi(-0.4\sqrt{6} + 2.33) = 0.09, \\ \pi(3.5) &= \frac{1}{2} - \Phi(2.33) = 0.01, & \pi(3.7) &= \frac{1}{2} - \Phi(-0.8\sqrt{6} + 2.33) = 0.36.\end{aligned}$$

În cazul testului Z unilateral stânga, avem

$$\pi(m) = \frac{1}{2} + \Phi\left(\frac{m_0 - m}{\frac{\sigma}{\sqrt{n}}} + z_\alpha\right),$$

deci

$$\begin{aligned}\pi(3.4) &= \frac{1}{2} + \Phi(0.4\sqrt{6} - 2.33) = 0.09, & \pi(3.6) &= \frac{1}{2} + \Phi(-0.4\sqrt{6} - 2.33) = 0, \\ \pi(3.5) &= \frac{1}{2} + \Phi(-2.33) = 0.01, & \pi(3.7) &= \frac{1}{2} + \Phi(-0.8\sqrt{6} - 2.33) = 0.\end{aligned}$$

Pentru determinarea volumului optim al selecției, când folosim testul unilateral stânga, cu alternativa $H_1 : m = 3.4$, se utilizează formula (6.3.4):

$$n = \frac{\sigma^2 (z_{1-\beta} - z_\alpha)^2}{(m_0 - m_1)^2} = \frac{0.25 (1.65 + 2.33)^2}{(3.5 - 3.4)^2},$$

de unde se obține $n = 396$.

Pentru determinarea volumului optim al selecției, când folosim testul unilateral dreapta, cu alternativa $H_1: m = 3.7$, se utilizează formula (6.3.5). Astfel volumul optim n al selecției se calculează cu formula:

$$n = \frac{\sigma^2 (z_\beta - z_{1-\alpha})^2}{(m_0 - m_1)^2} = \frac{0.25 (-1.65 - 2.33)^2}{(3.5 - 3.7)^2} = \frac{1584}{16},$$

de unde se obține $n = 99$.

Programul 6.3.20. Programul Matlab următor, efectuează calculele din exemplul precedent, iar în plus, reprezintă grafic curbele putere, în cazul celor două alternative. Punctele de pe curbele putere, care au fost calculate în exemplul precedent, sunt marcate pe Figura 6.2 prin cerușe.

```
clear,clf
x=[3.59,4.27,2.83,2.28,3.69,2.81,3.16,3.03,3.65,...
   3.84,4.32,3.21,3.20,3.32,3.82,3.27,3.77,3.56,...
   3.29,4.08,2.97,3.29,3.78,3.58];
z=(mean(x)-3.5)/(0.5/sqrt(24));
fprintf(' z=%6.3f\n',z)
z1=norminv(0,0,1); z2=norminv(0.99,0,1);
fprintf('      m gt 3.5\n')
fprintf(' z1=%6.3f, z2=%6.3f\n',z1,z2)
z1=norminv(0.01,0,1); z2=norminv(1,0,1);
fprintf('      m lt 3.5\n')
fprintf(' z1=%6.3f, z2=%6.3f\n',z1,z2)
m=3:0.01:4; m1=3.4:0.1:3.7;
c2=norminv(0.99,0,1);
pib=1-normcdf(c2+(3.5-m)/(0.5/sqrt(24)),0,1);
pib1(1,:)=1-normcdf(c2+(3.5-m1)/(0.5/sqrt(24)),0,1);
subplot(2,1,1)
plot(m,pib,'k-',m1,pib1(1,:), 'o'), grid on
c1=norminv(0.01,0,1);
pib=normcdf(c1+(3.5-m)/(0.5/sqrt(24)),0,1);
pib1(2,:)=normcdf(c1+(3.5-m1)/(0.5/sqrt(24)),0,1);
subplot(2,1,2)
plot(m,pib,'k-',m1,pib1(2,:), 'o'), grid on
len=length(m1);
for i=1:len
    for j=1:2
        fprintf(' pi(%3.1f)=%5.3f      ',m1(i),pib1(j,i))
    end
    fprintf('\n')
end
zb=norminv(0.95,0,1); za=norminv(0.01,0,1);
nr=ceil(0.5^2*(zb-za)^2/(3.5-3.4)^2);
zb=norminv(0.05,0,1); za=norminv(0.99,0,1);
nl=ceil(0.5^2*(zb-za)^2/(3.5-3.7)^2);
fprintf(' nr=%3d, nl=%3d',nr,nl)
```

În urma executării programului, pe lângă graficele curbelor de putere, sunt afișate și următoarele rezultate:

```
z=-0.567
  m gt 3.5
z1= -Inf, z2= 2.326
  m lt 3.5
z1=-2.326, z2= Inf
pi(3.4)=0.000    pi(3.4)=0.089
pi(3.5)=0.010    pi(3.5)=0.010
pi(3.6)=0.089    pi(3.6)=0.000
pi(3.7)=0.357    pi(3.7)=0.000
nr=395,  nl= 99
```

care reprezintă intervalele numerice ($z1, z2$) și volumul optim, când se consideră cele două alternative.

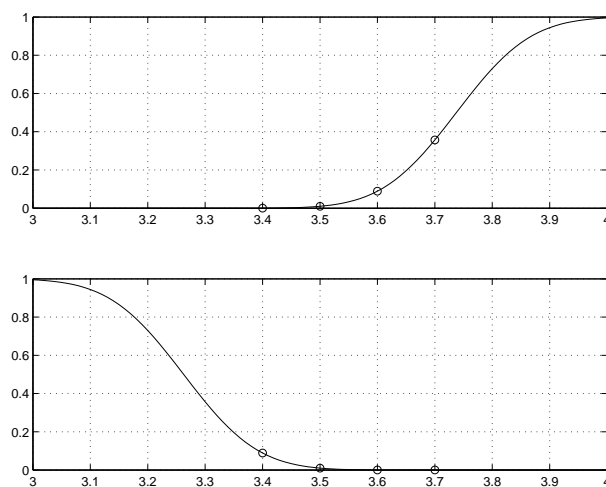


Figura 6.2: Curbele funcțiilor putere

6.4 Testul T (Student) privind media teoretică

Fie caracteristica X ce urmează legea normală $\mathcal{N}(m, \sigma)$, cu ambii parametri $m \in \mathbb{R}$ și $\sigma > 0$ necunoscuți. Relativ la valoarea medie a acestei caracteristici, se face ipoteza nulă $H_0 : m = m_0$, cu una din alternativele

$H_1 : m \neq m_0$, când obținem *testul T bilateral*,

$H_1 : m > m_0$, când obținem *testul T unilateral dreapta*,

$H_1 : m < m_0$, când obținem *testul T unilateral stânga*.

Pentru verificarea acestei ipoteze se consideră o selecție repetată de volum n , cu datele de selecție x_1, \dots, x_n și corespunzător variabilele de selecție X_1, \dots, X_n . Conform Proprietății 4.2.26, statistica

$$T = \frac{\bar{X} - m}{\frac{\bar{\sigma}}{\sqrt{n}}} = \frac{\bar{X} - m}{\sqrt{\frac{\bar{\mu}_2}{n-1}}},$$

unde

$$\bar{X} = \frac{1}{n} \sum_{k=1}^n X_k, \quad \bar{\sigma}^2 = \frac{1}{n-1} \sum_{k=1}^n (X_k - \bar{X})^2 = \frac{n}{n-1} \bar{\mu}_2,$$

urmează legea Student cu $n-1$ grade de libertate.

Prin urmare, pentru nivelul de semnificație $\alpha \in (0, 1)$ dat, se poate determina intervalul numeric (t_1, t_2) astfel încât

$$P(T \in (t_1, t_2) \mid H_0) = F_{n-1}(t_2) - F_{n-1}(t_1) = 1 - \alpha,$$

unde

$$F_m(t) = \frac{\Gamma\left(\frac{m+1}{2}\right)}{\sqrt{m\pi}\Gamma\left(\frac{m}{2}\right)} \int_{-\infty}^t \left(1 + \frac{x^2}{m}\right)^{-\frac{m+1}{2}} dx, \quad t \in \mathbb{R},$$

este funcția de repartiție pentru legea Student cu m grade de libertate (tabelată în *Anexa II*, pentru anumite valori).

Intervalul numeric (t_1, t_2) pentru statistica T nu este determinat în mod unic din condiția de mai sus. În funcție de alternativa H_1 aleasă, se consideră suplimentar:

$$t_1 = -t_2, \quad t_2 = t_{n-1, 1-\frac{\alpha}{2}}, \text{ dacă } H_1 : m \neq m_0,$$

$$t_1 = -\infty, \quad t_2 = t_{n-1, 1-\alpha}, \text{ dacă } H_1 : m > m_0,$$

$$t_1 = t_{n-1, \alpha}, \quad t_2 = +\infty, \text{ dacă } H_1 : m < m_0,$$

unde $t_{m, \gamma}$ este cuantila de ordin γ a legii Student cu m grade de libertate, adică

$$F_m(t_{m, \gamma}) = \gamma.$$

Corespunzător intervalului (t_1, t_2) se consideră respectiv regiunea critice \mathcal{U} definite prin:

$$\mathcal{U} = \left\{ (u_1, u_2, \dots, u_n) \in \mathbb{R}^n \mid \left| \frac{\bar{u} - m_0}{\frac{\bar{\sigma}_u}{\sqrt{n}}} \right| \geq t_{n-1, 1-\frac{\alpha}{2}} \right\},$$

$$\mathcal{U} = \left\{ (u_1, u_2, \dots, u_n) \in \mathbb{R}^n \mid \frac{\bar{u} - m_0}{\frac{\bar{\sigma}_u}{\sqrt{n}}} \geq t_{n-1, 1-\alpha} \right\},$$

$$\mathcal{U} = \left\{ (u_1, u_2, \dots, u_n) \in \mathbb{R}^n \mid \frac{\bar{u} - m_0}{\frac{\bar{\sigma}_u}{\sqrt{n}}} \leq t_{n-1, \alpha} \right\}.$$

Aici s-au utilizat notațiile $\bar{u} = \frac{1}{n} \sum_{k=1}^n u_k$, respectiv $\bar{\sigma}_u^2 = \frac{1}{n-1} \sum_{k=1}^n (u_k - \bar{u})^2$. Se verifică imediat că $P((X_1, X_2, \dots, X_n) \in \mathcal{U} \mid H_0) = \alpha$, iar cele trei moduri de definire a regiunii critice \mathcal{U} ne conduc respectiv la testul T bilateral, unilateral dreapta și unilateral stânga.

Odată construită regiunea critică \mathcal{U} , folosind datele de selecție x_1, x_2, \dots, x_n , ipoteza nulă H_0 va fi admisă dacă $(x_1, x_2, \dots, x_n) \notin \mathcal{U}$, iar în caz contrar va fi respinsă. Remarcăm de asemenea că regiunea critică \mathcal{U} corespunde mulțimii complementare intervalului (t_1, t_2) .

Etapele aplicării testului T

1. Se dau: α ; x_1, x_2, \dots, x_n ; m_0 ;
2. Se calculează intervalul (t_1, t_2) astfel încât $F_{n-1}(t_1) - F_{n-1}(t_2) = 1 - \alpha$, (după cum s-a prezentat înainte);
3. Se calculează

$$t = \frac{\bar{x} - m_0}{\frac{\bar{\sigma}}{\sqrt{n}}}, \quad \text{unde} \quad \bar{x} = \frac{1}{n} \sum_{k=1}^n x_k, \quad \bar{\sigma}^2 = \frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x})^2;$$

4. Concluzia: dacă $t \in (t_1, t_2)$ ipoteza H_0 este admisă, în caz contrar ipoteza este respinsă.

Observația 6.4.1. Când numărul gradelor de libertate tinde la infinit, conform teoremei limită centrală, avem că legea Student converge în repartiție la legea normală $\mathcal{N}(0, 1)$. Prin urmare, dacă volumul n al selecției este mare ($n > 30$) se poate utiliza testul Z pentru verificarea ipotezei nule $H_0 : m = m_0$, prin utilizarea statisticii T în loc de statistica Z . Toate rezultatele de la testul Z rămân, așadar, adevărate în acest caz.

Observația 6.4.2. Având în vedere că funcția de repartiție a legii Student este strict crescătoare, iar pentru determinarea intervalului (t_1, t_2) este necesară inversarea acestei funcții, se poate renunța la operația de inversare. Pentru aceasta se determină ceea ce se numește *valoare critică*, ca și la testul Z , și pe care o notăm prin c . Vom analiza pe rând cele trei alternative.

Dacă se consideră testul bilateral, făcând notația

$$1 - \frac{c}{2} = F_{n-1}(|t|), \quad \text{adică} \quad c = 2(1 - F_{n-1}(|t|)),$$

ipoteza nulă va fi respinsă dacă $1 - \frac{c}{2} \geq 1 - \frac{\alpha}{2}$, adică $c \leq \alpha$.

Pentru testul unilateral dreapta, se notează

$$1 - c = F_{n-1}(t), \quad \text{adică} \quad c = 1 - F_{n-1}(t),$$

iar ipoteza nulă este respinsă, dacă $1 - c \geq 1 - \alpha$, adică, la fel ca mai înainte, dacă $c \leq \alpha$.

Un raționament analog, pentru testul unilateral stânga ne conduce la valoarea critică, calculată după formulă $c = F_{n-1}(t)$. Drept urmare, ipoteza nulă este respinsă, dacă $c \leq \alpha$.

În rezumat, ipoteza nulă este respinsă, dacă $c \leq \alpha$, unde

$$c = \begin{cases} 2(1 - F_{n-1}(|t|)), & \text{dacă } H_1 : m \neq m_0, \\ 1 - F_{n-1}(t), & \text{dacă } H_1 : m > m_0, \\ F_{n-1}(t), & \text{dacă } H_1 : m < m_0. \end{cases}$$

6.4.1 Funcția `ttest`

Sistemul Matlab, prin *Statistics toolbox*, dispune de funcția `ttest`, cu aplicabilitate la testul T . Apelarea acestei funcții se face prin:

```
h=ttest(x,m0)
h=ttest(x,m0,alpha)
h=ttest(x,m0,alpha,tail)
[h,c,ci]=ttest(x,m0,alpha,tail)
```

În urma executării acestor instrucțiuni, se efectuează testul T asupra datelor conținute în vectorul `x`, folosind nivelul de semnificație `alpha`, care are valoarea implicită `alpha=0.05`.

Parametrul `tail` specifică una din cele trei alternative, care conduc la testul bilateral (`tail=0`, implicit), unilateral dreapta (`tail=1`) și unilateral stânga (`tail=-1`). Dacă `h=1`, atunci ipoteza nulă va fi respinsă, respectiv dacă `h=0`, ipoteza nu poate fi respinsă.

Ultima formă de apel permite, de asemenea, obținerea valorii critice `c`, precum și a intervalului de încredere pentru media teoretică, corespunzător probabilității de încredere `1-alpha`, obținut în vectorul cu două componente `ci`.

Programul 6.4.3. Programul Matlab, ce urmează, aplică testul T , pentru datele din Exemplul 6.2.2, considerând că dispersia este necunoscută. Ipoteza nulă este

$H: m = 16$, iar ca ipoteze alternative se consideră, pe rând, cele trei, care conduc, respectiv la testele bilateral, unilateral stânga și unilateral dreapta. Programul va determina, de asemenea, intervalele de încredere pentru media teoretică, în cele trei cazuri, folosind $\alpha = 0.01$.

```
x=[11*ones(1,4),13*ones(1,6),15*ones(1,12),...
  17*ones(1,10),20*ones(1,8)];
fprintf(' h      c      ci      \n')
fprintf(' _____ \n')
for i=-1:1
    [h,c,ci]=ttest(x,16,0.01,i);
    fprintf(' %d  %5.4f  (%4.2f,%4.2f) \n',h,c,ci)
end
```

În urma executării programului, se obțin rezultatele:

h	c	ci
0	0.3261	(-Inf, 16.87)
0	0.6522	(14.61, 16.99)
0	0.6739	(14.73, Inf)

Se observă că pentru $\text{tail}=-1$ și $\text{tail}=1$, se construiesc intervale nemărginite, respectiv la stânga și la dreapta.

6.5 Testul raportului verosimilităților

Fie caracteristica X cu legea de probabilitate $f(x; \theta)$, unde parametrul necunoscut $\theta \in \mathcal{A} \subset \mathbb{R}^p$. Relativ la parametrul θ , se consideră ipoteza nulă $H_0: \theta \in \mathcal{A}_0$ cu alternativa $H_1: \theta \in \mathcal{A} \setminus \mathcal{A}_0$. Se consideră o selecție repetată de volum n cu ajutorul căreia se construiește statistica

$$\hat{\Lambda} = \hat{\Lambda}(X_1, \dots, X_n) = \frac{\sup_{\theta \in \mathcal{A}_0} g(X_1, \dots, X_n; \theta)}{\sup_{\theta \in \mathcal{A}} g(X_1, \dots, X_n; \theta)},$$

unde

$$g(u_1, \dots, u_n; \theta) = \prod_{k=1}^n f(u_k; \theta)$$

este funcția de verosimilitate. Din felul cum a fost definită statistica $\hat{\Lambda}$, avem că $0 < \hat{\Lambda} < 1$, iar ipoteza H_0 va putea fi acceptată dacă $\hat{\Lambda}$ este apropiată de 1. Folosind această observație, pentru un nivel de semnificație $\alpha \in (0, 1)$ dat, se determină regiunea critică din relația $P(\hat{\Lambda} \leq \lambda_\alpha \mid H_0) = \alpha$, unde λ_α este cuantila de ordin α , pentru legea de probabilitate a statisticii $\hat{\Lambda}$.

Exemplul 6.5.1. Fie caracteristica X ce urmează legea normală $\mathcal{N}(m, \sigma)$, unde parametrii $m \in \mathbb{R}$ și $\sigma > 0$ sunt necunoscuți. Vrem să verificăm următoarea ipoteză nulă $H_0 : (m = m_0, \sigma > 0)$, cu alternativa $H_1 : (m \neq m_0, \sigma > 0)$. Pentru aceasta considerăm o selecție repetată de volum n și nivelul de semnificație $\alpha \in (0, 1)$ fixat. Funcția de verosimilitate este dată prin

$$g(u_1, \dots, u_n; m, \sigma) = \left(\frac{1}{\sigma\sqrt{2\pi}} \right)^n \exp \left[-\frac{1}{2\sigma^2} \sum_{k=1}^n (u_k - m)^2 \right].$$

Când ipoteza H_0 este adevărată, atunci estimatorul de verosimilitate maximă pentru σ este dat prin

$$\hat{\sigma}_1 = \sqrt{\frac{1}{n} \sum_{k=1}^n (X_k - m_0)^2}.$$

Obținem astfel, în acest caz, că

$$\begin{aligned} \sup_{\sigma > 0} g(X_1, \dots, X_n; m_0, \sigma) &= \left(\frac{1}{\hat{\sigma}_1 \sqrt{2\pi}} \right)^n \exp \left[-\frac{1}{2\hat{\sigma}_1^2} \sum_{k=1}^n (X_k - m_0)^2 \right] \\ &= \left(\frac{n}{2\pi e \sum_{k=1}^n (X_k - m_0)^2} \right)^{\frac{n}{2}}. \end{aligned}$$

Pe de altă parte, pentru determinarea numitorului statisticii $\hat{\Lambda}$, vom căuta estimatorii de verosimilitate maximă pentru parametrii $m \in \mathbb{R}$ și $\sigma > 0$ necunoscuți.

Sistemul de verosimilitate maximă

$$\begin{cases} \frac{\partial \ln g(X_1, X_2, \dots, X_n; m, \sigma)}{\partial m} = 0 \\ \frac{\partial \ln g(X_1, X_2, \dots, X_n; m, \sigma)}{\partial \sigma} = 0, \end{cases}$$

ne dă estimatorii de verosimilitate maximă pentru m și σ , anume

$$\begin{aligned} \hat{m} &= \frac{1}{n} \sum_{k=1}^n X_k = \bar{X}, \\ \hat{\sigma} &= \sqrt{\frac{1}{n} \sum_{k=1}^n (X_k - \bar{X})^2} = \sqrt{\bar{\mu}_2}. \end{aligned}$$

Obținem astfel că

$$\begin{aligned} \sup_{m \in \mathbb{R}, \sigma > 0} g(X_1, X_2, \dots, X_n; m, \sigma) &= \left(\frac{1}{\sqrt{2\pi\bar{\mu}_2}} \right)^n \exp \left[-\frac{1}{2\bar{\mu}_2} \sum_{k=1}^n (X_k - \bar{X})^2 \right] \\ &= \left(\frac{n}{2\pi e \sum_{k=1}^n (X_k - \bar{X})^2} \right)^{\frac{n}{2}}. \end{aligned}$$

Putem scrie acum statistica raportului verosimilităților

$$\hat{\Lambda} = \left(\sum_{k=1}^n (X_k - \bar{X})^2 \bigg/ \sum_{k=1}^n (X_k - m_0)^2 \right)^{\frac{n}{2}}.$$

Folosind formula lui König (Observația 3.3.23):

$$\sum_{k=1}^n (X_k - m_0)^2 = \sum_{k=1}^n (X_k - \bar{X})^2 + n(\bar{X} - m_0)^2,$$

avem că

$$\hat{\Lambda} = \left(1 + \frac{n(\bar{X} - m_0)^2}{\sum_{k=1}^n (X_k - \bar{X})^2} \right)^{-\frac{n}{2}}.$$

Dacă se introduce statistica T definită prin

$$T = \frac{\bar{X} - m_0}{\frac{\bar{\sigma}}{\sqrt{n}}} = \frac{\bar{X} - m_0}{\sqrt{\frac{1}{n-1} \sum_{k=1}^n (X_k - \bar{X})^2} \bigg/ \sqrt{n}},$$

atunci

$$\hat{\Lambda} = \left(\frac{1}{1 + \frac{T^2}{n-1}} \right)^{\frac{n}{2}}.$$

Regiunea critică \mathcal{U} , pentru $\alpha \in (0, 1)$, se obține din $P(\hat{\Lambda} \leq \lambda_\alpha \mid H_0) = \alpha$, care este același lucru cu $P(|T| < t \mid m = m_0) = 1 - \alpha$. Altfel spus, determinarea cuantilei λ_α revine la determinarea lui $t > 0$, astfel încât $P(|T| < t \mid m = m_0) = 1 - \alpha$. Dar statistica T , în ipoteza H_0 , urmează legea Student cu $n - 1$ grade de libertate (Proprietatea 4.2.26), deci $t = t_{n-1, 1-\frac{\alpha}{2}}$, adică cuantila de ordin $1 - \frac{\alpha}{2}$ a legii Student cu $n - 1$ grade de libertate. Putem scrie în final regiunea critică:

$$\mathcal{U} = \left\{ (u_1, u_2, \dots, u_n) \in \mathbb{R}^n \mid \frac{|\bar{u} - m_0|}{\frac{\bar{\sigma}}{\sqrt{n}}} \geq t_{n-1, 1-\frac{\alpha}{2}} \right\},$$

adică chiar regiunea critică de la testul T bilateral. Prin urmare, cu această metodă obținem testul T bilateral.

Teorema 6.5.2. Fie statistica $\hat{\Lambda}$ a raportului verosimilităților, atunci statistica $-2 \ln \hat{\Lambda}$ urmează legea χ^2 cu p grade de libertate, când $n \rightarrow \infty$, p fiind numărul parametrilor necunoscuți.

Demonstrație. Vom considera numai cazul $p = 1$.

Folosind formula lui Taylor cu doi termeni avem că

$$\begin{aligned} & \ln g(X_1, \dots, X_n; \theta_0) - \ln g(X_1, \dots, X_n; \hat{\Theta}) \\ &= (\theta_0 - \hat{\Theta}) \frac{\partial \ln g(X_1, \dots, X_n; \hat{\Theta})}{\partial \theta} + \frac{1}{2} (\theta_0 - \hat{\Theta})^2 \frac{\partial^2 \ln g(X_1, \dots, X_n; \xi)}{\partial \theta^2}, \end{aligned}$$

unde $\xi \in (\theta_0, \hat{\Theta})$ sau $\xi \in (\hat{\Theta}, \theta_0)$. Dar $\hat{\Theta}$ este estimator de verosimilitate maximă pentru θ , prin urmare

$$\frac{\partial \ln g(X_1, \dots, X_n; \hat{\Theta})}{\partial \theta} = 0$$

deci avem că

$$-2 \ln \hat{\Lambda} = -(\theta_0 - \hat{\Theta})^2 \frac{\partial^2 \ln g(X_1, \dots, X_n; \xi)}{\partial \theta^2}.$$

Pe de altă parte, dacă ipoteza $H_0 : \theta = \theta_0$ este adevărată, se știe că $\hat{\Theta} \xrightarrow{\text{a.s.}} \theta_0$, deci $\xi \xrightarrow{\text{a.s.}} \theta_0$, când $n \rightarrow \infty$. Putem scrie atunci că

$$\begin{aligned} & \frac{\partial^2 \ln g(X_1, \dots, X_n; \xi)}{\partial \theta^2} \cong \frac{\partial^2 \ln g(X_1, \dots, X_n; \theta_0)}{\partial \theta^2} \\ &= \sum_{k=1}^n \frac{\partial^2 \ln f(X_k; \theta_0)}{\partial \theta^2} = n \cdot \frac{1}{n} \sum_{k=1}^n \frac{\partial^2 \ln f(X_k; \theta_0)}{\partial \theta^2}. \end{aligned}$$

Folosind legea numerelor mari avem că

$$\frac{1}{n} \sum_{k=1}^n \frac{\partial^2 \ln f(X_k; \theta_0)}{\partial \theta^2} \xrightarrow{\text{a.s.}} E \left[\frac{\partial^2 \ln f(X; \theta_0)}{\partial \theta^2} \right] = -I_1(\theta_0),$$

de unde

$$\frac{\partial^2 \ln g(X_1, \dots, X_n; \xi)}{\partial \theta^2} \cong -n I_1(\theta_0) = -I_n(\theta_0),$$

sau

$$-2 \ln \hat{\Lambda} \cong (\theta_0 - \hat{\Theta})^2 I_n(\theta_0).$$

Conform Observației 5.3.11, statistica $\frac{\hat{\Theta} - \theta_0}{\sqrt{I_n^{-1}(\theta_0)}}$ urmează legea normală $\mathcal{N}(0, 1)$, când $n \rightarrow \infty$, de unde avem că statistica $(\theta_0 - \hat{\Theta})^2 I_n(\theta_0)$ urmează legea χ^2 cu un grad de libertate, ceea ce trebuia arătat. \square

6.6 Testul χ^2 privind dispersia teoretică

Fie caracteristica X ce urmează legea normală $\mathcal{N}(m, \sigma)$, unde $\sigma^2 = \text{Var}(X)$ este necunoscută și $m \in \mathbb{R}$, de asemenea necunoscut.

Relativ la dispersia teoretică se face ipoteza nulă $H_0 : \sigma^2 = \sigma_0^2$, cu una din alternativele:

$$H_1 : \sigma^2 \neq \sigma_0^2, \text{ când obținem testul } \chi^2 \text{ bilateral},$$

$$H_1 : \sigma^2 > \sigma_0^2, \text{ când obținem testul } \chi^2 \text{ unilateral dreapta},$$

$$H_1 : \sigma^2 < \sigma_0^2, \text{ când obținem testul } \chi^2 \text{ unilateral stânga}.$$

Pentru verificarea ipotezei nule H_0 , cu una din alternativele H_1 precizate, se consideră o selecție repetată de volum n , cu datele de selecție x_1, x_2, \dots, x_n și corespunzător variabilele de selecție X_1, X_2, \dots, X_n . Conform Proprietății 4.2.25, avem că statistica

$$\chi^2 = \frac{1}{\sigma^2} \sum_{k=1}^n (X_k - \bar{X})^2 = \frac{(n-1)\bar{\sigma}^2}{\sigma^2},$$

urmează legea χ^2 cu $n-1$ grade de libertate.

Folosind un nivel de semnificație $\alpha \in (0, 1)$ dat, se poate determina un interval numeric (χ_1^2, χ_2^2) astfel încât

$$P(\chi^2 \in (\chi_1^2, \chi_2^2) \mid H_0) = F_{n-1}(\chi_2^2) - F_{n-1}(\chi_1^2) = 1 - \alpha,$$

unde

$$F_m(x) = \frac{1}{2^{\frac{m}{2}} \Gamma(\frac{m}{2})} \int_0^x t^{\frac{m}{2}-1} e^{-\frac{t}{2}} dt, \quad x > 0,$$

este funcția de repartiție pentru legea χ^2 cu m grade de libertate (tabelată pentru anumite valori în *Anexa III*).

Intervalul numeric (χ_1^2, χ_2^2) pentru statistica χ^2 nu este determinat în mod unic din condiția de mai sus. În funcție de alternativa H_1 aleasă se consideră suplimentar:

$$\chi_1^2 = \chi_{n-1, \frac{\alpha}{2}}^2, \chi_2^2 = \chi_{n-1, 1-\frac{\alpha}{2}}^2, \quad \text{dacă } H_1 : \sigma^2 \neq \sigma_0^2,$$

$$\chi_1^2 = 0, \chi_2^2 = \chi_{n-1, 1-\alpha}^2, \quad \text{dacă } H_1 : \sigma^2 > \sigma_0^2,$$

$$\chi_1^2 = \chi_{n-1, \alpha}^2, \chi_2^2 = +\infty, \quad \text{dacă } H_1 : \sigma^2 < \sigma_0^2,$$

unde $\chi_{m, \gamma}^2$ este cuantila de ordin γ a legii χ^2 cu m grade de libertate, adică $F_m(\chi_{m, \gamma}^2) = \gamma$.

Cu ajutorul intervalului numeric (χ_1^2, χ_2^2) , astfel determinat, se consideră respectiv regiunile critice:

$$\begin{aligned}\mathcal{U} &= \left\{ (u_1, u_2, \dots, u_n) \in \mathbb{R}^n \mid \frac{1}{\sigma_0^2} \sum_{k=1}^n (u_k - \bar{u})^2 \notin \left(\chi_{n-1, \frac{\alpha}{2}}^2, \chi_{n-1, 1-\frac{\alpha}{2}}^2 \right) \right\}, \\ \mathcal{U} &= \left\{ (u_1, u_2, \dots, u_n) \in \mathbb{R}^n \mid \frac{1}{\sigma_0^2} \sum_{k=1}^n (u_k - \bar{u})^2 \geq \chi_{n-1, 1-\alpha}^2 \right\}, \\ \mathcal{U} &= \left\{ (u_1, u_2, \dots, u_n) \in \mathbb{R}^n \mid \frac{1}{\sigma_0^2} \sum_{k=1}^n (u_k - \bar{u})^2 \leq \chi_{n-1, \alpha}^2 \right\}.\end{aligned}$$

Ușor se verifică faptul că $P((X_1, X_2, \dots, X_n) \in \mathcal{U} \mid H_0) = \alpha$, iar cele trei moduri de definire a regiunii critice \mathcal{U} ne conduc respectiv la testul χ^2 bilateral, unilateral dreapta și unilateral stânga.

Odată construită regiunea critică \mathcal{U} , folosind datele de selecție, ipoteza nulă H_0 va fi admisă dacă $(x_1, x_2, \dots, x_n) \notin \mathcal{U}$, iar în caz contrar va fi respinsă. Remarcăm de asemenea că regiunea critică \mathcal{U} corespunde mulțimii complementare intervalului (χ_1^2, χ_2^2) .

Etapele aplicării testului χ^2

1. Se dau: α ; x_1, x_2, \dots, x_n ; $\sigma = \sigma_0$;
2. Se determină intervalul (χ_1^2, χ_2^2) astfel încât $F_{n-1}(\chi_2^2) - F_{n-1}(\chi_1^2) = 1 - \alpha$, (după cum s-a precizat înainte);
3. Se calculează $\chi^2 = \frac{1}{\sigma_0^2} \sum_{k=1}^n (x_k - \bar{x})^2$, unde $\bar{x} = \frac{1}{n} \sum_{k=1}^n x_k$;
4. Concluzia: dacă $\chi^2 \in (\chi_1^2, \chi_2^2)$ ipoteza H_0 este admisă, în caz contrar este respinsă.

Exemplul 6.6.1. Se consideră caracteristica X ce urmează legea normală $\mathcal{N}(m, \sigma)$ cu parametrii $m \in \mathbb{R}$ și $\sigma > 0$ necunoscuți. Relativ la caracteristica X se consideră o selecție repetată de volum $n = 15$. Fie datele de selecție

24.03	25.92	26.92	23.04	25.44	26.31	26.13	25.97	24.44	24.82
24.17	25.53	23.78	26.06	24.29					

Vrem să verificăm ipoteza nulă $H_0 : \sigma = 1.8$, cu fiecare din următoarele alternative $H_1 : \sigma \neq 1.8$, $H_1 : \sigma > 1.8$, $H_1 : \sigma < 1.8$, când se consideră $\alpha = 0.05$. Vom calcula apoi valorile funcției putere pentru $\sigma = 1.5, 1.6, 1.7, 1.8, 1.9$, pentru fiecare din alternativele considerate.

Când $H_1 : \sigma \neq 1.8$ avem testul χ^2 bilateral. Se calculează cuantilele

$$\chi_{n-1, \frac{\alpha}{2}}^2 = \chi_{14, 0.025}^2 = 5.63 \quad \text{și} \quad \chi_{n-1, 1-\frac{\alpha}{2}}^2 = \chi_{14, 0.975}^2 = 26.12.$$

Pe de altă parte avem

$$\bar{x} = \frac{1}{n} \sum_{k=1}^n x_k = \frac{1}{15} (24.03 + 25.92 + \dots + 24.29) = 25.11,$$

$$\bar{\sigma}^2 = \frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x})^2 = \frac{1}{14} [(24.03 - 25.11)^2 + \dots + (24.29 - 25.11)^2] = 1.27.$$

Deoarece

$$\chi^2 = \frac{(n-1) \bar{\sigma}^2}{\sigma_0^2} = \frac{14 \cdot 1.27}{1.8^2} = 5.49 \notin (5.63, 26.12),$$

ipoteza H_0 este respinsă.

Când $H_1 : \sigma > 1.8$, avem testul χ^2 unilateral dreapta și se calculează cuantila $\chi_{n-1, 1-\alpha}^2 = \chi_{14, 0.95}^2 = 23.69$. Astfel, având în vedere că $\chi^2 = 5.49 < 23.69$, ipoteza nulă este admisă.

Când $H_1 : \sigma < 1.8$, avem testul χ^2 unilateral stânga. În acest caz se calculează cuantila $\chi_{n-1, \alpha}^2 = \chi_{14, 0.05}^2 = 6.57$ și deoarece $\chi^2 = 5.49 < 6.57$, ipoteza nulă este respinsă.

Când se consideră testul χ^2 bilateral, puterea se obține din

$$\begin{aligned} 1 - \pi(\sigma) = \beta &= P \left(\chi_{n-1, \frac{\alpha}{2}}^2 < \frac{(n-1) \bar{\sigma}^2}{\sigma_0^2} < \chi_{n-1, 1-\frac{\alpha}{2}}^2 \mid H_1 \right) \\ &= P \left(\frac{\sigma_0^2}{\sigma^2} \chi_{n-1, \frac{\alpha}{2}}^2 < \frac{(n-1) \bar{\sigma}^2}{\sigma^2} < \frac{\sigma_0^2}{\sigma^2} \chi_{n-1, 1-\frac{\alpha}{2}}^2 \mid H_1 \right), \end{aligned}$$

deci

$$\pi(\sigma) = 1 - F_{n-1} \left(\frac{\sigma_0^2}{\sigma^2} \chi_{n-1, 1-\frac{\alpha}{2}}^2 \right) + F_{n-1} \left(\frac{\sigma_0^2}{\sigma^2} \chi_{n-1, \frac{\alpha}{2}}^2 \right).$$

În cazul de față, obținem

$$\pi(1.5) = 1 - F_{14} \left(\frac{1.8^2}{1.5^2} \cdot 26.14 \right) + F_{14} \left(\frac{1.8^2}{1.5^2} \cdot 5.63 \right) = 0.117,$$

$$\pi(1.6) = 1 - F_{14} \left(\frac{1.8^2}{1.6^2} \cdot 26.14 \right) + F_{14} \left(\frac{1.8^2}{1.6^2} \cdot 5.63 \right) = 0.073,$$

$$\pi(1.7) = 1 - F_{14} \left(\frac{1.8^2}{1.7^2} \cdot 26.14 \right) + F_{14} \left(\frac{1.8^2}{1.7^2} \cdot 5.63 \right) = 0.052,$$

$$\begin{aligned}\pi(1.8) &= 1 - F_{14} \left(\frac{1.8^2}{1.8^2} \cdot 26.14 \right) + F_{14} \left(\frac{1.8^2}{1.8^2} \cdot 5.63 \right) = 0.050, \\ \pi(1.9) &= 1 - F_{14} \left(\frac{1.8^2}{1.9^2} \cdot 26.14 \right) + F_{14} \left(\frac{1.8^2}{1.9^2} \cdot 5.63 \right) = 0.068.\end{aligned}$$

Pentru testul unilateral dreapta, în mod analog, se ajunge la

$$\pi(\sigma) = 1 - F_{n-1} \left(\frac{\sigma_0^2}{\sigma^2} \chi_{n-1, 1-\alpha}^2 \right),$$

de unde

$$\begin{aligned}\pi(1.5) &= 1 - F_{14} \left(\frac{1.8^2}{1.5^2} \cdot 23.69 \right) = 0.002, \\ \pi(1.6) &= 1 - F_{14} \left(\frac{1.8^2}{1.6^2} \cdot 23.69 \right) = 0.008, \\ \pi(1.7) &= 1 - F_{14} \left(\frac{1.8^2}{1.7^2} \cdot 23.69 \right) = 0.022, \\ \pi(1.8) &= 1 - F_{14} \left(\frac{1.8^2}{1.8^2} \cdot 23.69 \right) = 0.050, \\ \pi(1.9) &= 1 - F_{14} \left(\frac{1.8^2}{1.9^2} \cdot 23.69 \right) = 0.095.\end{aligned}$$

Când se consideră testul χ^2 unilateral stânga, se obține

$$\pi(\sigma) = F_{n-1} \left(\frac{\sigma_0^2}{\sigma^2} \chi_{n-1, \alpha}^2 \right),$$

deci

$$\begin{aligned}\pi(1.5) &= F_{14} \left(\frac{1.8^2}{1.5^2} \cdot 6.57 \right) = 0.200, & \pi(1.8) &= F_{14} \left(\frac{1.8^2}{1.8^2} \cdot 6.57 \right) = 0.050, \\ \pi(1.6) &= F_{14} \left(\frac{1.8^2}{1.6^2} \cdot 6.57 \right) = 0.128, & \pi(1.9) &= F_{14} \left(\frac{1.8^2}{1.9^2} \cdot 6.57 \right) = 0.031. \\ \pi(1.7) &= F_{14} \left(\frac{1.8^2}{1.7^2} \cdot 6.57 \right) = 0.080,\end{aligned}$$

Programul 6.6.2. Programul următor, efectuează calculele cerute în exemplul precedent. În plus, reprezintă grafic funcția putere, cea din Figura 6.3, în cazul testului bilateral, marcând prin cerușele punctele curbei putere, care au fost determinate în exemplu.

```

clf
x=[24.03,25.92,26.92,23.04,25.44,26.31,...
   26.13,25.97,24.44,24.82,24.17,25.53,...
   23.78,26.06,24.29];
x2=14*var(x)/1.8^2;
c1=chi2inv(0.025,14); c2=chi2inv(0.975,14);
fprintf(' x2=%6.3f\n',x2)
fprintf('      sigma ne 1.8\n')
fprintf(' c1=%6.3f, c2=%6.3f\n',c1,c2)
c1=chi2inv(0,14); c2=chi2inv(0.95,14);
fprintf('      sigma gt 1.8\n')
fprintf(' c1=%6.3f, c2=%6.3f\n',c1,c2)
c1=chi2inv(0.05,14); c2=chi2inv(1,14);
fprintf('      sigma lt 1.8\n')
fprintf(' c1=%6.3f, c2=%6.3f\n',c1,c2)
s=1.2:0.01:2.4; s1=1.5:0.1:1.9;
c1=1.8^2*chi2inv(0.025,14); c2=1.8^2*chi2inv(0.975,14);
pib=1-chi2cdf(c2./s.^2,14)+chi2cdf(c1./s.^2,14);
pib1(1,:)=1-chi2cdf(c2./s1.^2,14)+chi2cdf(c1./s1.^2,14);
plot(s,pib,'k-',s1,pib1(1,:), 'o'), grid on
c1=1.8^2*chi2inv(0,14); c2=1.8^2*chi2inv(0.95,14);
pib=1-chi2cdf(c2./s.^2,14)+chi2cdf(c1./s.^2,14);
pib1(2,:)=1-chi2cdf(c2./s1.^2,14)+chi2cdf(c1./s1.^2,14);
c1=1.8^2*chi2inv(0.05,14); c2=1.8^2*chi2inv(1,14);
pib=1-chi2cdf(c2./s.^2,14)+chi2cdf(c1./s.^2,14);
pib1(3,:)=1-chi2cdf(c2./s1.^2,14)+chi2cdf(c1./s1.^2,14);
len=length(s1);
for i=1:len
    for j=1:3
        fprintf(' pi(%3.1f)=%5.3f      ',s1(i),pib1(j,i))
    end
    fprintf('\n')
end

```

În urma executării programului s-au obținut rezultatele:

```

x2= 5.446
      sigma ne 1.8
c1= 5.629, c2=26.119
      sigma gt 1.8
c1= 0.000, c2=23.685
      sigma lt 1.8
c1= 6.571, c2=    Inf
pi(1.5)=0.117      pi(1.5)=0.002      pi(1.5)=0.200
pi(1.6)=0.073      pi(1.6)=0.008      pi(1.6)=0.128
pi(1.7)=0.052      pi(1.7)=0.022      pi(1.7)=0.080
pi(1.8)=0.050      pi(1.8)=0.050      pi(1.8)=0.050
pi(1.9)=0.068      pi(1.9)=0.095      pi(1.9)=0.031

```

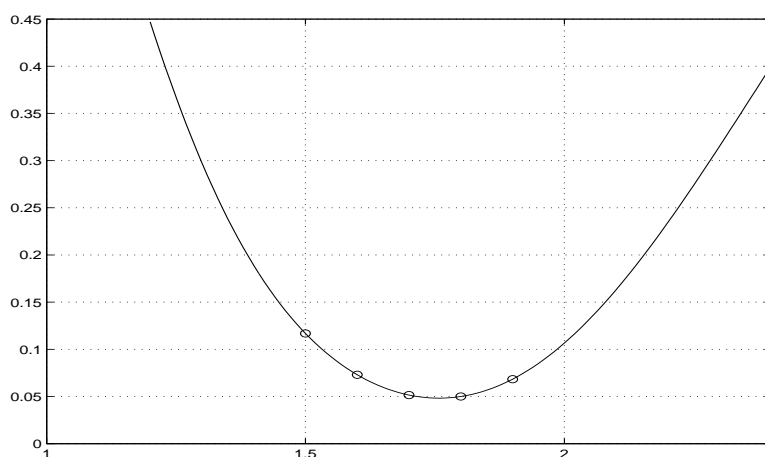


Figura 6.3: Curba funcției putere - cazul bilateral

Observația 6.6.3. Când caracteristica X nu urmează legea normală, pentru a verifica ipoteza nulă de forma dinainte, cu una din alternativele precizate, unde dispersia teoretică este $\sigma^2 = Var(X)$, se ține seama de faptul că statistica

$$S^2 = \frac{\bar{\sigma}^2 - \sigma^2}{\frac{\sigma^2 \sqrt{2}}{n}}, \quad \text{unde} \quad \bar{\sigma}^2 = \frac{1}{n-1} \sum_{k=1}^n (X_k - \bar{X})^2,$$

urmează legea normală $\mathcal{N}(0, 1)$, când $n \rightarrow \infty$.

De exemplu, dacă ipoteza alternativă este $H_1 : \sigma^2 \neq \sigma_0^2$, se va ajunge la regiunea critică

$$\mathcal{U} = \left\{ (u_1, u_2, \dots, u_n) \in \mathbb{R}^n \mid \sqrt{n} \frac{|\bar{\sigma}_u^2 - \sigma_0^2|}{\sigma_0^2 \sqrt{2}} \geq z_{1-\frac{\alpha}{2}} \right\},$$

unde $z_{1-\frac{\alpha}{2}}$ este cuantila de ordin $1 - \frac{\alpha}{2}$ de la legea normală $\mathcal{N}(0, 1)$.

Observația 6.6.4. Dacă se consideră statistica

$$H^2 = \frac{\bar{\sigma}^2}{\sigma^2} = \frac{1}{(n-1)\sigma^2} \sum_{k=1}^n (X_k - \bar{X})^2,$$

atunci între statisticile H^2 și χ^2 există relația $\chi^2 = (n-1)H^2$.

Deoarece se cunoaște legea de probabilitate pentru statistica χ^2 (legea χ^2 cu $n-1$ grade de libertate), se poate determina și legea de probabilitate a statisticii H^2 . Sunt manuale care conțin tabelată legea de probabilitate a lui H^2 . Descrierea testului χ^2 , folosind statistica H^2 , urmează aceeași cale ca cea când se folosește statistica χ^2 .

Observația 6.6.5. Când se cunoaște parametrul $m \in \mathbb{R}$ (ceea ce se întâmplă mai rar) se poate considera statistica

$$\chi^2 = \frac{1}{\sigma^2} \sum_{k=1}^n (X_k - m)^2 = \sum_{k=1}^n \left(\frac{X_k - m}{\sigma} \right)^2.$$

Deoarece $\frac{X_k - m}{\sigma}$ urmează legea normală $\mathcal{N}(0, 1)$, avem că statistica χ^2 urmează legea χ^2 cu n grade de libertate. Cele prezentate mai înainte pot fi rescrise cu această statistică.

6.7 Testul F (Fisher–Snedecor) pentru compararea dispersiilor

Se consideră două populații independente C' și C'' cercetate din punct de vedere al aceleași caracteristici. Această caracteristică este X' pentru C' și urmează legea normală $\mathcal{N}(m', \sigma')$ și respectiv X'' pentru C'' și urmează legea normală $\mathcal{N}(m'', \sigma'')$. Relativ la dispersiile teoretice ale celor două caracteristici se face ipoteza nulă compusă $H_0 : \sigma'^2 = \sigma''^2$, cu una din alternativele:

$$H_1 : \sigma'^2 \neq \sigma''^2, \text{ când obținem testul } F \text{ bilateral},$$

$$H_1 : \sigma'^2 > \sigma''^2, \text{ când obținem testul } F \text{ unilateral dreapta},$$

$$H_1 : \sigma'^2 < \sigma''^2, \text{ când obținem testul } F \text{ unilateral stânga}.$$

Pentru verificarea ipotezei nule H_0 , cu una din alternativele H_1 considerate, se efectuează câte o selecție repetată de volume respectiv n' și n'' din cele două populații C' și C'' . Notăm datele de selecție prin $x'_1, x'_2, \dots, x'_{n'}$ și respectiv $x''_1, x''_2, \dots, x''_{n''}$, cu variabilele de selecție $X'_1, X'_2, \dots, X'_{n'}$ și $X''_1, X''_2, \dots, X''_{n''}$. Conform Proprietății 4.2.30, statistica

$$F = \frac{\bar{\sigma}'^2}{\sigma'^2} \bigg/ \frac{\bar{\sigma}''^2}{\sigma''^2},$$

unde

$$\begin{aligned} \bar{\sigma}'^2 &= \frac{1}{n' - 1} \sum_{k=1}^{n'} (X'_k - \bar{X}')^2, & \bar{X}' &= \frac{1}{n'} \sum_{k=1}^{n'} X'_k, \\ \bar{\sigma}''^2 &= \frac{1}{n'' - 1} \sum_{k=1}^{n''} (X''_k - \bar{X}'')^2, & \bar{X}'' &= \frac{1}{n''} \sum_{k=1}^{n''} X''_k, \end{aligned}$$

urmează legea Fisher–Snedecor cu $m = n' - 1$ și $n = n'' - 1$ grade de libertate.

Pentru un nivel de semnificație $\alpha \in (0, 1)$ fixat, se poate determina un interval numeric (f_1, f_2) , astfel ca $P(F \in (f_1, f_2) \mid H_0) = F_{m,n}(f_2) - F_{m,n}(f_1) = 1 - \alpha$, unde

$$F_{m,n}(f) = \left(\frac{m}{n}\right)^{\frac{m}{2}} \frac{\Gamma\left(\frac{m+n}{2}\right)}{\Gamma\left(\frac{m}{2}\right)\Gamma\left(\frac{n}{2}\right)} \int_0^f x^{\frac{n}{2}-1} \left(1 + \frac{m}{n}x\right)^{-\frac{m+n}{2}} dx, \quad f > 0,$$

este funcția de repartiție pentru legea Fisher–Snedecor cu m și n grade de libertate (tabelată pentru anumite valori în *Anexa IV*). Intervalul de încredere (f_1, f_2) pentru statistica F nu este unic determinat. În funcție de alternativa H_1 aleasă, se consideră suplimentar:

$$f_1 = f_{m,n;\frac{\alpha}{2}}, \quad f_2 = f_{m,n;1-\frac{\alpha}{2}}, \quad \text{dacă } H_1: \sigma'^2 \neq \sigma''^2,$$

$$f_1 = 0, \quad f_2 = f_{m,n;1-\alpha}, \quad \text{dacă } H_1: \sigma'^2 > \sigma''^2,$$

$$f_1 = f_{m,n;\alpha}, \quad f_2 = +\infty, \quad \text{dacă } H_1: \sigma'^2 < \sigma''^2.$$

Cu ajutorul intervalului numeric (f_1, f_2) astfel determinat, se consideră respectiv regiunile critice

$$\mathcal{U} = \left\{ (u_1, \dots, u_{n'}; v_1, \dots, v_{n''}) \in \mathbb{R}^{n'+n''} \mid \frac{\bar{\sigma}_u^2}{\bar{\sigma}_v^2} \notin \left(f_{m,n;\frac{\alpha}{2}}, f_{m,n;1-\frac{\alpha}{2}}\right) \right\},$$

$$\mathcal{U} = \left\{ (u_1, \dots, u_{n'}; v_1, \dots, v_{n''}) \in \mathbb{R}^{n'+n''} \mid \frac{\bar{\sigma}_u^2}{\bar{\sigma}_v^2} \geq f_{m,n;1-\alpha} \right\},$$

$$\mathcal{U} = \left\{ (u_1, \dots, u_{n'}; v_1, \dots, v_{n''}) \in \mathbb{R}^{n'+n''} \mid \frac{\bar{\sigma}_u^2}{\bar{\sigma}_v^2} \leq f_{m,n;\alpha} \right\}.$$

S-au folosit notațiile

$$\bar{\sigma}_u^2 = \frac{1}{n' - 1} \sum_{k=1}^{n'} (u_k - \bar{u})^2, \quad \bar{u} = \frac{1}{n'} \sum_{k=1}^{n'} u_k,$$

$$\bar{\sigma}_v^2 = \frac{1}{n'' - 1} \sum_{k=1}^{n''} (v_k - \bar{v})^2, \quad \bar{v} = \frac{1}{n''} \sum_{k=1}^{n''} v_k.$$

Se verifică imediat că $P((X'_1, X'_2, \dots, X'_{n'}; X''_1, X''_2, \dots, X''_{n''}) \in \mathcal{U} \mid H_0) = \alpha$, iar cele trei alternative ne conduc la cele trei regiuni critice, care definesc respectiv testul F bilateral, unilateral dreapta și unilateral stânga.

Odată construită regiunea critică \mathcal{U} , folosind datele de selecție, ipoteza nulă va fi admisă dacă $(x'_1, x'_2, \dots, x'_{n'}; x''_1, x''_2, \dots, x''_{n''}) \notin \mathcal{U}$, iar în caz contrar va fi respinsă. Remarcăm de asemenea că regiunea critică \mathcal{U} corespunde mulțimii complementare intervalului (f_1, f_2) .

Etapele aplicării testului F

1. Se dau: α ; $x'_1, x'_2, \dots, x'_{n'}$; $x''_1, x''_2, \dots, x''_{n''}$;
2. Se determină intervalul (f_1, f_2) astfel încât $F_{m,n}(f_2) - F_{m,n}(f_1) = 1 - \alpha$ (după cum s-a prezentat înainte);
3. Se calculează $f = \frac{\bar{\sigma}'^2}{\bar{\sigma}''^2}$, unde

$$\bar{\sigma}'^2 = \frac{1}{n' - 1} \sum_{k=1}^{n'} (x'_k - \bar{x}')^2, \quad \bar{x}' = \frac{1}{n'} \sum_{k=1}^{n'} x'_k,$$

$$\bar{\sigma}''^2 = \frac{1}{n'' - 1} \sum_{k=1}^{n''} (x''_k - \bar{x}'')^2, \quad \bar{x}'' = \frac{1}{n''} \sum_{k=1}^{n''} x''_k;$$

4. Concluzia: dacă $f \in (f_1, f_2)$ ipoteza H_0 este admisă, în caz contrar este respinsă.

Observația 6.7.1. Dacă se notează prin $\theta = \frac{\sigma'}{\sigma''}$, atunci ipoteza nulă se poate rescrie $H_0 : \theta^2 = 1$, iar ipotezele alternative se scriu corespunzător $H_1 : \theta^2 \neq 1$, $H_1 : \theta^2 > 1$, respectiv $H_1 : \theta^2 < 1$, iar statistica F , se scrie sub forma $F = \frac{1}{\theta^2} \frac{\bar{\sigma}'^2}{\bar{\sigma}''^2}$.

Exemplul 6.7.2. Se cercetează precizia cu care două mașini produc conserve de același tip. Pentru aceasta, se consideră câte un eșantion din conservele produse de cele două mașini și se măsoară greutatea acestora. Fie X' greutatea în grame a unei conserve produsă de prima mașină, respectiv X'' pentru a doua mașină. Măsurătorile obținute sunt:

$X' :$ 1021 980 988 1017 1005 998 1014 985 995 1004 1030
 1015 995 1023 1008 1013
 $X'' :$ 1003 988 993 1013 1006 1002 1014 997 1002 1010 975

Considerând nivelul de semnificație $\alpha = 0.05$, vrem să verificăm ipoteza nulă compusă $H_0 : \sigma' = \sigma''$, respectiv cu fiecare din alternativele $H_1 : \sigma' \neq \sigma''$, $H_1 : \sigma' > \sigma''$, $H_1 : \sigma' < \sigma''$. Vom considera că cele două caracteristici X' și X'' sunt independente și că urmează legile normale $\mathcal{N}(m', \sigma')$ și respectiv $\mathcal{N}(m'', \sigma'')$.

Se consideră statistica

$$F = \frac{\bar{\sigma}'^2}{\sigma'^2} \bigg/ \frac{\bar{\sigma}''^2}{\sigma''^2},$$

ce urmează legea Fisher–Snedecor cu $m = n' - 1$ și $n = n'' - 1$ grade de libertate, care ne conduce la testul F . Pentru aceasta avem

$$\bar{x}' = \frac{1}{n'} \sum_{k=1}^{n'} x'_k = \frac{1}{16} (1021 + 980 + \cdots + 1013) = 1005.7,$$

$$\bar{x}'' = \frac{1}{n''} \sum_{k=1}^{n''} x''_k = \frac{1}{11} (1003 + 988 + \cdots + 975) = 1000.3,$$

$$\begin{aligned} \bar{\sigma}'^2 &= \frac{1}{n' - 1} \sum_{k=1}^{n'} (x'_k - \bar{x}')^2 = \frac{1}{15} [(1021 - 1005.7)^2 + \cdots + (1013 - 1005.7)^2] \\ &= 210.63, \end{aligned}$$

$$\begin{aligned} \bar{\sigma}''^2 &= \frac{1}{n'' - 1} \sum_{k=1}^{n''} (x''_k - \bar{x}'')^2 = \frac{1}{10} [(1003 - 1000.3)^2 + \cdots + (975 - 1000.3)^2] \\ &= 134.42. \end{aligned}$$

Valoarea calculată a statisticii F este

$$f = \frac{\bar{\sigma}'^2}{\bar{\sigma}''^2} = \frac{210.63}{134.42} = 1.57.$$

De asemenea, avem calculate, conform *Anexei IV*, cuantilele

$$f_{m,n;\frac{\alpha}{2}} = \frac{1}{f_{n,m;1-\frac{\alpha}{2}}} = \frac{1}{f_{10,15;0.975}} = \frac{1}{3.06} = 0.33,$$

$$f_{m,n;1-\frac{\alpha}{2}} = f_{15,10;0.975} = 3.52, \quad f_{m,n;1-\alpha} = f_{15,10;0.95} = 2.85,$$

$$f_{m,n;\alpha} = \frac{1}{f_{n,m;1-\alpha}} = \frac{1}{f_{10,15;0.95}} = \frac{1}{2.54} = 0.39.$$

Când se consideră alternativa $H_1 : \sigma' \neq \sigma''$, avem testul F bilateral. Deoarece

$$f = 1.57 \in (0.33, 3.52) = \left(f_{m,n;\frac{\alpha}{2}}, f_{m,n;1-\frac{\alpha}{2}} \right),$$

ipoteza H_0 este admisă.

Pentru alternativa $H_1 : \sigma' > \sigma''$, avem testul F unilateral dreapta. Astfel că

$$f = 1.57 < 2.85 = f_{m,n;1-\alpha},$$

prin urmare ipoteza H_0 este admisă.

De asemenea, pentru alternativa $H_1 : \sigma' < \sigma''$, se obține testul F unilateral stânga. Deoarece avem $f = 1.57 > 0.39 = f_{m,n;\alpha}$, ipoteza H_0 este admisă.

Programul 6.7.3. Programul Matlab, care urmează, efectuează calculele din exemplul precedent, după care afișează valoarea f a statisticii F , precum și intervalele numerice (f_1, f_2) , pentru fiecare din cele trei alternative.

```
x1=[1021,980,988,1017,1005,998,1014,985,995,...
    1004,1030,1015,995,1023,1008,1013];
x2=[1003,988,993,1013,1006,1002,1014,997,...
    1002,1010,975];
f=var(x1)/var(x2);
fprintf('    Bilateral\n')
f1=finv(0.025,15,10); f2=finv(0.975,15,10);
fprintf(' f=%7.3f, (f1,f2)=(%5.3f,%5.3f)\n',f,f1,f2)
fprintf('    Unilateral dreapta\n')
f1=finv(0,15,10); f2=finv(0.95,15,10);
fprintf(' f=%7.3f, (f1,f2)=(%5.3f,%5.3f)\n',f,f1,f2)
fprintf('    Unilateral stanga\n')
f1=finv(0.05,15,10); f2=finv(1,15,10);
fprintf(' f=%7.3f, (f1,f2)=(%5.3f,%5.3f)\n',f,f1,f2)
```

În urma executării programului se obțin rezultatele:

```
Bilateral
f= 1.567, (f1,f2)=(0.327,3.522)
Unilateral dreapta
f= 1.567, (f1,f2)=(0.000,2.845)
Unilateral stanga
f= 1.567, (f1,f2)=(0.393, Inf)
```

6.8 Teste pentru compararea mediilor

Se consideră două populații independente C' și C'' , cercetate din punct de vedere al aceleiași caracteristici. Această caracteristică este X' pentru C' și urmează legea normală $\mathcal{N}(m', \sigma')$ și respectiv X'' pentru C'' și urmează legea normală $\mathcal{N}(m'', \sigma'')$.

Relativ la mediile teoretice ale celor două caracteristici independente, se face ipoteza nulă $H_0 : m' = m''$, cu una din alternativele:

$H_1 : m' \neq m''$, când obținem un *test bilateral*,

$H_1 : m' > m''$, când obținem un *test unilateral dreapta*,

$H_1 : m' < m''$, când obținem un *test unilateral stânga*.

Ca și în cazul testului F se consideră câte o selecție repetată de volum n' și respectiv n'' . Păstrăm în continuare notațiile de la testul F .

Distingem trei cazuri.

6.8.1 Dispersii cunoscute

Dispersiile σ'^2 și σ''^2 sunt cunoscute. În acest caz se consideră statistica

$$(6.8.1) \quad Z = \frac{(\bar{X}' - \bar{X}'') - (m' - m'')}{\sqrt{\frac{\sigma'^2}{n'} + \frac{\sigma''^2}{n''}}},$$

care urmează legea normală $\mathcal{N}(0, 1)$.

Se aplică prin urmare testul Z pentru compararea celor două medii teoretice.

Pentru nivelul de semnificație $\alpha \in (0, 1)$ dat, se obțin regiunile critice corespunzătoare celor trei alternative:

$$\begin{aligned} \mathcal{U} &= \left\{ (u_1, \dots, u_{n'}; v_1, \dots, v_{n''}) \in \mathbb{R}^{n'+n''} \mid \frac{|\bar{u} - \bar{v}|}{\sqrt{\frac{\sigma'^2}{n'} + \frac{\sigma''^2}{n''}}} \geq z_{1-\frac{\alpha}{2}} \right\}, \\ \mathcal{U} &= \left\{ (u_1, \dots, u_{n'}; v_1, \dots, v_{n''}) \in \mathbb{R}^{n'+n''} \mid \frac{\bar{u} - \bar{v}}{\sqrt{\frac{\sigma'^2}{n'} + \frac{\sigma''^2}{n''}}} \geq z_{1-\alpha} \right\}, \\ \mathcal{U} &= \left\{ (u_1, \dots, u_{n'}; v_1, \dots, v_{n''}) \in \mathbb{R}^{n'+n''} \mid \frac{\bar{u} - \bar{v}}{\sqrt{\frac{\sigma'^2}{n'} + \frac{\sigma''^2}{n''}}} \leq z_{\alpha} \right\}. \end{aligned}$$

Etapele aplicării testului Z

1. Se dau: α ; $x'_1, x'_2, \dots, x'_{n'}$; $x''_1, x''_2, \dots, x''_{n''}$; σ' , σ'' ;
2. Se determină intervalul (z_1, z_2) astfel încât $\Phi(z_2) - \Phi(z_1) = 1 - \alpha$, unde $\Phi(x)$ este funcția lui Laplace tabelată în *Anexa I*. Intervalul numeric (z_1, z_2) este respectiv pentru cele trei alternative considerate: $(-z_{1-\frac{\alpha}{2}}, z_{1-\frac{\alpha}{2}})$, $(-\infty, z_{1-\alpha})$, $(z_{\alpha}, +\infty)$;
3. Se calculează $z = \frac{\bar{x}' - \bar{x}''}{\sqrt{\frac{\sigma'^2}{n'} + \frac{\sigma''^2}{n''}}}$, unde $\bar{x}' = \frac{1}{n'} \sum_{k=1}^{n'} x'_k$, $\bar{x}'' = \frac{1}{n''} \sum_{k=1}^{n''} x''_k$;
4. Concluzia: dacă $z \in (z_1, z_2)$ ipoteza H_0 este admisă, în caz contrar este respinsă.

Exemplul 6.8.1. La o unitate de îmbuteliere a laptelui există două mașini care efectuează această operație în sticle de 1 litru. Pentru a cerceta reglajul de îmbuteliere la cele două mașini s-au efectuat două selecții relative la sticlele îmbuteliate de cele două mașini și s-au obținut datele de selecție:

x'_k (în ml)	990	995	1000	1005	1010
f'_k	7	9	11	8	5

x''_k (în ml)	985	990	995	1000	1005	1010
f''_k	5	5	6	7	6	4

Folosind nivelul de semnificație $\alpha = 0.01$, vrem să verificăm dacă mediile de umplere a sticlelor de către cele două mașini sunt aceleași, în cazul în care abaterile standard sunt $\sigma' = 6$ ml și $\sigma'' = 7.5$ ml.

Caracteristicile X' și X'' , ce reprezintă cantitatea de lapte (în ml) conținută de o sticlă îmbuteliată de prima mașină, respectiv de a doua mașină, se consideră ca urmând legile de probabilitate normale $\mathcal{N}(m', 6)$ și $\mathcal{N}(m'', 7.5)$.

Verificarea ipotezei nule $H_0 : m' = m''$, cu alternativa $H_1 : m' \neq m''$, se va face cu testul Z , deoarece sunt cunoscute abaterile standard.

Folosind nivelul de semnificație $\alpha = 0.01$, se determină din *Anexa I* valoarea

$$z_{1-\frac{\alpha}{2}} = z_{0.995}, \quad \text{astfel încât} \quad \Phi\left(z_{1-\frac{\alpha}{2}}\right) = \frac{1-\alpha}{2} = 0.495.$$

Anume, se obține că $z_{0.995} = 2.58$, care ne dă intervalul $(-2.58; 2.58)$ pentru statistica Z , dată prin formula (6.8.1).

Se calculează succesiv:

$$\bar{x}' = \frac{1}{n'} \sum_{k=1}^{n'} x'_k = \frac{1}{40} (7 \cdot 990 + 9 \cdot 995 + \dots + 5 \cdot 1010) = 999.375;$$

$$\bar{x}'' = \frac{1}{n''} \sum_{k=1}^{n''} x''_k = \frac{1}{33} (5 \cdot 985 + 5 \cdot 990 + \dots + 4 \cdot 1010) = 997.424;$$

$$z = \frac{\bar{x}' - \bar{x}''}{\sqrt{\frac{\sigma'^2}{n'} + \frac{\sigma''^2}{n''}}} = \frac{999.375 - 997.424}{\sqrt{\frac{36}{40} + \frac{56.25}{33}}} = \frac{1.951}{\sqrt{2.6046}} = 1.209.$$

Deoarece $z = 1.209 \in (-2.58; 2.58)$, rezultă că mediile de umplere a sticlelor nu diferă semnificativ pentru cele două mașini.

Programul 6.8.2. Programul Matlab, care urmează, efectuează calculele din exemplul precedent.

```
x1=[990*ones(1,7),995*ones(1,9),1000*ones(1,11),...
    1005*ones(1,8),1010*ones(1,5)];
x2=[985*ones(1,5),990*ones(1,5),995*ones(1,6),...
    1000*ones(1,7),1005*ones(1,6),1010*ones(1,4)];
ma1=mean(x1); ma2=mean(x2);
z=(ma1-ma2)/sqrt(6^2/40+7.5^2/33);
z2=norminv(0.995); z1=-z2;
fprintf(' z=%7.3f\n z1=%6.3f\n z2=%6.3f',z,z1,z2)
```

Rezultatele obținute în urma executării programului sunt:

z = 1.209
z1 = -2.576
z2 = 2.576

6.8.2 Dispersii egale necunoscute

Dispersiile σ'^2 și σ''^2 sunt necunoscute și $\sigma'^2 = \sigma''^2 = \sigma^2$. Se consideră statistica

$$(6.8.2) \quad T = \frac{(\bar{X}' - \bar{X}'') - (m' - m'')}{\sqrt{(n' - 1) \bar{\sigma}'^2 + (n'' - 1) \bar{\sigma}''^2}} \sqrt{\frac{n' + n'' - 2}{\frac{1}{n'} + \frac{1}{n''}}},$$

care urmează, conform Proprietății 4.2.28, legea Student cu $n = n' + n'' - 2$ grade de libertate.

În acest caz se va aplica testul T . Pentru nivelul de semnificație $\alpha \in (0, 1)$ dat, se obțin regiunile critice corespunzătoare celor trei alternative:

$$\begin{aligned} \mathcal{U} &= \left\{ (u_1, \dots, u_{n'}; v_1, \dots, v_{n''}) \in \mathbb{R}^{n' + n''} \mid K |\bar{u} - \bar{v}| \geq t_{n, 1 - \frac{\alpha}{2}} \right\}, \\ \mathcal{U} &= \left\{ (u_1, \dots, u_{n'}; v_1, \dots, v_{n''}) \in \mathbb{R}^{n' + n''} \mid K (\bar{u} - \bar{v}) \geq t_{n, 1 - \alpha} \right\}, \\ \mathcal{U} &= \left\{ (u_1, \dots, u_{n'}; v_1, \dots, v_{n''}) \in \mathbb{R}^{n' + n''} \mid K (\bar{u} - \bar{v}) \leq t_{n, \alpha} \right\}, \end{aligned}$$

unde

$$\begin{aligned} \bar{u} &= \frac{1}{n'} \sum_{k=1}^{n'} u_k, \quad \bar{v} = \frac{1}{n''} \sum_{k=1}^{n''} v_k, \quad K = \frac{1}{\sqrt{(n' - 1) \bar{\sigma}_u^2 + (n'' - 1) \bar{\sigma}_v^2}} \sqrt{\frac{n}{\frac{1}{n'} + \frac{1}{n''}}}, \\ \bar{\sigma}_u^2 &= \frac{1}{n' - 1} \sum_{k=1}^{n'} (u_k - \bar{u})^2, \quad \bar{\sigma}_v^2 = \frac{1}{n'' - 1} \sum_{k=1}^{n''} (v_k - \bar{v})^2. \end{aligned}$$

Etapele aplicării testului T

1. Se dau: α ; $x'_1, x'_2, \dots, x'_{n'}; x''_1, x''_2, \dots, x''_{n''}$;
2. Se determină intervalul (t_1, t_2) astfel încât $F_n(t_2) - F_n(t_1) = 1 - \alpha$, unde $F_n(x)$ este funcția de repartiție pentru legea Student cu $n = n' + n'' - 2$ grade de libertate, iar intervalul numeric (t_1, t_2) este respectiv, pentru cele trei alternative considerate, $(-t_{n, 1 - \frac{\alpha}{2}}; t_{n, 1 - \frac{\alpha}{2}})$, $(-\infty; t_{n, 1 - \alpha})$, $(t_{n, \alpha}; +\infty)$;

3. Se calculează $t = \frac{\bar{x}' - \bar{x}''}{\sqrt{(n' - 1) \bar{\sigma}'^2 + (n'' - 1) \bar{\sigma}''^2} \sqrt{\frac{1}{n'} + \frac{1}{n''}}}$, unde

$$\bar{\sigma}'^2 = \frac{1}{n' - 1} \sum_{k=1}^{n'} (x'_k - \bar{x}')^2, \quad \bar{x}' = \frac{1}{n'} \sum_{k=1}^{n'} x'_k,$$

$$\bar{\sigma}''^2 = \frac{1}{n'' - 1} \sum_{k=1}^{n''} (x''_k - \bar{x}'')^2, \quad \bar{x}'' = \frac{1}{n''} \sum_{k=1}^{n''} x''_k;$$

4. Concluzia: dacă $t \in (t_1, t_2)$ ipoteza H_0 este admisă, în caz contrar este respinsă.

Exemplul 6.8.3. Se cercetează două loturi de ulei pentru automobile, din punct de vedere al vâscozității, obținându-se datele de selecție

x'_k	10.27	10.28	10.29	10.30	10.32
f'_k	3	2	1	1	1

pentru primul lot, respectiv

x''_k	10.26	10.27	10.29	10.30	10.31
f''_k	2	1	1	1	3

pentru al doilea lot.

Analizele făcându-se cu același aparat, se consideră că abaterea standard sunt aceleași. Considerând nivelul de semnificație $\alpha = 0.05$, să se verifice dacă mediile de vâscozitate pentru cele două loturi nu diferă semnificativ.

Caracteristicile X' și X'' , ce reprezintă vâscozitățile pentru cele două loturi de ulei, se consideră că urmează fiecare legea normală, respectiv $\mathcal{N}(m', \sigma)$ și $\mathcal{N}(m'', \sigma)$, cu abaterea standard $\sigma > 0$ necunoscută.

Verificarea ipotezei nule $H_0 : m' = m''$, cu alternativa $H_1 : m' \neq m''$, se va face cu testul T , deoarece abaterea standard σ este necunoscută.

Folosind nivelul de semnificație $\alpha = 0.05$, se determină, din *Anexa II*, valoarea $t_{n, 1-\frac{\alpha}{2}}$, astfel încât $F_n(t_{n, 1-\frac{\alpha}{2}}) = 1 - \frac{\alpha}{2}$, unde numărul gradelor de libertate este $n = n' + n'' - 2 = 8 + 8 - 2 = 14$. Adică, se determină $t_{14; 0.975}$ astfel încât $F_{14}(t_{14; 0.975}) = 0.975$, obținându-se $t_{14; 0.975} = 2.145$. În acest mod, s-a obținut intervalul $(-2.145; 2.145)$ pentru statistica dată prin formula (6.8.2), care urmează legea Student cu $n = n' + n'' - 2$ grade de libertate.

Se calculează pe rând

$$\bar{x}' = \frac{1}{n'} \sum_{k=1}^{n'} x'_k = \frac{1}{8} (3 \cdot 10.27 + 2 \cdot 10.28 + \dots + 1 \cdot 10.32) = 10.285;$$

$$\begin{aligned}\bar{x}'' &= \frac{1}{n''} \sum_{k=1}^{n''} x_k'' = \frac{1}{8} (2 \cdot 10.26 + 1 \cdot 10.27 + \dots + 3 \cdot 10.31) = 10.289; \\ \bar{\sigma}'^2 &= \frac{1}{n' - 1} \sum_{k=1}^{n'} (x_k' - \bar{x}')^2 = 3.14 \cdot 10^{-4}; \\ \bar{\sigma}''^2 &= \frac{1}{n'' - 1} \sum_{k=1}^{n''} (x_k'' - \bar{x}'')^2 = 4.98 \cdot 10^{-4}; \\ t &= \frac{\bar{x}' - \bar{x}''}{\sqrt{(n' - 1) \bar{\sigma}'^2 + (n'' - 1) \bar{\sigma}''^2}} \cdot \sqrt{\frac{n' + n'' - 2}{\frac{1}{n'} + \frac{1}{n''}}} \\ &= \frac{10.285 - 10.289}{\sqrt{(22.001 + 34.881) \cdot 10^{-4}}} \sqrt{\frac{14}{\frac{1}{8} + \frac{1}{8}}} = \frac{-4 \cdot 10^{-1}}{\sqrt{56.882}} \cdot \sqrt{56} = -0.37.\end{aligned}$$

Deoarece $t = -0.37 \in (-2.145; 2.145)$, rezultă că vâscozitățile medii ale celor două loturi de ulei nu diferă semnificativ.

Programul 6.8.4. Programul Matlab următor efectuează calculele de mai sus, după care afișează valoarea t a statisticii T , împreună cu capetele intervalului (t_1, t_2) .

```
x1=[10.27*ones(1,3),10.28*ones(1,2),10.29,10.3,10.32];
x2=[10.26,10.26,10.27,10.29,10.3,10.31*ones(1,3)];
ma1=mean(x1); ma2=mean(x2);
v1=var(x1); v2=var(x2);
t=(ma1-ma2)/sqrt(7*v1+7*v2)*sqrt(14/(1/8+1/8));
t2=tinvt(0.975,14); t1=-t2;
fprintf(' t=%7.3f\n t1=%6.3f\n t2=%6.3f',t,t1,t2)
```

În urma executării programului, se obțin rezultatele:

```
t= -0.372
t1=-2.145
t2= 2.145
```

6.8.3 Funcția `ttest2`

Sistemul Matlab, prin *Statistics toolbox*, dispune de funcția `ttest2`, cu aplicabilitate la testul T pentru compararea a două medii, când dispersiile sunt egale, dar necunoscute. Apelarea acestei funcții se face prin:

```
h=ttest2(x,y)
h=ttest2(x,y,alpha)
h=ttest2(x,y,alpha,tail)
[h,c,ci]=ttest(x,y,alpha,tail)
```

În urma executării acestor instrucțiuni, se efectuează testul T asupra datelor conținute în vectorii x și y , folosind nivelul de semnificație α , care are valoarea implicită $\alpha=0.05$.

Parametrul `tail` specifică una din cele trei alternative, care conduc la testul bilateral (`tail=0`, implicit), unilateral dreapta (`tail=1`) și unilateral stânga (`tail=-1`). Dacă $h=1$, atunci ipoteza nulă va fi respinsă, respectiv dacă $h=0$, ipoteza nu poate fi respinsă.

Ultima formă de apel permite, de asemenea, obținerea valorii critice c , precum și a intervalului de încredere pentru diferența mediilor teoretice, corespunzător probabilității de încredere $1-\alpha$, obținut în vectorul cu două componente ci . Valoarea critică c are aceeași semnificație cu cea precizată la Observația 6.4.2 și se calculează cu formulele prezentate acolo.

Programul 6.8.5. Programul Matlab, ce urmează, aplică testul T , pentru datele din Exemplul 6.8.3. Mai mult, se consideră și testele unilateral dreapta și unilateral stânga, precum și obținerea intervalelor de încredere, când $\alpha=0.05$.

```
x=[10.27*ones(1,3),10.28*ones(1,2),10.29,10.3,10.32];
y=[10.26*ones(1,2),10.27,10.29,10.3,10.31*ones(1,3)];
fprintf(' h      c      ci      \n')
fprintf(' _____ \n')
for i=-1:1
    [h,c,ci]=ttest2(x,y,0.05,i);
    fprintf(' %d %5.4f (%4.2f,%4.2f) \n',h,c,ci)
end
```

În urma executării programului, se obțin rezultatele:

h	c	ci
0	0.3577	(-Inf, 0.01)
0	0.7154	(-0.03, 0.02)
0	0.6423	(-0.02, Inf)

Se observă că pentru `tail=-1` și `tail=1`, se construiesc intervale nemărginite, respectiv la stânga și la dreapta.

6.8.4 Dispersii diferite necunoscute

Dispersiile σ'^2 și σ''^2 sunt necunoscute și diferite. Se va considera statistica

$$(6.8.3) \quad T = \frac{(\bar{X}' - \bar{X}'') - (m' - m'')}{\sqrt{\frac{\hat{\sigma}'^2}{n'} + \frac{\hat{\sigma}''^2}{n''}}},$$

care urmează legea Student cu n grade de libertate. Numărul n al gradelor de libertate se calculează cu formula

$$(6.8.4) \quad \frac{1}{n} = \frac{c^2}{n' - 1} + \frac{(1 - c)^2}{n'' - 1}, \quad \text{unde} \quad c = \frac{\bar{\sigma}'^2}{\bar{\sigma}'^2 + \bar{\sigma}''^2}.$$

S-a ajuns la testul T , care pentru nivelul de semnificație $\alpha \in (0, 1)$ dat, conduce la regiunile critice corespunzătoare alternativelor considerate, respectiv:

$$\begin{aligned} \mathcal{U} &= \left\{ (u_1, \dots, u_{n'}; v_1, \dots, v_{n''}) \in \mathbb{R}^{n' + n''} \mid \frac{|\bar{u} - \bar{v}|}{\sqrt{\frac{\bar{\sigma}_u^2}{n'} + \frac{\bar{\sigma}_v^2}{n''}}} \geq t_{n, 1 - \frac{\alpha}{2}} \right\}, \\ \mathcal{U} &= \left\{ (u_1, \dots, u_{n'}; v_1, \dots, v_{n''}) \in \mathbb{R}^{n' + n''} \mid \frac{\bar{u} - \bar{v}}{\sqrt{\frac{\bar{\sigma}_u^2}{n'} + \frac{\bar{\sigma}_v^2}{n''}}} \geq t_{n, 1 - \alpha} \right\}, \\ \mathcal{U} &= \left\{ (u_1, \dots, u_{n'}; v_1, \dots, v_{n''}) \in \mathbb{R}^{n' + n''} \mid \frac{\bar{u} - \bar{v}}{\sqrt{\frac{\bar{\sigma}_u^2}{n'} + \frac{\bar{\sigma}_v^2}{n''}}} \leq t_{n, \alpha} \right\}. \end{aligned}$$

Etapele aplicării testului T

1. Se dau: α ; $x'_1, x'_2, \dots, x'_{n'}$; $x''_1, x''_2, \dots, x''_{n''}$;
2. Se determină intervalul (t_1, t_2) astfel încât $F_n(t_2) - F_n(t_1) = 1 - \alpha$, unde $F_n(x)$ este funcția de repartiție pentru legea Student cu n grade de libertate calculat cu formula (6.8.4), iar intervalul numeric (t_1, t_2) este respectiv, pentru cele trei alternative considerate, $(-t_{n, 1 - \frac{\alpha}{2}}; t_{n, 1 - \frac{\alpha}{2}})$, $(-\infty; t_{n, 1 - \alpha})$, $(t_{n, \alpha}; +\infty)$;

3. Se calculează $t = \frac{\bar{x}' - \bar{x}''}{\sqrt{\frac{\bar{\sigma}'^2}{n'} + \frac{\bar{\sigma}''^2}{n''}}}$, unde

$$\begin{aligned} \bar{\sigma}'^2 &= \frac{1}{n' - 1} \sum_{k=1}^{n'} (x'_k - \bar{x}')^2, & \bar{x}' &= \frac{1}{n'} \sum_{k=1}^{n'} x'_k, \\ \bar{\sigma}''^2 &= \frac{1}{n'' - 1} \sum_{k=1}^{n''} (x''_k - \bar{x}'')^2, & \bar{x}'' &= \frac{1}{n''} \sum_{k=1}^{n''} x''_k; \end{aligned}$$

4. Concluzia: dacă $t \in (t_1, t_2)$ ipoteza H_0 este admisă, în caz contrar este respinsă.

Exemplul 6.8.6. Se cercetează capacitatea fiolelor farmaceutice de 100 ml, care provin de la două fabrici. În acest scop, se consideră câte o selecție pentru două loturi de fiole provenite respectiv de la cele două fabrici. Selecțiile obținute au distribuțiile empirice de selecție

$$X' \left(\begin{array}{cccccccccc} 100 & 101 & 102 & 103 & 104 & 105 & 106 & 107 & 108 & 109 \\ 1 & 1 & 2 & 3 & 4 & 5 & 4 & 1 & 3 & 1 \end{array} \right),$$

respectiv, pentru X'' : 110, 101, 112, 120, 117, 105, 109, 111, 118, 113, 106, 108, 115, 113, 112, 100, 116, 112, 114, 112.

Folosind nivelul de semnificație $\alpha = 0.02$, vom compara dispersiile celor două caracteristici. Apoi, pe baza acestui rezultat, vom compara mediile celor două caracteristici, utilizând același nivel de semnificație $\alpha = 0.02$.

Vom considera că cele două caracteristici X' și X'' sunt repartizate normal, respectiv $\mathcal{N}(m'; \sigma')$ și $\mathcal{N}(m''; \sigma'')$. Se poate aplica testul F , pentru compararea dispersiilor teoretice σ'^2 și σ''^2 .

Calculăm pe rând:

$$\bar{x}' = \frac{1}{n'} \sum_{k=1}^{n'} x'_k = \frac{1}{25} (1 \cdot 100 + \dots + 1 \cdot 109) = 104.76; \quad \bar{x}'' = \frac{1}{n''} \sum_{k=1}^{n''} x''_k = 111.2;$$

$$\bar{\sigma}'^2 = \frac{1}{n' - 1} \sum_{k=1}^{n'} (x'_k - \bar{x}')^2 = 5.19; \quad \bar{\sigma}''^2 = \frac{1}{n'' - 1} \sum_{k=1}^{n''} (x''_k - \bar{x}'')^2 = 27.537.$$

Deoarece $\bar{\sigma}'^2 < \bar{\sigma}''^2$, se consideră statistica $F = \frac{\bar{\sigma}''^2}{\bar{\sigma}'^2}$, care urmează legea Fisher-Snedecor cu $(m, n) = (n'' - 1, n' - 1) = (19, 24)$ grade de libertate.

Dacă se consideră ipoteza nulă $H_0 : \sigma'^2 = \sigma''^2$, cu alternativa $H_1 : \sigma'^2 \neq \sigma''^2$, avem că $f = \frac{\bar{\sigma}''^2}{\bar{\sigma}'^2} = \frac{27.537}{5.19} = 5.31$.

Pe de altă parte, pentru $\alpha = 0.02$, avem din *Anexa IV* că

$$f_{m,n;1-\frac{\alpha}{2}} = f_{19,24;0.99} = 2.76, \quad f_{m,n;\frac{\alpha}{2}} = f_{19,24;0.01} = \frac{1}{f_{24,19;0.99}} = \frac{1}{2.92} = 0.34.$$

În acest fel, am obținut intervalul $(0.34; 2.76)$, pentru statistica F .

Deoarece $f = 5.31 \notin (0.34; 2.76)$, respingem ipoteza că $\sigma'^2 = \sigma''^2$.

Având în vedere că dispersiile teoretice σ'^2 și σ''^2 sunt necunoscute, iar conform rezultatului precedent diferă în mod semnificativ, folosim testul T pentru compararea mediilor m' și m'' . Statistica T , ce se consideră, în acest caz, este cea dată prin formula (6.8.3), care urmează legea Student cu n grade de libertate, unde n se calculează din relația (6.8.4).

Astfel, pentru determinarea lui n , avem succesiv

$$c = \frac{\frac{5.19}{25}}{\frac{5.19}{25} + \frac{27.537}{20}} = 0.131 \quad \text{și} \quad \frac{1}{n} = \frac{0.131^2}{24} + \frac{(1 - 0.131)^2}{19} = 0.0404604,$$

de unde $n = 25$. Folosind *Anexa II*, se obține că $t_{25;0.99} = 2.485$, prin urmare intervalul pentru statistica T este $(-2.485; 2.485)$.

Pe de altă parte, avem că

$$t = \frac{\bar{x}' - \bar{x}''}{\sqrt{\frac{\bar{\sigma}'^2}{n'} + \frac{\bar{\sigma}''^2}{n''}}} = \frac{104.76 - 111.2}{\sqrt{0.2076 + 1.3768}} = -\frac{6.44}{1.26} = -5.11.$$

Deoarece $t = -5.11 \notin (-2.485; 2.485)$, respingem ipoteza că mediile teoretice pentru fiolele produse de cele două fabrici nu diferă semnificativ.

Programul 6.8.7. Calculele din exemplul de mai sus, se pot efectua cu următorul program Matlab:

```
x1=[100,101,102*ones(1,2),103*ones(1,3),...
    104*ones(1,4),105*ones(1,5),106*ones(1,4),...
    107,108*ones(1,3),109];
x2=[110,101,112,120,117,105,109,111,118,113,...
    106,108,115,113,112,100,116,112,114,112];
f=var(x2)/var(x1);
c1=finv(0.01,19,24); c2=finv(0.99,19,24);
fprintf(' f=%6.2f, c1=%5.2f, c2=%5.2f\n',f,c1,c2)
ma1=mean(x1); ma2=mean(x2);
v1=var(x1); v2=var(x2); c=(v1/25)/(v1/25+v2/20);
n=c^2/24+(1-c)^2/19; n=ceil(1/n);
t=(ma1-ma2)/sqrt(v1/25+v2/20);
t2=tinv(0.99,n); t1=-t2;
fprintf(' t=%7.3f, t1=%6.3f, t2=%6.3f',t,t1,t2)
```

Rezultatele obținute, în urma executării programului, sunt:

```
f= 5.31, c1= 0.34, c2= 2.76
t= -5.116, t1=-2.485, t2= 2.485
```

Observația 6.8.8. Dacă se notează prin $\theta = m' - m''$, atunci ipoteza compusă nulă devine $H_0: \theta = 0$, iar ipotezele alternative se rescriu după cum urmează: $H_1: \theta \neq 0$, $H_1: \theta > 0$ și respectiv $H_1: \theta < 0$, iar statisticile care au fost date mai înainte se pot rescrie cu acest nou parametru necunoscut θ .

Observația 6.8.9. Când selecțiile sunt de volum mare, $n', n'' > 50$, pentru legi de probabilitate oarecari, respectiv $n', n'' > 20$, pentru legi de probabilitate normale, statistica

$$Z = \frac{(\bar{X}' - \bar{X}'') - (m' - m'')}{\sqrt{\frac{\bar{\sigma}'^2}{n'} + \frac{\bar{\sigma}''^2}{n''}}},$$

poate fi considerată ca urmând legea normală $\mathcal{N}(0, 1)$. Așadar în acest caz se poate aplica testul Z .

6.8.5 Observații perechi

Dacă cele două selecții sunt de tip pereche, adică variabilele de selecție sunt perechile (X'_k, X''_k) , $k = \overline{1, n}$, atunci se pot considera diferențele $D_k = X'_k - X''_k$, pentru care $E(D_k) = m = m' - m''$, oricare a fi $k = \overline{1, n}$.

Problema poate fi reformulată.

Se consideră o populație \mathcal{C} , pentru care se cercetează caracteristica D , care urmează legea normală $\mathcal{N}(m, \sigma)$, cu $\sigma > 0$ necunoscut.

Vrem să verificăm ipoteza nulă $H_0 : m = 0$, cu una din alternativele

$H_1 : m \neq 0$, când obținem *testul T bilateral*,

$H_1 : m > 0$, când obținem *testul T unilateral dreapta*,

$H_1 : m < 0$, când obținem *testul T unilateral stânga*.

Conform Proprietății 4.2.26, statistica

$$T = \frac{\bar{D} - m}{\frac{\bar{\sigma}_D}{\sqrt{n}}},$$

unde

$$\bar{D} = \frac{1}{n} \sum_{k=1}^n D_k = \sum_{k=1}^n (X'_k - X''_k), \quad \bar{\sigma}^2 = \frac{1}{n-1} \sum_{k=1}^n (D_k - \bar{D})^2,$$

urmează legea Student cu $n - 1$ grade de libertate.

Prin urmare, pentru nivelul de semnificație $\alpha \in (0, 1)$ dat, se poate determina intervalul numeric (t_1, t_2) astfel încât

$$P(T \in (t_1, t_2) \mid H_0) = F_{n-1}(t_2) - F_{n-1}(t_1) = 1 - \alpha,$$

unde

$$F_m(t) = \frac{\Gamma\left(\frac{m+1}{2}\right)}{\sqrt{m\pi}\Gamma\left(\frac{m}{2}\right)} \int_{-\infty}^t \left(1 + \frac{x^2}{m}\right)^{-\frac{m+1}{2}} dx, \quad t \in \mathbb{R},$$

este funcția de repartiție pentru legea Student cu m grade de libertate (tabelată în *Anexa II*, pentru anumite valori).

Intervalul numeric (t_1, t_2) pentru statistica T nu este determinat în mod unic din condiția de mai sus. În funcție de alternativa H_1 aleasă, se consideră suplimentar:

$$t_1 = -t_2, \quad t_2 = t_{n-1, 1-\frac{\alpha}{2}}, \text{ dacă } H_1 : m \neq 0,$$

$$t_1 = -\infty, \quad t_2 = t_{n-1, 1-\alpha}, \text{ dacă } H_1 : m > 0,$$

$$t_1 = t_{n-1, \alpha}, \quad t_2 = +\infty, \text{ dacă } H_1 : m < 0,$$

unde $t_{n-1, \gamma}$ este cuantila de ordin γ a legii Student cu $n - 1$ grade de libertate.

Etapele aplicării testului T

1. Se dau: α ; x'_1, x'_2, \dots, x'_n ; $x''_1, x''_2, \dots, x''_n$;
2. Se calculează intervalul (t_1, t_2) astfel încât $F_{n-1}(t_1) - F_{n-1}(t_2) = 1 - \alpha$, (după cum s-a prezentat înainte);
3. Se calculează

$$t = \frac{\bar{d}}{\frac{\bar{\sigma}_D}{\sqrt{n}}}, \quad \bar{d} = \frac{1}{n} \sum_{k=1}^n d_k = \frac{1}{n} \sum_{k=1}^n (x'_k - x''_k), \quad \bar{\sigma}_D^2 = \frac{1}{n-1} \sum_{k=1}^n (d_k - \bar{d})^2;$$

4. Concluzia: dacă $t \in (t_1, t_2)$ ipoteza H_0 este admisă, în caz contrar ipoteza este respinsă.

Programul 6.8.10. Vom scrie un program, care generează n vectori aleatori, care urmează legea normală bidimensională. Folosind testul T pentru perechi, vom verifica ipoteza nulă, că mediile celor două variabile sunt egale, precum și intervalul de încredere pentru diferența mediilor, respectiv când se aplică testele bilateral, unilateral dreapta și unilateral stânga. Nivelul de semnificație folosit este $\alpha=0.05$.

```
mu(1)=input('m1='); mu(2)=input('m2=');
v(1,1)=input('sigma1^2=');
v(2,2)=input('sigma2^2=');
v(1,2)=input('Cov(X,Y)='); v(2,1) =v(1,2);
if det(v) <= 0
    error('Matricea v nu e pozitiv definita!')
end
n=input('n='); Z=mvnrnd(mu,v,n);
d=diff(Z,1,2);
fprintf(' h      c      ci      \n')
fprintf(' _____\n')
for i=1:1
    [h,c,ci]=ttest(d,0,0.05,i);
    fprintf(' %d  %5.4f    (%4.2f,%4.2f) \n',h,c,ci)
end
```

Rezultatele obținute, pentru $\mu_1=\mu_2=10$, $v = \begin{pmatrix} 2 & 2 \\ 2 & 3 \end{pmatrix}$ și $n=30$, sunt

h	c	ci
0	0.6715	(-Inf, 0.35)
0	0.6570	(-0.26, 0.41)
0	0.3285	(-0.20, Inf)

Deoarece $h=0$, pentru fiecare alternativă în parte, testul T nu respinge ipoteza că cele două medii sunt egale. Acest lucru se vede și din faptul că valoarea critică satisface inegalitatea $c > \alpha = 0.05$.

Programul mai calculează intervalul de încredere pentru diferența celor două medii. Se observă că acesta diferă, în funcție de parametrul $tail = -1, 0, 1$. De asemenea, se observă că am considerat nivelul de semnificație $\alpha = 0.05$, dar acesta poate fi ușor modificat în apelul funcției `ttest`.

6.9 Testul χ^2 pentru concordanță

Se consideră colectivitatea \mathcal{C} cercetată din punct de vedere al caracteristicii X cantitativă sau calitativă. Fie k numărul claselor caracteristicii X și E_i evenimentul ca un individ luat la întâmplare din colectivitatea \mathcal{C} să aparțină clasei cu numărul de ordine i .

Notând $p_i = P(E_i)$, $i = \overline{1, k}$, atunci $\sum_{i=1}^k p_i = 1$.

Relativ la caracteristica X facem ipoteza nulă $H_0 : p_i = p_i^{(0)}$, $i = \overline{1, k}$, cu ipoteza alternativă H_1 : există i_0 astfel încât $p_{i_0} \neq p_{i_0}^{(0)}$.

Pentru a verifica această ipoteză se consideră o selecție repetată de volum n . Fie datele de selecție x_1, x_2, \dots, x_n . Folosind aceste date de selecție se obțin frecvențele absolute ale claselor caracteristicii X . Vom nota prin n_i frecvența absolută a clasei cu numărul de ordine i . Altfel spus, n_i numără de câte ori a apărut evenimentul E_i în selecția considerată.

Corespunzător frecvențelor absolute n_i , $i = \overline{1, k}$, avem variabilele aleatoare (de selecție) N_i , $i = \overline{1, k}$, ce iau aceste valori.

Așadar, pornind de la caracteristica X , care poate fi și calitativă, s-a ajuns la variabilele de selecție N_i , $i = \overline{1, k}$, care sunt componentele unui vector aleator k -dimensional $N = (N_1, N_2, \dots, N_k)$, ce urmează legea multinomială. Anume, avem:

$$P(N_1 = n_1, N_2 = n_2, \dots, N_k = n_k) = \frac{n!}{n_1! n_2! \dots n_k!} p_1^{n_1} p_2^{n_2} \dots p_k^{n_k},$$

unde $n = n_1 + n_2 + \dots + n_k$, $n_i = \overline{0, n}$, $i = \overline{1, k}$, $p_1 + p_2 + \dots + p_k = 1$.

Dacă privim ipoteza nulă H_0 , constatăm că aceasta se referă la parametrul unei legi multinomiale.

Proprietatea 6.9.1. *Cu notațiile dinainte avem că statistica*

$$\chi^2 = \sum_{i=1}^k \frac{(N_i - np_i)^2}{np_i},$$

urmează legea χ^2 cu $k - 1$ grade de libertate, când $n \rightarrow \infty$.

Demonstrație. Se pornește de la formula lui Stirling $n! \cong \sqrt{2\pi n} n^n e^{-n}$. Astfel, avem că

$$P(N_1 = n_1, N_2 = n_2, \dots, N_k = n_k) \cong \frac{\sqrt{2\pi n} n^n e^{-n}}{\prod_{i=1}^k \left(\sqrt{2\pi n_i} n_i^{n_i} e^{-n_i} \right)} p_1^{n_1} p_2^{n_2} \dots p_k^{n_k},$$

sau

$$P(N_1 = n_1, N_2 = n_2, \dots, N_k = n_k) \cong \prod_{i=1}^k \left(\frac{np_i}{n_i} \right)^{n_i + \frac{1}{2}} \cdot K,$$

unde K este o constantă pozitivă. Logaritmând această relație obținem

$$\ln P(N_1 = n_1, N_2 = n_2, \dots, N_k = n_k) \cong \ln K + \sum_{i=1}^k \left(n_i + \frac{1}{2} \right) \ln \frac{np_i}{n_i},$$

și notând

$$x_i = \frac{n_i - np_i}{\sqrt{np_i}}, \quad \text{adică} \quad \frac{n_i}{np_i} = 1 + \frac{x_i}{\sqrt{np_i}},$$

rezultă că

$$\ln P(N_1 = n_1, N_2 = n_2, \dots, N_k = n_k) \cong \ln K - \sum_{i=1}^k \left(n_i + \frac{1}{2} \right) \ln \left(1 + \frac{x_i}{\sqrt{np_i}} \right).$$

Păstrând primii doi termeni din dezvoltarea în serie a logaritmului avem

$$\ln \left(1 + \frac{x_i}{\sqrt{np_i}} \right) \cong \frac{x_i}{\sqrt{np_i}} - \frac{x_i^2}{2np_i},$$

deci

$$\begin{aligned} \ln P(N_1 = n_1, N_2 = n_2, \dots, N_k = n_k) &\cong \ln K - \sum_{i=1}^k \left(n_i + \frac{1}{2} \right) \left(\frac{x_i}{\sqrt{np_i}} - \frac{x_i^2}{2np_i} \right) \\ &= \ln K - \sum_{i=1}^k \left(np_i + x_i \sqrt{np_i} + \frac{1}{2} \right) \left(\frac{x_i}{\sqrt{np_i}} - \frac{x_i^2}{2np_i} \right) \\ &\cong \ln K - \sum_{i=1}^k \left(x_i \sqrt{np_i} + \frac{x_i^2}{2} \right). \end{aligned}$$

Dar avem că

$$\sum_{i=1}^k x_i \sqrt{np_i} = \sum_{i=1}^k (n_i - np_i) = n - n = 0,$$

deci

$$\ln P(N_1 = n_1, N_2 = n_2, \dots, N_k = n_k) \cong \ln K - \frac{1}{2} \sum_{i=1}^k x_i^2$$

sau

$$P(N_1 = n_1, N_2 = n_2, \dots, N_k = n_k) \cong K e^{-\frac{1}{2} \sum_{i=1}^k x_i^2}.$$

Dacă se notează $X_i = \frac{N_i - np_i}{\sqrt{np_i}}$, obținem

$$P(X_1 = x_1, X_2 = x_2, \dots, X_k = x_k) \cong K e^{-\frac{1}{2} \sum_{i=1}^k x_i^2},$$

adică vectorul aleator

$$(X_1, X_2, \dots, X_k) = \left(\frac{N_1 - np_1}{\sqrt{np_1}}, \frac{N_2 - np_2}{\sqrt{np_2}}, \dots, \frac{N_k - np_k}{\sqrt{np_k}} \right),$$

pentru $n \rightarrow \infty$, urmează o lege normală k dimensională degenerată, deoarece fiecare componentă X_i se poate exprima ca și combinație liniară a celorlalte componente.

Din teoria probabilităților se cunoaște că suma pătratelor componentelor unui vector aleator ce urmează legea normală și între care există o legătură liniară este o variabilă aleatoare ce urmează legea χ^2 cu numărul gradelor de libertate dat de numărul componentelor vectorului mai puțin unu. Ceea ce trebuie demonstrat. \square

Pe baza proprietății precedente, pentru nivelul de semnificație $\alpha \in (0, 1)$ dat, se poate determina cuantila $\chi_{k-1, 1-\alpha}^2$ astfel încât

$$F_{k-1}(\chi_{k-1, 1-\alpha}^2) = P(\chi^2 \leq \chi_{k-1, 1-\alpha}^2 \mid H_0) = 1 - \alpha,$$

unde F_{k-1} notează funcția de repartiție a legii χ^2 cu $k - 1$ grade de libertate.

Regiunea critică \mathcal{U} se poate scrie și în acest caz, anume

$$\mathcal{U} = \left\{ (u_1, u_2, \dots, u_k) \in \mathbb{R}^k \mid \sum_{i=1}^k \frac{(u_i - np_i^{(0)})^2}{np_i^{(0)}} \geq \chi_{k-1, 1-\alpha}^2, \sum_{i=1}^k u_i = n \right\}.$$

Astfel, pentru n mare, are loc relația $P((N_1, N_2, \dots, N_k) \in \mathcal{U} \mid H_0) = \alpha$, deci ipoteza H_0 va fi admisă dacă $(n_1, n_2, \dots, n_k) \notin \mathcal{U}$, respectiv va fi respinsă în caz contrar.

Etapele aplicării testului χ^2

1. Se dau: α ; x_1, x_2, \dots, x_n ; $p_i = p_i^{(0)}$, $i = \overline{1, k}$. (șirul datelor de selecție x_1, x_2, \dots, x_n este un șir având ca elemente evenimentele E_i , $i = \overline{1, k}$, din care se obțin frecvențele absolute n_1, n_2, \dots, n_k);
2. Se calculează $\chi_{k-1, 1-\alpha}^2$ astfel încât $F_{k-1}(\chi_{k-1, 1-\alpha}^2) = 1 - \alpha$;
3. Se calculează $x^2 = \sum_{i=1}^k \frac{(n_i - np_i^{(0)})^2}{np_i^{(0)}}$;
4. Concluzia: dacă $x^2 < \chi_{k-1, 1-\alpha}^2$ ipoteza H_0 este admisă, în caz contrar ipoteza este respinsă.

Observația 6.9.2. La utilizarea testului χ^2 privind concordanța trebuie ca să fie îndeplinite condițiile $np_i > 4$, când $k > 4$, iar dacă numărul claselor $k \leq 4$, atunci np_i să fie mult mai mare decât 4. Când aceste condiții nu sunt îndeplinite, se efectuează o regrupare a datelor de selecție.

Exemplul 6.9.3. S-a aruncat un zar de 60 de ori și s-au obținut următoarele rezultate:

Fața	1	2	3	4	5	6
Frecvența	15	7	4	11	6	17

Folosind nivelul de semnificație $\alpha = 0.05$, să verificăm dacă zarul respectiv este fals sau nu.

Se aplică testul χ^2 , privind parametrii legii multinomiale. Într-adevăr, evenimentul E_i reprezintă apariția feței cu numărul i , $i = \overline{1, 6}$. Se face ipoteza nulă

$$H_0 : p_i = \frac{1}{6}, \quad i = \overline{1, 6}, \quad \text{adică zarul nu este fals,}$$

cu alternativa

$$H_1 : \exists i_0 \text{ astfel încât } p_{i_0} \neq \frac{1}{6}, \quad \text{adică zarul este fals.}$$

Se calculează valoarea statisticii χ^2 , care are $k - 1 = 5$ grade de libertate, anume

$$x^2 = \sum_{i=1}^k \frac{(n_i - np_i^{(0)})^2}{np_i^{(0)}} = \frac{(15 - 60 \cdot \frac{1}{6})^2}{60 \cdot \frac{1}{6}} + \frac{(7 - 60 \cdot \frac{1}{6})^2}{60 \cdot \frac{1}{6}} + \dots + \frac{(17 - 60 \cdot \frac{1}{6})^2}{60 \cdot \frac{1}{6}} = 13.6.$$

Pe de altă parte, din *Anexa III*, avem cuantila $\chi_{k-1, 1-\alpha}^2 = \chi_{5, 0.95}^2 = 11.1$, deci rezultă $x^2 = 13.6 > 11.1 = \chi_{5, 0.95}^2$, ceea ce conduce la respingerea ipotezei nule, adică acceptăm ipoteza că zarul este fals.

Programul 6.9.4. Programul Matlab, care urmează, efectuează calculele din exemplul precedent:

```
x=1:6; f=[15,7,4,11,6,17]; p=1/6*ones(1,6);
x2=sum((f-60*p).^2./(60*p));
cuant=chi2inv(0.95,5);
fprintf(' x^2= %6.1f\n cuant= %5.2f',x2,cuant)
```

iar rezultatele obținute sunt:

```
x^2=      13.6
cuant= 11.07
```

Aplicația 6.9.5. (Testul χ^2 neparametric privind concordanța). Fie caracteristica X , care urmează legea de probabilitate cu funcția de repartiție F necunoscută. Relativ la legea de probabilitate, se face ipoteza nulă $H_0 : F = F_0$, cu ipoteza alternativă $H_1 : F \neq F_0$. Când considerăm testul χ^2 neparametric, funcția de repartiție F_0 nu depinde de nici un parametru necunoscut.

Dacă domeniul valorilor caracteristicii X este, să zicem, intervalul (a, b) și dacă vom considera clasele obținute prin punctele $a = a_0 < a_1 < \dots < a_k = b$, atunci apar parametrii necunoscuți

$$p_i = P(a_{i-1} < X \leq a_i) = F(a_i) - F(a_{i-1}), \quad i = \overline{1, k}.$$

De asemenea, evenimentul E_i va fi evenimentul ca un individ luat la întâmplare din colectivitatea cercetată să aparțină clasei $[a_{i-1}, a_i)$. Prin urmare, s-a ajuns la aplicarea testului χ^2 privind concordanța. Ipoteza nulă, mai sus precizată, devine în acest fel $H_0 : p_i = p_i^{(0)}, i = \overline{1, k}$, iar ipoteza alternativă se va rescrie H_1 : există i_0 astfel încât $p_{i_0} \neq p_{i_0}^{(0)}$, unde $p_i^{(0)} = P(a_{i-1} < X \leq a_i | H_0) = F_0(a_i) - F_0(a_{i-1})$.

Etapele aplicării testului χ^2 neparametric

1. Se dau: $\alpha; x_1, x_2, \dots, x_n; a_0, a_1, \dots, a_k; F_0$;
2. Se calculează frecvențele absolute $n_i, i = \overline{1, k}$, și probabilitățile

$$p_i^{(0)} = F_0(a_i) - F_0(a_{i-1}), \quad i = \overline{1, k};$$

3. Se calculează $\chi_{k-1, 1-\alpha}^2$ astfel încât $F_{k-1}(\chi_{k-1, 1-\alpha}^2) = 1 - \alpha$, F_{k-1} fiind funcția de repartiție a legii χ^2 cu $k - 1$ grade de libertate;

4. Se calculează $x^2 = \sum_{i=1}^k \frac{(n_i - np_i^{(0)})^2}{np_i^{(0)}}$;

5. Concluzia: dacă $x^2 < \chi_{k-1,1-\alpha}^2$ se admite ipoteza că legea de probabilitate a caracteristicii X este dată de funcția de repartiție F_0 , în caz contrar ipoteza este respinsă.

Exemplul 6.9.6. Fie caracteristica X ce reprezintă numărul fetelor dintr-o familie cu patru copii. Pentru verificarea ipotezei că X urmează o lege binomială de parametru $p = \frac{1}{2}$, s-a efectuat o selecție de volum $n = 32$. Distribuția selecției lui X este

$$X \begin{pmatrix} 0 & 1 & 2 & 3 & 4 \\ 4 & 10 & 8 & 7 & 3 \end{pmatrix}.$$

Folosind nivelul de semnificație $\alpha = 0.05$, să verificăm dacă X urmează legea binomială. Dacă X urmează legea binomială cu $p = \frac{1}{2}$, atunci are distribuția teoretică

$$X \begin{pmatrix} 0 & 1 & 2 & 3 & 4 \\ \frac{1}{16} & \frac{1}{4} & \frac{3}{8} & \frac{1}{4} & \frac{1}{16} \end{pmatrix}.$$

Ipoteza privind faptul că X urmează legea binomială, se scrie sub forma

$$H_0 : p_0 = p_4 = \frac{1}{16}, p_1 = p_3 = \frac{1}{4}, p_2 = \frac{3}{8}.$$

Se va aplica testul χ^2 neparametric, unde numărul gradelor de libertate este dat de $k - 1 = 5 - 1 = 4$. Valoarea calculată a caracteristicii χ^2 este

$$x^2 = \sum_{i=0}^4 \frac{(n_i - np_i^{(0)})^2}{np_i^{(0)}} = \frac{(4 - 32 \cdot \frac{1}{16})^2}{32 \cdot \frac{1}{16}} + \frac{(10 - 32 \cdot \frac{1}{4})^2}{32 \cdot \frac{1}{4}} + \dots + \frac{(3 - 32 \cdot \frac{1}{16})^2}{32 \cdot \frac{1}{16}} = 4.5.$$

Cum $x^2 = 4.5 < 9.49 = \chi_{4,0.95}^2 = \chi_{k-1,1-\alpha}^2$, avem că ipoteza nulă este admisă.

Programul 6.9.7. Prin executarea programului Matlab:

```
x=0:4; f=[4,10,8,7,3];
p=[1/16,1/4,3/8,1/4,1/16];
x2=sum((f-32*p).^2./(32*p));
cuant=chi2inv(0.95,4);
fprintf(' x^2= %6.1f\n cuant= %5.2f',x2,cuant)
```

se obține

```
x^2=      4.5
cuant=    9.49
```

Aplicația 6.9.8. (Testul χ^2 neparametric privind exponențialitatea). Relativ la caracteristica X se face ipoteza nulă $H_0 : F = F_0$, unde $F_0(x) = 1 - e^{-\frac{x}{\mu}}$, $x > 0$ și parametrul $\mu > 0$ cunoscut, adică X urmează legea exponențială.

Deoarece domeniul valorilor caracteristicii X este intervalul $(0, +\infty)$, se consideră clasele caracteristicii X date prin $0 = a_0 < a_1 < \dots < a_k = +\infty$. Astfel avem:

$$p_i^{(0)} = F_0(a_i) - F_0(a_{i-1}) = e^{-\frac{1}{\mu}a_{i-1}} - e^{-\frac{1}{\mu}a_i}, \quad i = \overline{1, k-1} \quad \text{și} \quad p_k^{(0)} = e^{-\frac{1}{\mu}a_{k-1}}.$$

Etapele aplicării testului

1. Se dau: $\alpha; x_1, x_2, \dots, x_n; a_1, a_2, \dots, a_{k-1}; \mu;$
2. Se calculează frecvențele absolute $n_i, i = \overline{1, k}$, și probabilitățile

$$p_1^{(0)} = 1 - e^{-\frac{1}{\mu}a_1}, \quad p_i^{(0)} = e^{-\frac{1}{\mu}a_{i-1}} - e^{-\frac{1}{\mu}a_i}, \quad i = \overline{2, k-1}, \quad p_k^{(0)} = e^{-\frac{1}{\mu}a_{k-1}};$$

3. Se calculează $\chi_{k-1, 1-\alpha}^2$ astfel încât $F_{k-1}(\chi_{k-1, 1-\alpha}^2) = 1 - \alpha;$

$$4. \text{ Se calculează } x^2 = \sum_{i=1}^k \frac{(n_i - np_i^{(0)})^2}{np_i^{(0)}};$$

5. Concluzia: dacă $x^2 < \chi_{k-1, 1-\alpha}^2$ se acceptă ipoteza că X urmează legea exponențială de parametru λ , în caz contrar ipoteza este respinsă.

Aplicația 6.9.9. (Testul χ^2 parametric privind concordanța). Se consideră caracteristica X care urmează legea de probabilitate cu funcția de repartiție F necunoscută. Relativ la legea de probabilitate se face ipoteza nulă $H_0 : F = F_0$ cu ipoteza alternativă $H_1 : F \neq F_0$. Față de cazul neparametric, funcția de repartiție F_0 se consideră că depinde de s parametri necunoscuți. Fie acești parametri $\theta_1, \theta_2, \dots, \theta_s$, adică $F_0 = F_0(x; \theta_1, \theta_2, \dots, \theta_s)$.

În acest caz, față de cazul neparametric, la început se estimează parametrii necunoscuți, folosind datele de selecție considerate, cu ajutorul metodei verosimilității maxime. Fie estimațiile de verosimilitate maximă ale parametrilor mai sus precizați, respectiv $\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_s$.

În continuare, cele expuse la cazul neparametric rămân neschimbate cu două observații. În primul rând, probabilitățile claselor considerate se calculează prin formula

$$p_i^{(0)} = P(a_{i-1} < X \leq a_i | H_0) = F_0(a_i; \hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_s) - F_0(a_{i-1}; \hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_s),$$

pentru $i = \overline{1, k}$. În al doilea rând, legea χ^2 , în acest caz, are $k - s - 1$ grade de libertate.

Etapele aplicării testului χ^2 parametric

1. Se dau: α ; x_1, x_2, \dots, x_n ; a_0, a_1, \dots, a_k ; $F_0 = F_0(x; \theta_1, \theta_2, \dots, \theta_s)$;
2. Se determină estimațiile de verosimilitate maximă $\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_s$ pentru parametrii $\theta_1, \theta_2, \dots, \theta_s$;
3. Se calculează frecvențele absolute $n_i, i = \overline{1, k}$, și probabilitățile

$$p_i^{(0)} = F_0(a_i; \hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_s) - F_0(a_{i-1}; \hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_s), \quad i = \overline{1, k};$$

4. Se calculează $\chi_{k-s-1, 1-\alpha}^2$ astfel încât $F_{k-s-1}(\chi_{k-s-1, 1-\alpha}^2) = 1 - \alpha$;

$$5. \text{ Se calculează } x^2 = \sum_{i=1}^k \frac{(n_i - np_i^{(0)})^2}{np_i^{(0)}};$$

6. Concluzia: dacă $x^2 < \chi_{k-s-1, 1-\alpha}^2$ se acceptă ipoteza că X urmează legea de probabilitate dată prin funcția de repartiție $F_0 = F_0(x; \hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_s)$, în caz contrar ipoteza este respinsă.

Exemplul 6.9.10. O substanță radioactivă este observată în 2608 intervale de timp de lungimi egale (7.5 secunde). Pentru fiecare interval de timp este înregistrat numărul particulelor emise. Rezultatele obținute sunt:

Nr. part.	0	1	2	3	4	5	6	7	8	9	≥ 10
n_i	57	203	383	525	532	408	273	139	45	27	16

Notând cu X caracteristica ce reprezintă numărul de particule emise într-un interval de timp de lungime 7.5 secunde, vrem să verificăm dacă X urmează legea lui Poisson, folosind nivelul de semnificație $\alpha = 0.01$.

Folosim testul χ^2 parametric privind concordanța, deoarece parametrul legii lui Poisson, $\lambda = E(X)$, este necunoscut. Așadar, numărul gradelor de libertate va fi $k - s - 1 = 11 - 1 - 1 = 9$. Avem că estimatorul de verosimilitate maximă pentru λ este media de selecție. Prin urmare avem

$$\hat{\lambda} = \frac{1}{n} \sum_{i=0}^k n_i x_i = \frac{1}{2608} (0 \cdot 57 + 1 \cdot 203 + \dots + 10 \cdot 16) = 3.87.$$

Se calculează apoi probabilitățile de la legea lui Poisson, cu parametrul $\hat{\lambda} = 3.87$, adică $p_i^{(0)} = \frac{\hat{\lambda}^i}{i!} e^{-\hat{\lambda}}, i = 0, 1, 2, \dots$. Astfel se obțin valorile:

i	0	1	2	3	4	5	6	7	8	9	10
$p_i^{(0)}$	0.021	0.081	0.156	0.202	0.195	0.151	0.097	0.054	0.026	0.011	0.006

de unde

$$x^2 = \sum_{i=0}^{10} \frac{(n_i - np_i^{(0)})^2}{np_i^{(0)}} = 14.9.$$

Pe de altă parte, folosind *Anexa III*, avem cuantila $\chi_{9,0.99}^2 = 21.7$ și prin urmare $x^2 = 14.9 < 21.7 = \chi_{9,0.99}^2$, deci ipoteza că X urmează legea lui Poisson este acceptată.

Programul 6.9.11. Ilustrarea calculelor este făcută prin programul Matlab

```
x=0:10;
f=[57,203,383,525,532,408,273,139,45,27,16];
la=sum(f.*x)/2608; p=poisspdf(x,la);
x2=sum((f-2608*p).^2./(2608*p));
cuant=chi2inv(0.99,9);
fprintf(' x^2= %6.1f\n cuant= %5.2f',x2,cuant)
```

care prin executare conduce la următoarele rezultate:

```
x^2=      14.9
cuant=  21.67
```

Aplicația 6.9.12. (Testul χ^2 parametric privind normalitatea). Relativ la caracteristica X se consideră ipoteza nulă $H_0 : F = F_0$, unde

$$F_0(x) = F_0(x; m, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{(t-m)^2}{2\sigma^2}} dt, \quad x \in \mathbb{R},$$

$m \in \mathbb{R}$ și $\sigma > 0$ sunt parametri necunoscuți, adică faptul că X urmează legea normală cu cei doi parametri necunoscuți.

Folosind metoda verosimilității maxime se determină, pe baza datelor de selecție x_1, x_2, \dots, x_n , estimațiile de verosimilitate maximă pentru m și respectiv σ , anume

$$\hat{m} = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x} \quad \text{și} \quad \hat{\sigma} = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} = \sqrt{\mu_2}$$

(a se vedea Exemplul 5.3.7). Deoarece domeniul valorilor caracteristicii X este \mathbb{R} se consideră clasele date prin $-\infty = a_0 < a_1 < \dots < a_k = +\infty$.

Pe de altă parte probabilitățile corespunzătoare acestor clase se obțin cu formulele

$$p_i^{(0)} = F_0(a_i; \bar{x}, \sqrt{\mu_2}) - F_0(a_{i-1}; \bar{x}, \sqrt{\mu_2}), \quad i = \overline{1, k}.$$

Având în vedere că

$$F_0(x; m, \sigma) = \frac{1}{2} + \Phi\left(\frac{x - m}{\sigma}\right),$$

unde Φ este funcția lui Laplace, adică

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

rezultă că

$$p_i^{(0)} = \Phi\left(\frac{a_i - \bar{x}}{\sqrt{\mu_2}}\right) - \Phi\left(\frac{a_{i-1} - \bar{x}}{\sqrt{\mu_2}}\right), \quad i = \overline{1, k}.$$

Etapele aplicării testului

1. Se dau: α ; x_1, x_2, \dots, x_n ; a_1, a_2, \dots, a_{k-1} ;
2. Se calculează frecvențele absolute $n_i, i = \overline{1, k}$, estimațiile

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \bar{\mu}_2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

și probabilitățile

$$p_1^{(0)} = \frac{1}{2} + \Phi\left(\frac{a_1 - \bar{x}}{\sqrt{\bar{\mu}_2}}\right), \quad p_k^{(0)} = \frac{1}{2} - \Phi\left(\frac{a_{k-1} - \bar{x}}{\sqrt{\bar{\mu}_2}}\right),$$

$$p_i^{(0)} = \Phi\left(\frac{a_i - \bar{x}}{\sqrt{\bar{\mu}_2}}\right) - \Phi\left(\frac{a_{i-1} - \bar{x}}{\sqrt{\bar{\mu}_2}}\right), \quad i = \overline{2, k-1};$$

3. Se calculează $\chi_{k-3, 1-\alpha}^2$ astfel încât $F_{k-3}(\chi_{k-3, 1-\alpha}^2) = 1 - \alpha$;

$$4. \text{ Se calculează } x^2 = \sum_{i=1}^k \frac{(n_i - np_i^{(0)})^2}{np_i^{(0)}};$$

5. Concluzia: dacă $x^2 < \chi_{k-3, 1-\alpha}^2$ se acceptă ipoteza că X urmează legea normală $\mathcal{N}(\bar{x}, \sqrt{\bar{\mu}_2})$, în caz contrar ipoteza este respinsă.

Exemplul 6.9.13. Se consideră caracteristica X , ce reprezintă rezistența, în $K\Omega$, a unor tronsoane de ceramică acoperite cu carbon. Să se verifice normalitatea lui X , folosind o selecție de volum $n = 124$, pentru care s-au obținut datele de selecție

clasa	$(-\infty; 1.65]$	$(1.65; 1.70]$	$(1.70; 1.75]$	$(1.75; 1.80]$	$(1.80; 1.85]$
frecv.	11	14	17	17	18

$(1.85; 1.90]$	$(1.90; 1.95]$	$(1.95; 2.00]$	$(2.0; +\infty)$
16	13	10	8

utilizând testul de concordanță χ^2 , cu nivelul de semnificație $\alpha = 0.05$.

Prima dată se estimează parametrii de la legea normală $\mathcal{N}(m, \sigma)$, adică media teoretică $m = E(X)$ și abaterea standard teoretică $\sigma = \sqrt{\text{Var}(X)}$, folosind metoda de verosimilitate maximă. Se cunoaște că estimațiile de verosimilitate maximă pentru m și σ sunt respectiv

$$\hat{m} = \frac{1}{n} \sum_{i=1}^k x_i = \bar{x}, \quad \hat{\sigma} = \sqrt{\frac{1}{n} \sum_{i=1}^k (x_i - \bar{x})^2} = \sqrt{\mu_2},$$

(a se vedea Exemplul 5.3.7). Avem distribuția empirică de selecție pentru caracteristica X

$$X \begin{pmatrix} 1.625 & 1.675 & 1.725 & 1.775 & 1.825 & 1.875 & 1.925 & 1.975 & 2.125 \\ 11 & 14 & 17 & 17 & 18 & 16 & 13 & 10 & 8 \end{pmatrix}$$

de unde calculăm

$$\hat{m} = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{124} (11 \cdot 1.625 + 14 \cdot 1.675 + \dots + 8 \cdot 2.125) = 1.82;$$

$$\hat{\sigma} = \sqrt{\mu_2} = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} = 0.129.$$

Se consideră valoarea numerică

$$x^2 = \sum_{i=1}^k \frac{(f_i - n\hat{p}_i)^2}{n\hat{p}_i},$$

unde k este numărul claselor ($k = 9$, în cazul de față), f_i este frecvența clasei i , iar \hat{p}_i este dat prin

$$\hat{p}_i = \Phi\left(\frac{a_i - \bar{x}}{\hat{\sigma}}\right) - \Phi\left(\frac{a_{i-1} - \bar{x}}{\hat{\sigma}}\right),$$

subintervalul $[a_{i-1}, a_i)$ definind clasa i . După cum este cunoscut, x^2 este valoarea unei variabile aleatoare care urmează legea χ^2 cu $k - s - 1$ grade de libertate, s fiind numărul parametrilor estimați. În cazul de față avem că $s = 2$, deci avem $9 - 2 - 1 = 6$ grade de libertate.

Se determină intervalul $(0, \chi_{k-3, 1-\alpha}^2)$, pentru statistica χ^2 , folosind *Anexa III*. Anume, se obține că $\chi_{k-3, 1-\alpha}^2 = \chi_{6; 0.95}^2 = 12.59$, adică intervalul numeric pentru statistica χ^2 este $(0; 12.59)$.

Calculule pentru valoarea numerică x^2 se aranjează în următorul tabel:

a_i	f_i	$\frac{a_i - \bar{x}}{\hat{\sigma}}$	$\Phi\left(\frac{a_i - \bar{x}}{\hat{\sigma}}\right)$	\hat{p}_i	$n\hat{p}_i$	$\frac{(f_i - n\hat{p}_i)^2}{n\hat{p}_i}$
1.65	11	-1.40	-0.4188	0.0812	10.0647	0.0869
1.70	14	-0.97	-0.3238	0.0861	10.6711	1.0385
1.75	17	-0.53	-0.2030	0.1298	16.0891	0.0516
1.80	17	-0.10	-0.0402	0.1628	20.1848	0.5025
1.85	18	0.33	0.1297	0.1699	21.0714	0.4477
1.90	16	0.76	0.2773	0.1476	18.3036	0.2899
1.95	13	1.20	0.3840	0.1067	13.2299	0.0040
2.00	10	1.63	0.4482	0.0642	7.9568	0.5247
$+\infty$	8	$+\infty$	0.5000	0.0518	6.4286	0.3841
$n = 124$						$x^2 = 3.3$

Valorile funcției lui Laplace Φ se iau din *Anexa I* și se are în vedere că funcția Φ este impară, adică $\Phi(-x) = -\Phi(x)$. De asemenea, facem observația că

$$\begin{aligned}\hat{p}_1 &= \Phi\left(\frac{a_1 - \bar{x}}{\hat{\sigma}}\right) - \Phi\left(\frac{a_0 - \bar{x}}{\hat{\sigma}}\right) = \Phi(-1.4) - \Phi(-\infty) \\ &= -0.4188 + 0.5 = 0.0812.\end{aligned}$$

Deoarece $x^2 = 3.3 \in (0; 12.59)$, rezultă că se acceptă ipoteza normalității caracteristicii X .

Programul 6.9.14. Toate calculele din tabelul precedent pot fi efectuate prin programul Matlab:

```
x=1.625:0.05:2.025; f=[11,14,17,17,18,16,13,10,8];
ma=sum(f.*x)/124; s=sqrt(sum(f.*(x-ma).^2)/124);
a=[-inf,1.65:0.05:2,inf];
F=normcdf(a,ma,s); p=diff(F);
x2=sum((f-124*p).^2/(124*p));
cuant=chi2inv(0.95,6);
fprintf(' x^2= %6.1f\n cuant= %5.2f',x2,cuant)
```

Rezultatele obținute în urma executării programului sunt:

```
x^2=      3.3
cuant= 12.59
```

6.10 Testul χ^2 pentru compararea mai multor caracteristici

Se consideră colectivitățile $C_i, i = \overline{1, r}$, independente, cercetate din punct de vedere al aceleiași caracteristici (calitative sau cantitative). Fie această caracteristică X_i pentru colectivitatea C_i . Relativ la caracteristicile $X_i, i = \overline{1, r}$, se face ipoteza nulă H_0 , că acestea urmează aceeași lege de probabilitate. Se notează cu s numărul claselor caracteristicii cercetate și E_j evenimentul ca un individ luat la întâmplare din una din colectivitățile cercetate să aparțină clasei cu numărul de ordine j .

Dacă se notează prin $p_{ij} = P(E_j | C_i)$, adică probabilitatea ca un individ luat la întâmplare din colectivitatea C_i să aparțină clasei cu numărul de ordine j , atunci ipoteza nulă se poate rescrie $H_0 : p_{ij} = p_j, i = \overline{1, r}, j = \overline{1, s}$.

Probabilitățile astfel definite satisfac relațiile

$$\sum_{j=1}^s p_{ij} = 1, \quad i = \overline{1, r},$$

iar când ipoteza H_0 este satisfăcută, avem de asemenea $\sum_{j=1}^s p_j = 1$.

Pentru a verifica această ipoteză se consideră câte o selecție repetată de volum $n_i, i = \overline{1, r}$. Datele de selecție sunt trecute în tabelul următor.

Selecția	Vol. selecției	E_1	E_2	\dots	E_s	Total
S_1	n_1	n_{11}	n_{12}	\dots	n_{1s}	$n_{1.} = n_1$
S_2	n_2	n_{21}	n_{22}	\dots	n_{2s}	$n_{2.} = n_2$
\vdots	\vdots	\vdots	\vdots		\vdots	\vdots
S_r	n_r	n_{r1}	n_{r2}	\dots	n_{rs}	$n_{r.} = n_r$
Total	n	$n_{.1}$	$n_{.2}$	\dots	$n_{.s}$	$n_{..} = n$

Elementul n_{ij} al tabelului reprezintă frecvența absolută a apariției clasei i în selecția S_j , iar

$$n_{i.} = \sum_{j=1}^s n_{ij}, \quad i = \overline{1, r}, \quad n_{.j} = \sum_{i=1}^r n_{ij}, \quad j = \overline{1, s}, \quad n_{..} = \sum_{i=1}^r n_{i.} = \sum_{j=1}^s n_{.j}.$$

Proprietatea 6.10.1. *Estimațiile de verosimilitate maximă pentru parametrii necunoscuți $p_j, j = \overline{1, s}$, sunt date prin*

$$\hat{p}_j = \frac{n_{.j}}{n}, \quad j = \overline{1, s}.$$

Demonstrație. Funcția de verosimilitate are expresia

$$g = \prod_{i=1}^r \prod_{j=1}^s p_{ij}^{n_{ij}} = \prod_{i=1}^r \prod_{j=1}^s p_j^{n_{ij}} = \prod_{j=1}^s \left(\prod_{i=1}^r p_j^{n_{ij}} \right) = \prod_{j=1}^s p_j^{n_{\cdot j}},$$

când ipoteza nulă H_0 este adevărată. Având în vedere relația dintre probabilitățile p_j , $j = \overline{1, s}$, se obține că

$$g = \left(1 - \sum_{j=1}^{s-1} p_j \right)^{n_{\cdot s}} \prod_{j=1}^{s-1} p_j^{n_{\cdot j}}$$

sau prin logaritmare, se ajunge la

$$\ln g = n_{\cdot s} \ln \left(1 - \sum_{j=1}^{s-1} p_j \right) + \sum_{j=1}^{s-1} n_{\cdot j} \ln p_j.$$

Sistemul ecuațiilor de verosimilitate maximă va fi

$$\frac{\partial \ln g}{\partial p_k} = -\frac{n_{\cdot s}}{1 - \sum_{j=1}^{s-1} p_j} + \frac{n_{\cdot k}}{p_k} = 0, \quad k = \overline{1, s-1},$$

de unde rezultă că

$$\frac{n_{\cdot k}}{p_k} = \frac{n_{\cdot s}}{p_s}, \quad k = \overline{1, s},$$

sau, având în vedere proprietăți ale șirului de rapoarte egale, se obține

$$\frac{n_{\cdot k}}{p_k} = \frac{n}{1}, \quad k = \overline{1, s}.$$

S-a ajuns astfel la estimațiile

$$\hat{p}_k = \frac{n_{\cdot k}}{n}, \quad k = \overline{1, s}.$$

□

Observația 6.10.2. Dacă se are în vedere că n_{ij} sunt valorile unor variabile aleatoare binomiale N_{ij} , rezultă că statistica

$$\chi^2 = \sum_{i=1}^r \left[\sum_{j=1}^s \frac{(N_{ij} - N_{i\cdot} \hat{p}_j)^2}{N_{i\cdot} \hat{p}_j} \right] = \sum_{i=1}^r \sum_{j=1}^s \frac{\left(N_{ij} - \frac{N_{i\cdot} N_{\cdot j}}{n} \right)^2}{\frac{N_{i\cdot} N_{\cdot j}}{n}},$$

cu

$$\sum_{j=1}^s N_{ij} = N_{i\cdot}, \quad i = \overline{1, r}, \quad \sum_{i=1}^r N_{ij} = N_{\cdot j}, \quad j = \overline{1, s}, \quad \sum_{j=1}^s N_{\cdot j} = \sum_{i=1}^r N_{i\cdot} = n_{\cdot\cdot} = n.$$

urmează o lege χ^2 . Pentru a stabili numărul m al gradelor de libertate se ține seama de faptul că numărul variabilelor aleatoare N_{ij} este rs , numărul legăturilor între p_{ij} este r , iar numărul parametrilor estimați este $s - 1$. Prin urmare, numărul gradelor de libertate este $m = rs - r - (s - 1) = (r - 1)(s - 1)$.

Etapele aplicării testului

1. Se dau: $\alpha; n_{ij}, i = \overline{1, r}, j = \overline{1, s}$;
2. Se calculează $n_{i.} = \sum_{j=1}^s n_{ij}, i = \overline{1, r}, n_{.j} = \sum_{i=1}^r n_{ij}, j = \overline{1, s}$;
3. Se calculează $\chi_{m, 1-\alpha}^2$, astfel încât $F_m(\chi_{m, 1-\alpha}^2) = 1 - \alpha$;
4. Se calculează $x^2 = \sum_{i=1}^r \sum_{j=1}^s \frac{(n_{ij} - \frac{n_{i.}n_{.j}}{n})^2}{\frac{n_{i.}n_{.j}}{n}}$;
5. Concluzia: dacă $x^2 < \chi_{m, 1-\alpha}^2$ se acceptă ipoteza nulă H_0 , în caz contrar se respinge.

Observația 6.10.3. Când probabilitățile $p_j, j = \overline{1, r}$, sunt cunoscute (caz mai rar întâlnit), deci nu se cere a fi estimate, atunci statistica

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^s \frac{(N_{ij} - N_{i.}p_j)^2}{N_{i.}p_j}$$

urmează legea χ^2 cu $m = rs - r = r(s - 1)$ grade de libertate. Aplicarea testului, în acest caz, se adaptează în mod corespunzător.

Exemplul 6.10.4. S-au considerat trei loturi de bolnavi de emfizem pulmonar, în funcție de numărul țigărilor fumate zilnic: mai puțin de un pachet, unul sau două pachete și respectiv mai mult de două pachete. Pentru emfizemul pulmonar sunt considerate 4 stadii notate de la I la IV. Rezultatele cercetărilor sunt trecute în următorul tabel sistematizat

Categoria \ Stadiul	I	II	III	IV	
< 1 pachet	41	28	25	6	100
1 – 2 pachete	24	116	46	14	200
> 2 pachete	4	50	34	12	100
	69	194	105	32	400

Vrem să verificăm, cu nivelul de semnificație $\alpha = 0.01$, ipoteza că numărul de țigări fumate zilnic nu influențează repartizarea bolnavilor în cele patru stadii ale bolii. Verificăm ipoteza făcută cu ajutorul testului χ^2 privind compararea aceluiași atribut, dar pentru populații diferite. Numărul gradelor de libertate este dat prin formula $(r-1)(s-1) = (3-1)(4-1) = 6$, pentru care se obține cuantila $\chi_{k,1-\alpha}^2 = \chi_{6,0.99}^2 = 16.8$.

Se calculează apoi valoarea

$$\begin{aligned} x^2 &= \sum_{i=1}^r \sum_{j=1}^s \frac{\left(n_{ij} - \frac{n_{i.}n_{.j}}{n}\right)^2}{\frac{n_{i.}n_{.j}}{n}} \\ &= \frac{\left(41 - \frac{100 \cdot 69}{400}\right)^2}{\frac{100 \cdot 69}{400}} + \frac{\left(28 - \frac{100 \cdot 194}{400}\right)^2}{\frac{100 \cdot 194}{400}} + \dots + \frac{\left(12 - \frac{100 \cdot 32}{400}\right)^2}{\frac{100 \cdot 32}{400}} = 64.41. \end{aligned}$$

Deoarece $x^2 = 64.41 > 16.8 = \chi_{k,1-\alpha}^2$, se respinge ipoteza că repartizarea pe cele patru stadii ale bolii nu depinde de numărul țigărilor fumate zilnic.

Programul 6.10.5. Programul Matlab

```
f=[41,28,25,6;24,116,46,14;4,50,34,12];
fip=sum(f'); fjp=sum(f); x2=0;
for i=1:3
    for j=1:4
        t=fip(i)*fjp(j)/400;
        x2=x2+(f(i,j)-t)^2/t;
    end
end
cuant=chi2inv(0.99,6);
fprintf(' x^2= %6.1f\n cuant= %5.2f',x2,cuant)
```

În urma executării, conduce la rezultatele din exemplul precedent:

```
x^2=      64.4
cuant= 16.81
```

6.11 Testul χ^2 pentru tabele de contingență

Se consideră colectivitatea \mathcal{C} cercetată din punct de vedere a două caracteristici X și Y (calitative sau cantitative). Se pune problema independenței celor două caracteristici. Fie r numărul claselor pentru caracteristica X și respectiv s pentru caracteristica Y . Notăm cu A_i evenimentul ca un individ luat la întâmplare din colectivitatea \mathcal{C} să aparțină clasei cu numărul de ordine i în raport cu X și notăm cu B_j evenimentul ca un individ luat la întâmplare din \mathcal{C} să aparțină clasei cu numărul de ordine j în raport

cu Y . Fie $p_{ij} = P(A_i \cap B_j)$, $i = \overline{1, r}$, $j = \overline{1, s}$. Dacă se notează

$$p_{i.} = \sum_{j=1}^s p_{ij}, \quad i = \overline{1, r} \quad \text{și} \quad p_{.j} = \sum_{i=1}^r p_{ij}, \quad j = \overline{1, s},$$

atunci ipoteza nulă relativă la independența celor două caracteristici se scrie

$$H_0 : p_{ij} = p_{i.} p_{.j}, \quad i = \overline{1, r}, \quad j = \overline{1, s},$$

cu ipoteza alternativă H_1 , că H_0 este falsă.

Pentru verificarea acestei ipoteze se consideră datele de selecție (x_i, y_i) , $i = \overline{1, n}$. Datele de selecție sunt sistematizate în tabelul de contingență

$X \setminus Y$	B_1	B_2	\dots	B_s	
A_1	n_{11}	n_{12}	\dots	n_{1s}	$n_{1.}$
A_2	n_{21}	n_{22}	\dots	n_{2s}	$n_{2.}$
\vdots	\vdots	\vdots		\vdots	\vdots
A_r	n_{r1}	n_{r2}	\dots	n_{rs}	$n_{r.}$
	$n_{.1}$	$n_{.2}$	\dots	$n_{.s}$	$n_{..} = n$

unde n_{ij} reprezintă frecvența absolută a clasei (i, j) , iar

$$n_{i.} = \sum_{j=1}^s n_{ij}, \quad i = \overline{1, r}, \quad n_{.j} = \sum_{i=1}^r n_{ij}, \quad j = \overline{1, s}.$$

Proprietatea 6.11.1. *Estimațiile de verosimilitate maximă pentru parametrii necunoscuți $p_{i.}$, $i = \overline{1, r}$, și $p_{.j}$, $j = \overline{1, s}$, sunt date prin formulele:*

$$\hat{p}_{i.} = \frac{n_{i.}}{n}, \quad i = \overline{1, r}, \quad \hat{p}_{.j} = \frac{n_{.j}}{n}, \quad j = \overline{1, s}.$$

Demonstrație. Funcția de verosimilitate are expresia

$$\begin{aligned} g &= \prod_{i=1}^r \prod_{j=1}^s (p_{ij})^{n_{ij}} = \prod_{i=1}^r \prod_{j=1}^s (p_{i.} p_{.j})^{n_{ij}} = \left(\prod_{i=1}^r p_{i.}^{n_{i.}} \right) \left(\prod_{j=1}^s p_{.j}^{n_{.j}} \right) \\ &= \left(1 - \sum_{i=1}^{r-1} p_{i.} \right)^{n_{r.}} \left(\prod_{i=1}^{r-1} p_{i.}^{n_{i.}} \right) \left(\prod_{j=1}^s p_{.j}^{n_{.j}} \right), \end{aligned}$$

când ipoteza nulă H_0 este adevărată. Prin logaritmare se obține

$$\ln g = n_{r.} \ln \left(1 - \sum_{i=1}^{r-1} p_{i.} \right) + \sum_{i=1}^{r-1} n_{i.} \ln p_{i.} + \ln \left(\prod_{j=1}^s p_{.j}^{n_{.j}} \right),$$

de unde se ajunge la sistemul ecuațiilor de verosimilitate maximă

$$\frac{\partial \ln g}{\partial p_k} = -\frac{n_{r.}}{1 - \sum_{i=1}^{r-1} p_{i.}} + \frac{n_{k.}}{p_{k.}} = 0, \quad k = \overline{1, r-1},$$

sau

$$\frac{n_{r.}}{p_{r.}} = \frac{n_{k.}}{p_{k.}}, \quad k = \overline{1, r}.$$

Având în vedere proprietățile unui șir de rapoarte egale se obține că

$$\frac{n_{k.}}{p_{k.}} = \frac{n}{1}, \quad k = \overline{1, r}, \quad \text{deci} \quad \hat{p}_{k.} = \frac{n_{k.}}{n}, \quad k = \overline{1, r}.$$

În mod analog se obține și faptul că

$$\hat{p}_{.k} = \frac{n_{.k}}{n}, \quad k = \overline{1, s}.$$

□

Observația 6.11.2. Valoarea numerică

$$x^2 = \sum_{i=1}^r \sum_{j=1}^s \frac{(n_{ij} - n_{i.}\hat{p}_{.j})^2}{n\hat{p}_{.j}} = \sum_{i=1}^r \sum_{j=1}^s \frac{(n_{ij} - \frac{n_{i.}n_{.j}}{n})^2}{\frac{n_{i.}n_{.j}}{n}},$$

este cea a unei variabile aleatoare ce urmează legea χ^2 cu $m = (r-1)(s-1)$ grade de libertate. Numărul m al gradelor de libertate s-a determinat având în vedere că au fost estimați $(r-1) + (s-1)$ parametri, iar între $p_{i.}$ și $p_{.j}$ există o singură legătură

$$\sum_{i=1}^r p_{i.} = \sum_{j=1}^s p_{.j} = 1, \quad \text{adică} \quad \sum_{i=1}^r \sum_{j=1}^s p_{ij} = 1.$$

Așadar $m = rs - (r-1) - (s-1) - 1 = (r-1)(s-1)$.

Etapele aplicării testului

1. Se dau: $\alpha; n_{ij}, i = \overline{1, r}, j = \overline{1, s};$
2. Se calculează: $n_{i.} = \sum_{j=1}^s n_{ij}, i = \overline{1, r}, n_{.j} = \sum_{i=1}^r n_{ij}, j = \overline{1, s};$
3. Se calculează $\chi_{m, 1-\alpha}^2$, astfel încât $F_m(\chi_{m, 1-\alpha}^2) = 1 - \alpha;$

4. Se calculează $x^2 = \sum_{i=1}^r \sum_{j=1}^s \frac{(n_{ij} - \frac{n_{i.}n_{.j}}{n})^2}{\frac{n_{i.}n_{.j}}{n}};$

5. Concluzia: dacă $x^2 < \chi_{m,1-\alpha}^2$ se acceptă ipoteza nulă H_0 , în caz contrar se respinge.

Observația 6.11.3. Când probabilitățile $p_{i.}$ și $p_{.j}$ sunt cunoscute (caz mai rar întâlnit), avem că

$$x^2 = \sum_{i=1}^r \sum_{j=1}^s \frac{(n_{ij} - np_{ij})^2}{np_{ij}} = \sum_{i=1}^r \frac{(n_{i.} - np_{i.})^2}{np_{i.}} + \sum_{j=1}^s \frac{(n_{.j} - np_{.j})^2}{np_{.j}}$$

este valoarea unei variabile aleatoare ce urmează legea χ^2 cu $m = (r-1)(s-1)$ grade de libertate. Aceasta rezultă din faptul că prima sumă corespunde unei legi χ^2 cu $rs - 1$ grade de libertate, a doua sumă corespunde unei legi χ^2 cu $r - 1$ grade de libertate, iar a treia unei legi χ^2 cu $s - 1$ grade de libertate. Prin urmare

$$m + (r - 1) + (s - 1) = rs - 1, \quad \text{de unde} \quad m = (r - 1)(s - 1).$$

Aplicarea testului, în acest caz, se face prin adaptarea corespunzătoare a cazului precedent.

Observația 6.11.4. Testul χ^2 pentru tabele de contingență și respectiv pentru compararea mai multor caracteristici, pare să fie același. Dar, remarcăm faptul că problemele de la care se pornește sunt complet diferite.

Exemplul 6.11.5. Pentru cercetarea dependenței dintre culoarea A a ochilor și culoarea B a părului, s-a considerat un eșantion format din 6800 indivizi. Rezultatele sunt trecute în tabelul sistematizat următor

A \ B	Blonzi	Șateni	Bruneți	Roșcați	
Albaștri	1768	807	189	47	2811
Gri sau verzi	946	1387	746	53	3132
Căprui	115	438	288	16	857
	2829	2632	1223	116	6800

Folosind nivelul de semnificație $\alpha = 0.01$, vom verifica ipoteza că cele două caracteristici (atribute) sunt independente.

Se folosește testul χ^2 pentru tabele de contingență cu numărul gradelor de libertate $k = (r - 1)(s - 1) = (3 - 1)(4 - 1) = 6$.

Pe de o parte avem că

$$x^2 = \sum_{i=1}^r \sum_{j=1}^s \frac{\left(n_{ij} - \frac{n_{i.} \cdot n_{.j}}{n}\right)^2}{\frac{n_{i.} \cdot n_{.j}}{n}}$$

$$= \frac{\left(1768 - \frac{2811 \cdot 2829}{6800}\right)^2}{\frac{2811 \cdot 2829}{6800}} + \frac{\left(807 - \frac{2811 \cdot 2632}{6800}\right)^2}{\frac{2811 \cdot 2632}{6800}} + \dots + \frac{\left(16 - \frac{857 \cdot 116}{6800}\right)^2}{\frac{857 \cdot 116}{6800}} = 1075.$$

Pe de altă parte, avem cuantila $\chi_{6,0.99}^2 = 16.8 < 1075 = x^2$, deci se respinge ipoteza că cele două atribute sunt independente.

Programul 6.11.6. Calculele din exemplul precedent pot fi efectuate cu programul Matlab

```
f=[1768,807,189,47;946,1387,746,53;115,438,288,16];
fip=sum(f'); fipj=sum(f); x2=0;
for i=1:3
    for j=1:4
        t=fip(i)*fipj(j)/6800;
        x2=x2+(f(i,j)-t)^2/t;
    end
end
cuant=chi2inv(0.99,6);
fprintf(' x^2= %6.1f\n cuant= %5.2f',x2,cuant)
```

În urma executării acestui program, se obțin rezultatele:

```
x^2= 1073.5
cuant= 16.81
```

6.12 Testul de concordanță al lui Kolmogorov

Se consideră caracteristica X de tip continuu cu funcția teoretică de repartiție F necunoscută. Relativ la funcția F se face ipoteza nulă $H_0 : F = F_0$, cu una din

- (1) $H_1 : F \neq F_0$ (testul lui Kolmogorov bilateral),
- (2) $H_1 : F > F_0$ (testul lui Kolmogorov unilateral dreapta),
- (3) $H_1 : F < F_0$ (testul lui Kolmogorov unilateral stânga).

Pentru verificarea ipotezei nule H_0 , cu una din alternativele precizate, se consideră o selecție repetată de volum n , cu datele de selecție x_1, x_2, \dots, x_n , respectiv variabilele de selecție X_1, X_2, \dots, X_n corespunzătoare. Se consideră statisticile

$$D_n = \sup_{x \in \mathbb{R}} \left\{ \left| \bar{F}_n(x) - F_0(x) \right| \right\},$$

$$D_n^+ = \sup_{x \in \mathbb{R}} \{ \bar{F}_n(x) - F_0(x) \},$$

$$D_n^- = \sup_{x \in \mathbb{R}} \{ F_0(x) - \bar{F}_n(x) \},$$

unde $\bar{F}_n(x)$ este funcția de repartiție de selecție.

Conform teoremei lui Kolmogorov (Teorema 4.2.38) statistica

$$D_n = \sup_{x \in \mathbb{R}} | \bar{F}_n(x) - F_0(x) |,$$

urmează o lege de probabilitate dată prin legea lui Kolmogorov, când $n \rightarrow \infty$. Anume, avem că

$$\lim_{n \rightarrow \infty} P(\sqrt{n} D_n \leq x \mid H_0) = K(x) = \sum_{k=-\infty}^{k=+\infty} (-1)^k e^{-2k^2 x^2}, \quad x > 0,$$

funcția $K(x)$ a lui Kolmogorov fiind tabelată pentru anumite valori în *Anexa V*.

De asemenea,

$$\lim_{n \rightarrow \infty} P(\sqrt{n} D_n^+ \leq x) = \lim_{n \rightarrow \infty} P(\sqrt{n} D_n^- \leq x) = K^\pm(x) = 1 - e^{-2x^2}, \quad x > 0,$$

ce poartă numele de lege χ cu două grade de libertate.

Pentru un nivel de semnificație $\alpha \in (0, 1)$ fixat, se pot determina cuantilele $k_{1-\alpha}$ și $k_{1-\alpha}^\pm$ astfel încât

$$P(\sqrt{n} D_n \leq k_{1-\alpha}) = 1 - \alpha, \quad \text{adică} \quad K(k_{1-\alpha}) = 1 - \alpha,$$

pentru testul bilateral al lui Kolmogorov,

$$P(\sqrt{n} D_n^+ \leq k_{1-\alpha}^\pm) = 1 - \alpha \quad \text{și} \quad P(\sqrt{n} D_n^- \leq k_{1-\alpha}^\pm) = 1 - \alpha,$$

adică $K^\pm(k_{1-\alpha}^\pm) = 1 - \alpha$, pentru testele unilaterale dreapta și stânga ale lui Kolmogorov.

Corespunzător celor trei alternative, ipoteza H_0 va fi admisă când valorile calculate d_n, d_n^+, d_n^- , pe baza datelor de selecție x_1, x_2, \dots, x_n ale statisticilor D_n, D_n^+, D_n^- , satisfac respectiv condițiile

- (1) $\sqrt{n} d_n < k_{1-\alpha}$, când $H_1 : F = F_0$,
- (2) $\sqrt{n} d_n^+ < k_{1-\alpha}^\pm$, când $H_1 : F > F_0$,
- (3) $\sqrt{n} d_n^- < k_{1-\alpha}^\pm$, când $H_1 : F < F_0$,

iar în caz contrar va fi respinsă.

Prezentăm, în cele ce urmează, împreună, cele trei teste (bilateral, unilateral dreapta, unilateral stânga) ale lui Kolmogorov. Remarcăm că cele trei teste sunt teste distincte, alegerea unuia dintre ele se face apriori.

Etapale aplicării testului lui Kolmogorov

1. Se dau: $\alpha; x_1, x_2, \dots, x_n; F = F_0$;

2. Se determină

$$(1) \quad k_{1-\alpha} \text{ astfel încât } K(k_{1-\alpha}) = 1 - \alpha, \text{ când } H_1 : F \neq F_0,$$

$$(2) \quad k_{1-\alpha}^{\pm} \text{ astfel încât } K^{\pm}(k_{1-\alpha}^{\pm}) = 1 - \alpha, \text{ când } H_1 : F > F_0,$$

$$(3) \quad k_{1-\alpha}^{\pm} \text{ astfel încât } K^{\pm}(k_{1-\alpha}^{\pm}) = 1 - \alpha, \text{ când } H_1 : F < F_0;$$

3. Se calculează

$$(1) \quad k = \sqrt{n} d_n, \text{ când } H_1 : F \neq F_0,$$

$$(2) \quad k = \sqrt{n} d_n^+, \text{ când } H_1 : F > F_0,$$

$$(3) \quad k = \sqrt{n} d_n^-, \text{ când } H_1 : F < F_0;$$

4. Concluzia: dacă

$$(1) \quad k < k_{1-\alpha} \text{ ipoteza } H_0 \text{ este admisă, când } H_1 : F \neq F_0,$$

$$(2) \quad k < k_{1-\alpha}^{\pm} \text{ ipoteza } H_0 \text{ este admisă, când } H_1 : F > F_0,$$

$$(3) \quad k < k_{1-\alpha}^{\pm} \text{ ipoteza } H_0 \text{ este admisă, când } H_1 : F < F_0,$$

Observația 6.12.1. Având în vedere că funcțiile (de repartiție) K și K^{\pm} sunt monoton crescătoare și că pentru determinarea cuantilelor $k_{1-\alpha}$ și $k_{1-\alpha}^{\pm}$ este necesară inversarea acestor funcții de repartiție, se poate evita această operație. Pentru aceasta se determină valorile funcțiilor K și K^{\pm} pe valorile calculate ale statisticilor D_n , D_n^+ , D_n^- , adică se calculează respectiv $1 - c = K(\sqrt{n}d_n)$, $1 - c = K^{\pm}(\sqrt{n}d_n^+)$, $1 - c = K^{\pm}(\sqrt{n}d_n^-)$. Astfel, se va respinge ipoteza nulă H_0 , dacă $c \leq \alpha$, c numindu-se *valoare critică*.

Observația 6.12.2. Pentru calculul valorilor statisticilor D_n , D_n^+ , D_n^- putem considera că datele de selecție sunt ordonate crescător, adică $x_1 < x_2 < \dots < x_n$, altfel se poate face o astfel de ordonare. De asemenea, remarcăm faptul că are loc $D_n = \max\{D_n^+, D_n^-\}$.

Folosind aceste precizări, precum și faptul că atât funcția de repartiție de selecție \bar{F}_n , cât și funcția de repartiție F_0 sunt funcții nedescrescătoare, obținem următoarele

formule de calcul

$$\begin{aligned}
 d_n^+ &= \sup_{x \in \mathbb{R}} \{ \bar{F}_n(x) - F_0(x) \} = \max_{k=1, n} \{ \bar{F}_n(x_k) - F_0(x_k) \} \\
 &= \max_{k=1, n} \left\{ \frac{k}{n} - F_0(x_k) \right\}, \\
 d_n^- &= \sup_{x \in \mathbb{R}} \{ F_0(x) - \bar{F}_n(x) \} = \max_{k=1, n} \{ F_0(x_k) - \bar{F}_n(x_k - 0) \} \\
 &= \max_{k=1, n} \left\{ F_0(x_k) - \frac{k-1}{n} \right\}, \\
 d_n &= \max \{ d_n^+, d_n^- \}.
 \end{aligned}$$

Observația 6.12.3. Când datele sunt grupate, având clasele date prin punctele

$$a = a_0 < a_1 < \dots < a_m = b,$$

atunci se folosesc formulele

$$\begin{aligned}
 d_n &= \max_{k=0, m} | \bar{F}_n(a_k) - F_0(a_k) |, \\
 d_n^+ &= \max_{k=0, m} \{ \bar{F}_n(a_k) - F_0(a_k) \}, \\
 d_n^- &= \max_{k=0, m} \{ F_0(a_k) - \bar{F}_n(a_k) \}.
 \end{aligned}$$

Exemplul 6.12.4. Rezultatele măsurătorilor asupra diametrului X , pentru 1000 de piese de același tip (în mm), sunt cele ce urmează

Diam.	97.75–98.25	98.25–98.75	98.75–99.25	99.25–99.75	99.75–100.25
Frecv.	21	47	87	158	181

100.25–100.75	100.75–101.25	101.25–101.75	101.75–102.25	102.25–102.75
201	142	97	41	25

Folosind nivelul de semnificație $\alpha = 0.05$, știind $m = E(X) = 100.25$ mm și abaterea standard $\sigma = \sqrt{\text{Var}(X)} = 1$ mm, se cere verificarea normalității caracteristicii X , cu ajutorul testului lui Kolmogorov bilateral, iar apoi cu ajutorul testului χ^2 neparametric.

Ipoteza care se face este că funcția de repartiție a lui X este

$$F_0(x) = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{(t-m)^2}{2\sigma^2}} dt, \quad x \in \mathbb{R},$$

unde $m = 100.25$, $\sigma = 1$. Pentru verificarea acestei ipoteze, folosind criteriul lui Kolmogorov, se calculează

$$d_n = \sup_{x \in \mathbb{R}} |\bar{F}_n(x) - F_0(x)| = \sup_{i=\overline{0,10}} |\bar{F}_n(a_i) - F_0(a_i)|,$$

unde $a_i = 97.75 + 0.5i$, $i = \overline{0,10}$, iar $\bar{F}_n(x) = \bar{F}_{1000}(x)$ este funcția de repartiție de selecție.

Pentru calculul valorilor funcției de repartiție F_0 se ține seama de faptul că

$$F_0(x) = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{(t-m)^2}{2\sigma^2}} dt = \frac{1}{2} + \Phi\left(\frac{x-m}{\sigma}\right) = \frac{1}{2} + \Phi(x - 100.25),$$

cu funcția lui Laplace $\Phi(x)$ tabelată în *Anexa I*.

Pe de altă parte, valorile funcției de repartiție de selecție se calculează cu formula

$$\bar{F}_{1000}(a_i) = \frac{1}{1000} \sum_{j=1}^i n_j, \quad \text{pentru } i = \overline{0,10}.$$

Calcululele sunt aranjate în următorul tabel.

a_i	n_i	$a_i - m$	$F_0(a_i)$	$\bar{F}_{1000}(a_i)$	$ \bar{F}_{1000}(a_i) - F_0(a_i) $
97.75	—	-2.5	0.0062	0.000	0.0062
98.25	21	-2.0	0.0228	0.021	0.0018
98.75	47	-1.5	0.0668	0.068	0.0012
99.25	87	-1.0	0.1587	0.155	0.0037
99.75	158	-0.5	0.3085	0.313	0.0045
100.25	181	0.0	0.5000	0.494	0.0060
100.75	201	0.5	0.6915	0.695	0.0035
101.25	142	1.0	0.8413	0.837	0.0043
101.75	97	1.5	0.9332	0.934	0.0008
102.25	41	2.0	0.9772	0.975	0.0022
102.75	25	2.5	0.9938	1.000	0.0062

Din calculele făcute avem că $d_{1000} = 0.0062$, de unde rezultă că

$$\sqrt{n} d_n = \sqrt{1000} d_{1000} = 10\sqrt{10} \cdot 0.0062 = 0.196.$$

Din *Anexa V*, se determină cuantila $k_{1-\alpha} = k_{0.95} = 1.36$ și admitem ipoteza de normalitate pentru caracteristica X , deoarece $\sqrt{n} d_n = 0.196 < 1.36 = k_{1-\alpha}$.

Pentru a verifica normalitatea caracteristicii X , cu testul χ^2 neparametric, calculăm prima dată probabilitățile

$$p_i^{(0)} = P(a_{i-1} < X \leq a_i | F = F_0) = \Phi(a_i - m) - \Phi(a_{i-1} - m), \quad i = \overline{1, 10},$$

unde $a_0 = -\infty$ și $a_{10} = +\infty$. Astfel avem

X	$(-\infty, a_1]$	$(a_1, a_2]$	$(a_2, a_3]$	$(a_3, a_4]$	$(a_4, a_5]$	$(a_5, a_6]$	$(a_6, a_7]$
$p_i^{(0)}$	0.0228	0.0440	0.0919	0.1498	0.1915	0.1915	0.1498

	$(a_7, a_8]$	$(a_8, a_9]$	$(a_9, +\infty)$
	0.0919	0.0440	0.0228

de unde

$$x^2 = \sum_{i=1}^{10} \frac{(n_i - np_i^{(0)})^2}{np_i^{(0)}} = \frac{(21 - 22.8)^2}{22.8} + \frac{(47 - 44)^2}{44} + \dots + \frac{(25 - 22.8)^2}{22.8} = 3.21.$$

Pe de altă parte, avem cuantila $\chi_{k-1, 1-\alpha}^2 = \chi_{9, 0.95}^2 = 16.9$, conform *Anexei III*. Deoarece $x^2 = 3.21 < 16.9 = \chi_{k-1, 1-\alpha}^2$, rezultă că ipoteza de normalitate este acceptată.

Programul 6.12.5. Vom scrie un program Matlab, pentru aplicarea testului lui Kolmogorov și a testului χ^2 , în cazul datelor din exemplul precedent.

```
% Testul lui Kolmogorov
a=97.75:0.5:102.75;
n=[0,21,47,87,158,181,201,142,97,41,25];
F0=normcdf(a,100.25,1); Fn=cumsum(n)/1000;
dn=max(abs(F0-Fn)); ks=sqrt(1000)*dn;
fprintf(' ks=    %6.4f\n',ks)
% Testul chi^2
F0(1)=0; F0(end)=1;n=n(2:end);
p0=diff(F0); x2=sum((n-1000*p0).^2./(1000*p0));
cuant=chi2inv(0.95,9);
fprintf(' x^2=    %6.4f\n',x2)
fprintf(' cuant=%6.2f',cuant)
```

Rezultatele obținute în urma executării programului sunt:

```
ks=    0.1964
x^2=   3.2117
cuant= 16.92
```

Ceea ce confirmă rezultatele mai sus obținute.

Exemplul 6.12.6. Se țin sub observație $n = 50$ motoare electrice până la defectarea ultimului dintre ele. Se consideră caracteristica X , ce reprezintă numărul miilor de ore de funcționare până la defectare. Rezultatele observațiilor sunt date în tabelul următor

Durata	0–30	30–60	60–90	90–120	120–150	150–180	180–300
Frecv.	13	10	5	5	5	5	7

Să se cerceteze exponențialitatea caracteristicii X , folosind testul χ^2 cu nivelul de semnificație $\alpha = 0.05$, iar apoi folosind testul lui Kolmogorov cu același nivel de semnificație.

Legea exponențială are funcția de repartiție

$$F(x; \mu) = 1 - e^{-\frac{x}{\mu}}, \quad x > 0, \quad \mu > 0 \text{ parametru necunoscut.}$$

Folosind metoda verosimilității maxime avem estimația pentru μ

$$\hat{\mu} = \bar{x} = \frac{1}{50} (13 \cdot 15 + 10 \cdot 45 + \dots + 7 \cdot 240) = 94.5.$$

Se calculează prima dată probabilitățile:

$$p_i^{(0)} = P(a_{i-1} < X \leq a_i \mid F = F_0) = F_0(a_i; \hat{\mu}) - F_0(a_{i-1}; \hat{\mu}), \quad i = \overline{1, 7},$$

unde

$$F_0(x; \hat{\mu}) = 1 - e^{-\frac{x}{\hat{\mu}}} = 1 - e^{-\frac{x}{94.5}}.$$

Rezultatele sunt trecute în tabelul următor

X	$(0, a_1]$	$(a_1, a_2]$	$(a_2, a_3]$	$(a_3, a_4]$	$(a_4, a_5]$	$(a_5, a_6]$	(a_6, ∞)
$p_i^{(0)}$	0.2720	0.1980	0.1442	0.1049	0.0764	0.0556	0.1489

și care conduce la

$$x^2 = \sum_{i=1}^k \frac{(n_i - np_i^{(0)})^2}{np_i^{(0)}} = \frac{(13 - 50 \cdot 0.2720)^2}{50 \cdot 0.2720} + \dots + \frac{(7 - 50 \cdot 0.1489)^2}{50 \cdot 0.1489} = 2.877.$$

Din *Anexa III* obținem cuantila $\chi_{k-2, 1-\alpha}^2 = \chi_{5, 0.95}^2 = 11.07$. Deoarece avem valoarea calculată $x^2 = 2.877 < 11.07 = \chi_{k-2, 1-\alpha}^2$ acceptăm ipoteza exponențialității.

Pentru aplicarea testului lui Kolmogorov, din *Anexa V* se determină cuantila $k_{1-\alpha} = k_{0.95} = 1.36$. Ipoteza exponențialității lui X va fi acceptată dacă

$$\sqrt{n} d_n = \sqrt{n} \sup_{i=\overline{0,7}} |\bar{F}_n(a_i) - F_0(a_i; \hat{\mu})| < k_{1-\alpha}.$$

Aici $a_7 = +\infty$, iar \bar{F}_n este funcția de repartiție de selecție. Calculele sunt efectuate în tabelul următor.

a_i	30	60	90	120	150	180	300
n_i	13	10	5	5	5	5	7
$\bar{F}_n(a_i)$	0.26	0.46	0.56	0.66	0.76	0.86	1
$F_0(a_i; \hat{\mu})$	0.2720	0.4700	0.6142	0.7191	0.7955	0.8511	0.9582
$ \bar{F}_n(a_i) - F_0(a_i; \hat{\mu}) $	0.0120	0.0100	0.0542	0.0591	0.0355	0.0089	0.0418

Așadar avem că $\sqrt{n}d_n = \sqrt{50} \cdot 0.0591 = 0.418 < 1.36 = k_{1-\alpha}$, deci exponențialitatea lui X este admisă.

Programul 6.12.7. Vom scrie un program Matlab, care efectuează calculele din exemplul precedent:

```
% Testul chi^2
clear
a=[0:30:180,300]; f=[13,10,5*ones(1,4),7];
mij=[15:30:165,240]; mu=sum(mij.*f)/50;
F0=expcdf(a,mu);F0(end)=1; p0=diff(F0);
x2=sum((f-50*p0).^2./(50*p0));
cuant=chi2inv(0.95,5);
fprintf(' x^2= %6.4f\n',x2)
fprintf(' cuant= %4.2f\n',cuant)
% Testul lui Kolmogorov
F0=F0(2:end);F0(end)=expcdf(a(end),mu);
Fn=cumsum(f)/50;
dn=max(abs(F0-Fn)); ks=sqrt(50)*dn;
fprintf(' ks= %6.4f',ks)
```

În urma executării programului se obține:

```
x^2= 2.8770
cuant= 11.07
ks= 0.4181
```

6.12.1 Funcția kstest

Statistics toolbox conține funcția `kstest`, care efectuează testul lui Kolmogorov. Apelul funcției se poate face cu una din instrucțiunile:

```
h=kstest(x)
h=kstest(x,cdf)
h=kstest(x,cdf,alpha)
h=kstest(x,cdf,alpha,tail)
[h,c,ks]=kstest(x,cdf,alpha,tail)
```

Comenzile de acest tip lansează execuția testului lui Kolmogorov bilateral (când `tail=0`, valoare implicită), unilateral dreapta (când `tail=1`), respectiv unilateral stânga (când `tail=-1`), prin considerarea datelor conținute de vectorul `x`, iar funcția

de repartiție ipotezată fiind precizată prin parametrul `cdf`. Implicit se consideră `cdf` corespunzând legii normale standard.

Dacă parametrul `cdf` este prezent, trebuie să fie o matrice cu două coloane, care conține în coloana a doua valorile funcției de repartiție ipotezate pe punctele precizate în prima coloană. Ar fi de dorit ca prima coloană a matricei `cdf` să coincidă cu vectorul `x`, altfel funcția va efectua un proces de interpolare pentru calculul valorilor funcției ipotezate pe punctele lui `x`, folosind punctele precizate prin matricea `cdf`. Din acest motiv, se impune ca toate componentele lui `x` să aibă valorile conținute în intervalul determinat de valoare minimă și valoarea maximă a componentelor primei coloane a lui `cdf`.

Parametrul `alpha` (implicit `alpha=0.05`) reprezintă nivelul de semnificație.

Rezultatele obținute au următoarele semnificații. Dacă `h=1` ipoteza nulă se respinge, ceea ce corespunde faptului că valoarea critică `c` satisface relația $c \leq \alpha$, iar dacă `h=0`, ipoteza nulă nu poate fi respinsă. Mai remarcăm faptul că `ks` va conține respectiv valoarea calculată a statisticilor D_n , D_n^+ și D_n^- , și verifică relațiile $c = 1 - K(\sqrt{n}d_n)$, pentru testul bilateral, respectiv $c = 1 - K^\pm(\sqrt{n}d_n^\pm)$ și, $c = 1 - K^\pm(\sqrt{n}d_n^\pm)$ pentru testele unilaterale.

Programul 6.12.8. Să aplicăm testele lui Kolmogorov (bilateral, unilateral dreapta și unilateral stânga) pentru datele `x=-2:1:4`, considerând legea ipotezată ca fiind legea normală standard și nivelul de semnificație `alpha=0.05`. Programul ce urmează, să reprezinte grafic și funcția de repartiție de selecție, împreună cu funcția de repartiție ipotezată.

```
x=-2:4;
fprintf(' h      c      ks\n')
fprintf('_____\n')
for i=-1:1
    [h,c,ks]=kstest(x,[],0.05,i);
    fprintf(' %d %5.4f %4.2f \n',h,c,ks)
end
cdfplot(x), hold on, title('')
t=-3:0.01:3; plot(t,normcdf(t,0,1),'k--')
```

În urma executării programului, se obțin rezultatele:

h	c	ks
0	0.0682	0.41
0	0.1363	0.41
0	0.7753	0.13

respectiv graficele din Figura 6.4 Se observă, în mod surprinzător, ipoteza nulă nu poate fi respinsă. Acest lucru se datorează faptului că testul lui Kolmogorov este bazat pe un rezultat asimptotic, adică se aplică pentru valori mari ale volumului selecției.

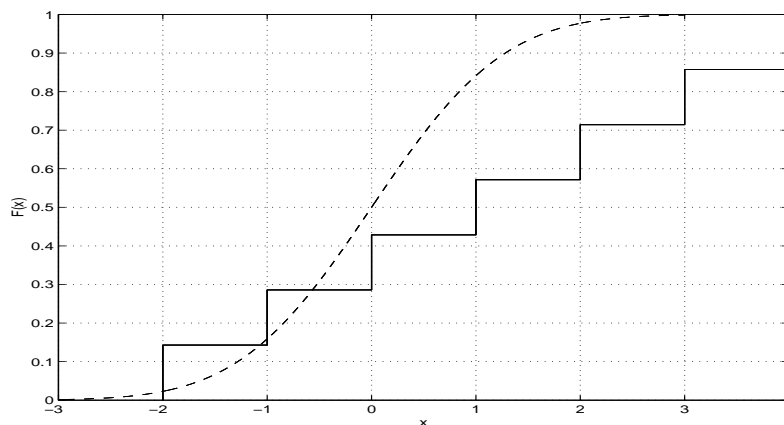


Figura 6.4: Testul lui Kolmogorov

6.12.2 Funcția `lillietest`

Funcția `lillietest` efectuează testul Lilliefors, care reprezintă o modificare a testului lui Kolmogorov, în care ipoteza nulă presupune că funcția de repartiție este cea de la legea normală, dar parametrii nu mai sunt precizați, ci sunt estimați folosind datele de selecție.

Apelul funcției se poate face cu una din comenzile

```
h=lillietest(x)
h=lillietest(x,alpha)
[h,c,ls]=lillietest(x,alpha)
```

Comenzile de acest tip lansează execuția testului Lilliefors, prin considerarea datelor conținute de vectorul `x`, iar funcția de repartiție ipotezată fiind cea de la legea normală $\mathcal{N}(\bar{x}, \bar{\sigma})$.

Parametrul `alpha` (implicit `alpha=0.05`) reprezintă nivelul de semnificație.

Rezultatele obținute au următoarele semnificații. Dacă `h=1` ipoteza nulă se respinge, ceea ce corespunde faptului că valoarea critică `c` satisface relația `c ≤ alpha`, iar dacă `h=0`, ipoteza nulă nu poate fi respinsă.

Mai remarcăm faptul că `ls` va conține respectiv valoarea calculată a statisticii D_n , iar `c` este *valoarea critică*, care se calculează cu ajutorul lui `ls`, folosindu-se tabelele construite de Lilliefors. Dacă în aceste tabele nu se află valoarea corespunzătoare, atunci `c=NaN`, dar `h` indică încă dacă ipoteza nulă este respinsă sau nu.

Programul 6.12.9. Să aplicăm testul lui Lilliefors pentru `n` numere aleatoare ce urmează legea normală $\mathcal{N}(\mu, \sigma)$, conținute de vectorul `x`, cu nivelul de semnificație

$\alpha=0.05$. Programul ce urmează să reprezinte grafic și funcția de repartiție de selecție, împreună cu funcția de repartiție ipotezată, adică cea de la legea normală $\mathcal{N}(\bar{x}, \bar{\sigma})$.

```
clf
n=input('n='); mu=input('mu='); s=input('sigma=');
x=normrnd(mu,s,1,n);
[h,c,ls]=lillietest(x);
fprintf(' h=%d\n c= %5.4f\n ls=%5.4f',h,c,ls)
ma=mean(x); va=sqrt(var(x));
cdfplot(x), hold on, title('')
t=ma-3*va:0.01:ma+3*va;
plot(t,normcdf(t,ma,va),'k--')
```

În urma executării programului pentru $n=10$, $\mu=100$ și $\sigma=2$, se obțin rezultatele:

```
h=0
c= 0.1433
ls=0.1093
```

și graficele din Figura 6.5

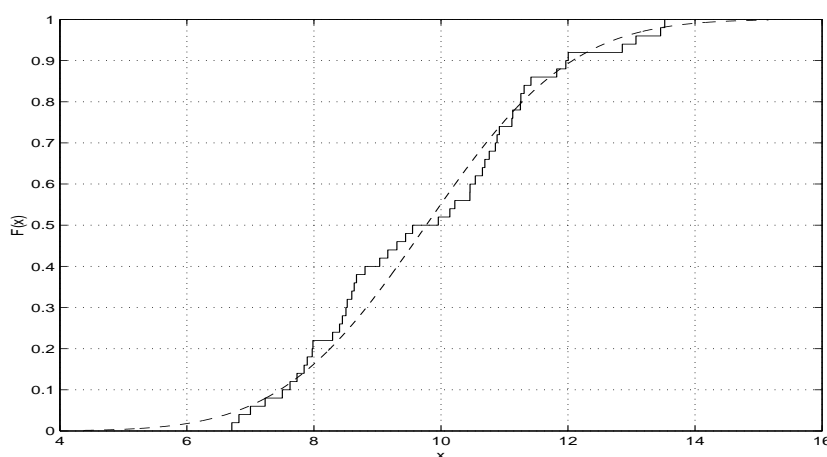


Figura 6.5: Testul Lilliefors

6.13 Testul Kolmogorov–Smirnov

Se consideră două caracteristici independente X și Y de tip continuu cu funcțiile teoretice de repartiție F_X și F_Y necunoscute. Privind cele două caracteristici vrem să verificăm dacă sunt identic repartizate. Astfel, considerăm ipoteza nulă dată prin $H_0: F_X = F_Y$, cu una din alternativele:

- (1) $H_1 : F_X \neq F_Y$ (testul Kolmogorov–Smirnov bilateral),
- (2) $H_1 : F_X > F_Y$ (testul Kolmogorov–Smirnov unilateral dreapta),
- (3) $H_1 : F_X < F_Y$ (testul Kolmogorov–Smirnov unilateral stânga).

Pentru verificarea ipotezei nule H_0 , cu una din alternativele precizate, se consideră câte o selecție repetată de volum n_1 și respectiv n_2 , cu datele de selecție x_1, x_2, \dots, x_{n_1} și variabilele de selecție X_1, X_2, \dots, X_{n_1} , corespunzătoare pentru caracteristica X și datele de selecție y_1, y_2, \dots, y_{n_2} și variabilele de selecție Y_1, Y_2, \dots, Y_{n_2} , corespunzătoare pentru caracteristica Y .

Se consideră statisticile

$$\begin{aligned} D_{n_1 n_2} &= \sup_{x \in \mathbb{R}} \{ |\bar{F}_X(x) - \bar{F}_Y(x)| \}, \\ D_{n_1 n_2}^+ &= \sup_{x \in \mathbb{R}} \{ \bar{F}_X(x) - \bar{F}_Y(x) \}, \\ D_{n_1 n_2}^- &= \sup_{x \in \mathbb{R}} \{ \bar{F}_Y(x) - \bar{F}_X(x) \}, \end{aligned}$$

unde \bar{F}_X și \bar{F}_Y sunt funcțiile de repartiție de selecție.

Pentru $n_1, n_2 \rightarrow \infty$, funcțiile de repartiție ale acestor funcții de selecție au următoarele comportări asimptotice:

$$\lim_{n_1, n_2 \rightarrow \infty} P \left(\sqrt{\frac{n_1 n_2}{n_1 + n_2}} D_{n_1 n_2} \leq x \right) = K(x) = \sum_{k=-\infty}^{k=+\infty} (-1)^k e^{-2k^2 x^2}, \quad x > 0,$$

funcția lui Kolmogorov $K(x)$, fiind tabelată pentru anumite valori în *Anexa V*, respectiv

$$\begin{aligned} \lim_{n_1, n_2 \rightarrow \infty} P \left(\sqrt{\frac{n_1 n_2}{n_1 + n_2}} D_{n_1 n_2}^+ \leq x \right) &= \lim_{n_1, n_2 \rightarrow \infty} P \left(\sqrt{\frac{n_1 n_2}{n_1 + n_2}} D_{n_1 n_2}^- \leq x \right) \\ &= K^\pm(x) = 1 - e^{-2x^2}, \quad x > 0, \end{aligned}$$

ce poartă numele de lege χ cu două grade de libertate.

Pentru un nivel de semnificație $\alpha \in (0, 1)$ fixat, se pot determina cuantilele $k_{1-\alpha}$, $k_{1-\alpha}^+$ și $k_{1-\alpha}^-$ astfel încât

$$P \left(\sqrt{\frac{n_1 n_2}{n_1 + n_2}} D_{n_1 n_2} \leq k_{1-\alpha} \mid H_0 \right) = 1 - \alpha, \quad \text{adică} \quad K(k_{1-\alpha}) = 1 - \alpha,$$

pentru testul bilateral Kolmogorov–Smirnov,

$$P \left(\sqrt{\frac{n_1 n_2}{n_1 + n_2}} D_{n_1 n_2}^+ \leq k_{1-\alpha}^+ \mid H_0 \right) = 1 - \alpha, \quad \text{adică} \quad K^+(k_{1-\alpha}^+) = 1 - \alpha,$$

pentru testul unilateral dreapta Kolmogorov–Smirnov,

$$P\left(\sqrt{\frac{n_1 n_2}{n_1 + n_2}} D_{n_1 n_2}^- \leq k_{1-\alpha}^- \mid H_0\right) = 1 - \alpha, \text{ adică } K^\pm(k_{1-\alpha}^-) = 1 - \alpha,$$

pentru testul unilateral stânga Kolmogorov–Smirnov,

Corespunzător celor trei alterantive, ipoteza H_0 va fi admisă când valorile calculate $d_{n_1 n_2}$, $d_{n_1 n_2}^+$, $d_{n_1 n_2}^-$ pe baza datelor de selecție x_1, \dots, x_{n_1} și y_1, \dots, y_{n_2} , ale statisticilor $D_{n_1 n_2}$, $D_{n_1 n_2}^+$, satisfac respectiv condițiile

$$(1) \quad \sqrt{\frac{n_1 n_2}{n_1 + n_2}} d_{n_1 n_2} < k_{1-\alpha}, \text{ când } H_1 : F_X \neq F_Y,$$

$$(2) \quad \sqrt{\frac{n_1 n_2}{n_1 + n_2}} d_{n_1 n_2}^+ < k_{1-\alpha}^+, \text{ când } H_1 : F_X > F_Y,$$

$$(3) \quad \sqrt{\frac{n_1 n_2}{n_1 + n_2}} d_{n_1 n_2}^- < k_{1-\alpha}^-, \text{ când } H_1 : F_X < F_Y,$$

iar în caz contrar va fi respinsă.

Prezentăm, în cele ce urmează, împreună cele trei teste (bilateral, unilateral dreapta și unilateral stânga) Kolmogorov–Smirnov. Remarcăm că cele trei teste sunt teste distincte, alegerea unuia dintre ele se face apriori.

Etapele aplicării testului lui Kolmogorov–Smirnov

1. Se dau: α ; x_1, x_2, \dots, x_{n_1} ; y_1, y_2, \dots, y_{n_2} ;

2. Se determină

$$(1) \quad k_{1-\alpha} \text{ astfel încât } K(k_{1-\alpha}) = 1 - \alpha, \text{ când } H_1 : F_X \neq F_Y,$$

$$(2) \quad k_{1-\alpha}^+ \text{ astfel încât } K^\pm(k_{1-\alpha}^+) = 1 - \alpha, \text{ când } H_1 : F_X > F_Y;$$

$$(3) \quad k_{1-\alpha}^- \text{ astfel încât } K^\pm(k_{1-\alpha}^-) = 1 - \alpha, \text{ când } H_1 : F_X < F_Y;$$

3. Se calculează

$$(1) \quad k = \sqrt{\frac{n_1 n_2}{n_1 + n_2}} d_{n_1 n_2}, \text{ când } H_1 : F_X \neq F_Y,$$

$$(2) \quad k = \sqrt{\frac{n_1 n_2}{n_1 + n_2}} d_{n_1 n_2}^+, \text{ când } H_1 : F_X > F_Y;$$

$$(3) \quad k = \sqrt{\frac{n_1 n_2}{n_1 + n_2}} d_{n_1 n_2}^-, \text{ când } H_1 : F_X < F_Y;$$

4. Concluzia: dacă

- (1) $k < k_{1-\alpha}$ ipoteza H_0 este admisă, când $H_1 : F_X \neq F_Y$,
- (2) $k < k_{1-\alpha}^+$ ipoteza H_0 este admisă, când $H_1 : F_X > F_Y$,
- (3) $k < k_{1-\alpha}^-$ ipoteza H_0 este admisă, când $H_1 : F_X < F_Y$,

Observația 6.13.1. Având în vedere că funcțiile (de repartiție) K și K^\pm sunt monoton crescătoare și că pentru determinarea cuantilelor $k_{1-\alpha}$, $k_{1-\alpha}^+$ și $k_{1-\alpha}^-$ este necesară inversarea acestor funcții de repartiție, se poate evita această operație. Pentru aceasta se determină valorile funcțiilor K și K^\pm pe valorile calculate ale statisticilor $D_{n_1 n_2}$, $D_{n_1 n_2}^+$, $D_{n_1 n_2}^-$, adică se calculează, respectiv $1 - c = K\left(\sqrt{\frac{n_1 n_2}{n_1 + n_2}} d_{n_1 n_2}\right)$, $1 - c = K^\pm\left(\sqrt{\frac{n_1 n_2}{n_1 + n_2}} d_{n_1 n_2}^\pm\right)$, $1 - c = K^\pm\left(\sqrt{\frac{n_1 n_2}{n_1 + n_2}} d_{n_1 n_2}^\pm\right)$. Astfel, se va respinge ipoteza nulă H_0 dacă valoarea critică c satisface inegalitatea $c \leq \alpha$.

6.13.1 Funcția `kstest2`

Statistics toolbox conține funcția `kstest2`, care efectuează testul Kolmogorov–Smirnov. Apelul funcției se poate face cu una din instrucțiunile:

```
h=kstest(x,y)
h=kstest(x,y,alpha)
h=kstest(x,y,alpha,tail)
[h,c,ks]=kstest(x,y,alpha,tail)
```

Comenzile de acest tip lansează execuția testului Kolmogorov–Smirnov bilateral (când `tail=0`, valoare implicită), unilateral dreapta (când `tail=1`), respectiv unilateral stânga (când `tail=-1`), prin considerarea datelor conținute în vectorii `x` și `y`.

Parametrul `alpha` (implicit `alpha=0.05`) reprezintă nivelul de semnificație.

Rezultatele obținute au următoarele semnificații. Dacă `h=1`, ipoteza nulă se respinge, ceea ce corespunde faptului că valoarea critică c satisface relația $c \leq \alpha$, iar dacă `h=0`, ipoteza nulă nu poate fi respinsă. Mai remarcăm faptul că `ks` va conține respectiv valoarea calculată a statisticilor D_{n_1, n_2} , D_{n_1, n_2}^+ și D_{n_1, n_2}^- , și verifică relațiile $c = 1 - K(\sqrt{n}d_n)$, pentru testul bilateral, respectiv $c = 1 - K^\pm(\sqrt{n}d_n^\pm)$ și $c = 1 - K^\pm(\sqrt{n}d_n^\pm)$, pentru testele unilaterale.

Programul 6.13.2. Să aplicăm testele Kolmogorov–Smirnov (bilateral, unilateral dreapta și unilateral stânga) pentru datele `x=-1:1:5` și `y` reprezentând 20 de numere aleatoare, ce urmează legea normală standard, iar nivelul de semnificație să fie `alpha=0.05`. Programul ce urmează, să reprezinte grafic și cele două funcții de repartiție de selecție.

```
clf
x=-1:5; y=randn(1,20);
fprintf(' h      c      ks\n')
fprintf('_____ \n')
for i=-1:1
```



```

[h,c,ks]=kstest2(x,y,[],i);
fprintf(' %d %5.4f %4.2f\n',h,c,ks)
end
h1=cdfplot(x); hold on, h2=cdfplot(y);
title('Funcțiile de repartitie de selecție')
set(h1,'linestyle','--')
set(h2,'linestyle','--')

```

În urma executării programului, se obțin rezultatele:

h	c	ks
1	0.0110	0.61
1	0.0219	0.61
0	0.9783	0.04

respectiv graficele din Figura 6.6. Testul Kolmogorov–Smirnov este bazat pe un re-

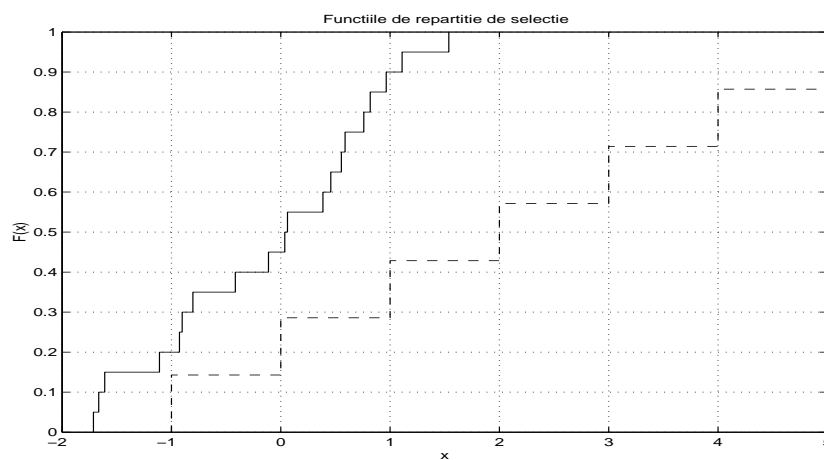


Figura 6.6: Testul Kolmogorov–Smirnov

zultat asimptotic, adică se aplică pentru valori mari ale volului selecției, dar aici s-au folosit selecții de volum mic, în mod special pentru a urmări pe figură modul de calcul a celor trei statistici d_{n_1, n_2} , d_{n_1, n_2}^+ și d_{n_1, n_2}^- .

Anexa I (Funcția lui Laplace) $\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_0^x e^{-\frac{t^2}{2}} dt, \Phi(-x) = -\Phi(x)$

[illegible]

$$\text{Anexa II (Legea Student)} F_n(t_{n;\gamma}) = \frac{\Gamma\left(\frac{n+1}{2}\right)}{\sqrt{n\pi}\Gamma\left(\frac{n}{2}\right)} \int_{-\infty}^{t_{n;\gamma}} \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}} dx = \gamma$$

γ	0.70	0.75	0.80	0.85	0.90	0.95	0.975	0.99	0.995
n	$t_{n;\gamma}$								
1	0.727	1.000	1.376	1.963	3.078	6.314	12.706	31.821	63.657
2	0.617	0.817	1.061	1.386	1.886	2.920	4.303	6.965	9.925
3	0.584	0.765	0.979	1.250	1.638	2.353	3.182	4.541	5.841
4	0.569	0.741	0.941	1.190	1.533	2.132	2.776	3.747	4.604
5	0.559	0.727	0.920	1.156	1.476	2.015	2.571	3.365	4.032
6	0.553	0.718	0.906	1.134	1.440	1.943	2.447	3.143	3.707
7	0.549	0.711	0.896	1.119	1.415	1.895	2.365	2.998	3.500
8	0.546	0.706	0.889	1.108	1.397	1.860	2.306	2.897	3.355
9	0.544	0.703	0.883	1.100	1.383	1.833	2.262	2.821	3.250
10	0.542	0.700	0.879	1.093	1.372	1.813	2.228	2.764	3.169
11	0.540	0.697	0.876	1.088	1.363	1.796	2.201	2.718	3.106
12	0.539	0.696	0.873	1.083	1.356	1.782	2.179	2.681	3.055
13	0.538	0.694	0.870	1.080	1.350	1.771	2.160	2.650	3.012
14	0.537	0.692	0.868	1.076	1.345	1.761	2.145	2.625	2.977
15	0.536	0.691	0.866	1.074	1.341	1.753	2.131	2.603	2.947
16	0.535	0.690	0.865	1.071	1.337	1.746	2.120	2.584	2.921
17	0.534	0.689	0.863	1.069	1.333	1.740	2.110	2.567	2.898
18	0.534	0.688	0.862	1.067	1.330	1.734	2.101	2.552	2.878
19	0.533	0.688	0.861	1.066	1.328	1.729	2.093	2.540	2.861
20	0.533	0.687	0.860	1.064	1.325	1.725	2.086	2.528	2.845
21	0.533	0.686	0.859	1.063	1.323	1.721	2.080	2.518	2.831
22	0.532	0.686	0.858	1.061	1.321	1.717	2.074	2.508	2.819
23	0.532	0.685	0.858	1.060	1.320	1.714	2.069	2.500	2.807
24	0.531	0.685	0.857	1.059	1.318	1.711	2.064	2.492	2.797
25	0.531	0.684	0.856	1.058	1.316	1.708	2.060	2.485	2.787
26	0.531	0.684	0.856	1.058	1.315	1.706	2.056	2.479	2.779
27	0.531	0.684	0.855	1.057	1.314	1.703	2.052	2.473	2.771
28	0.530	0.683	0.855	1.056	1.313	1.701	2.048	2.467	2.763
29	0.530	0.683	0.854	1.055	1.311	1.699	2.045	2.462	2.756
30	0.530	0.683	0.854	1.055	1.310	1.697	2.042	2.457	2.750
35	0.529	0.682	0.852	1.052	1.306	1.690	2.030	2.438	2.724
40	0.529	0.681	0.851	1.050	1.303	1.684	2.021	2.423	2.705
50	0.528	0.679	0.849	1.047	1.299	1.676	2.009	2.403	2.678
60	0.527	0.679	0.848	1.046	1.296	1.671	2.000	2.390	2.660
80	0.526	0.679	0.847	1.044	1.291	1.671	2.000	2.390	2.660
120	0.526	0.677	0.845	1.041	1.289	1.658	1.980	2.358	2.617
∞	0.524	0.674	0.842	1.036	1.282	1.645	1.956	2.326	2.576

Anexa III (Legea χ^2)

$$F_n(\chi_{n;\gamma}^2) = \frac{1}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})} \int_0^{\chi_{n;\gamma}^2} x^{\frac{n}{2}-1} e^{-\frac{x}{2}} dx = \gamma$$

γ	0.005	0.010	0.025	0.050	0.100	0.900	0.950	0.975	0.990	0.995
n	$\chi_{n;\gamma}^2$									
1	0.000	0.000	0.001	0.004	0.016	2.71	3.84	5.02	6.63	7.88
2	0.010	0.020	0.051	0.103	0.211	4.61	5.99	7.38	9.21	10.60
3	0.072	0.115	0.216	0.352	0.584	6.25	7.81	9.35	11.34	12.84
4	0.207	0.297	0.484	0.711	1.06	7.78	9.49	11.14	13.28	14.86
5	0.412	0.554	0.831	1.15	1.61	9.24	11.07	12.83	15.09	16.75
6	0.676	0.872	1.24	1.64	2.20	10.64	12.59	14.45	16.81	18.55
7	0.989	1.24	1.69	2.17	2.83	12.02	14.07	16.01	18.48	20.28
8	1.34	1.65	2.18	2.73	3.49	13.36	15.51	17.53	20.09	21.95
9	1.73	2.09	2.70	3.33	4.17	14.68	16.92	19.02	21.67	23.59
10	2.16	2.56	3.25	3.94	4.87	15.99	18.31	20.48	23.21	25.19
11	2.60	3.05	3.82	4.57	5.58	17.28	19.68	21.92	24.72	26.76
12	3.07	3.57	4.40	5.23	6.30	18.55	21.03	23.34	26.22	28.30
13	3.57	4.11	5.01	5.89	7.04	19.81	22.36	24.74	27.69	29.82
14	4.07	4.66	5.63	6.57	7.79	21.06	23.68	26.12	29.14	31.32
15	4.60	5.23	6.26	7.26	8.55	22.31	25.00	27.49	30.58	32.80
16	5.14	5.81	6.91	7.96	9.31	23.54	26.30	28.85	32.00	34.27
17	5.70	6.41	7.56	8.67	10.09	24.77	27.59	30.19	33.41	35.72
18	6.26	7.01	8.23	9.39	10.86	25.99	28.87	31.53	34.81	37.16
19	6.84	7.63	8.91	10.1	11.65	27.20	30.14	32.85	36.19	38.58
20	7.43	8.26	9.59	10.9	12.44	28.41	31.41	34.17	37.57	40.00
21	8.03	8.90	10.3	11.6	13.24	29.62	32.67	35.48	38.93	41.40
22	8.64	9.54	11.0	12.3	14.04	30.81	33.92	36.78	40.29	42.80
23	9.26	10.2	11.7	13.1	14.85	32.01	35.17	38.08	41.64	44.18
24	9.89	10.9	12.4	13.8	15.66	33.20	36.42	39.36	42.98	45.56
25	10.5	11.5	13.1	14.6	16.47	34.38	37.65	40.65	44.31	46.93
26	11.2	12.2	13.8	15.4	17.29	35.56	38.89	41.92	45.64	48.29
27	11.8	12.9	14.6	16.2	18.11	36.74	40.11	43.19	46.96	49.64
28	12.5	13.6	15.3	16.9	18.94	37.92	41.34	44.46	48.28	50.99
29	13.1	14.3	16.0	17.7	19.77	39.09	42.56	45.72	49.59	52.34
30	13.8	15.0	16.8	18.5	20.60	40.26	43.77	46.98	50.89	53.67
35	17.2	18.5	20.6	22.5	24.8	46.1	49.8	53.2	57.3	60.3
40	20.7	22.2	24.4	26.5	29.1	51.8	55.8	59.3	63.7	6.68
60	35.5	37.5	40.5	43.2	46.5	74.4	79.1	83.3	88.4	91.9

Anexa IV (Legea Fisher-Snedecor)

$$F_{m,n}(f_{m,n;\gamma}) = \left(\frac{m}{n}\right)^{\frac{m}{2}} \frac{\Gamma\left(\frac{m+n}{2}\right)}{\Gamma\left(\frac{m}{2}\right) \Gamma\left(\frac{n}{2}\right)} \int_0^{f_{m,n;\gamma}} x^{\frac{m}{2}-1} \left(1 + \frac{m}{n}x\right)^{-\frac{m+n}{2}} dx = \gamma$$

γ	m	1	2	3	4	5	6	7	8
	n	$f_{m,n;\gamma}$							
0.95	1	161.4	199.5	216	225	230	234	237	239
0.975		648	800	864	900	922	937	948	957
0.99		4052	4999	5403	5625	5764	5859	5930	5981
0.95	2	18.51	19.00	19.16	19.25	19.30	19.33	19.35	19.37
0.975		38.5	39.0	39.2	39.2	39.3	39.3	39.4	39.4
0.99		98.49	99.00	99.17	99.25	99.30	99.33	99.35	99.36
0.95	3	10.13	9.55	9.28	9.12	9.01	8.94	8.89	8.85
0.975		17.4	16.0	15.1	15.4	14.9	14.7	14.6	14.5
0.99		34.12	30.84	29.46	28.71	28.24	27.91	27.7	27.49
0.95	4	17.71	6.94	6.59	6.39	6.26	6.16	6.09	6.04
0.975		12.2	10.6	9.98	9.60	9.36	9.20	9.07	8.98
0.99		21.20	18.00	16.69	15.98	15.52	15.21	15.0	14.80
0.95	5	6.61	5.79	5.41	5.19	5.05	4.95	4.88	4.82
0.975		10.0	8.43	7.76	7.39	7.15	6.98	6.85	6.76
0.99		16.26	13.27	12.06	11.39	10.97	10.67	10.5	10.27
0.95	6	5.99	5.14	4.76	4.53	4.39	4.28	4.21	4.15
0.975		8.81	7.26	6.60	6.23	5.99	5.82	5.70	5.60
0.99		13.74	10.91	9.78	9.15	8.75	8.47	8.26	8.10
0.95	7	5.59	4.74	4.35	4.12	3.97	3.87	3.79	3.73
0.975		8.07	6.54	5.89	5.52	5.29	5.12	4.99	4.90
0.99		12.25	9.55	8.45	7.85	7.45	7.19	6.99	6.84
0.95	8	5.32	4.46	4.07	3.84	3.69	3.58	3.50	3.44
0.975		7.57	6.06	5.42	5.05	4.82	4.65	4.53	4.43
0.99		11.26	8.65	7.59	7.01	6.63	6.37	6.18	6.03
0.95	9	5.12	4.26	3.86	3.63	3.48	3.37	3.29	3.23
0.975		7.21	5.71	5.08	4.72	4.48	4.32	4.20	4.10
0.99		10.56	8.02	6.99	6.42	6.06	5.80	5.61	5.47
0.95	10	4.96	4.10	3.71	3.48	3.33	3.22	3.14	3.07
0.975		6.94	5.46	4.83	4.47	4.24	4.07	3.95	3.85
0.99		10.04	7.56	6.55	5.99	5.64	5.39	5.20	5.06
0.95	11	4.84	3.98	3.59	3.36	3.20	3.09	3.01	2.95
0.975		6.72	5.26	4.63	4.28	4.04	3.88	3.76	3.66
0.99		9.65	7.20	6.22	5.67	5.32	5.07	4.89	4.74

Anexa IV (continua)

γ	m	1	2	3	4	5	6	7	8
γ	n	$f_{m,n,\gamma}$							
0.95	12	4.75	3.88	3.49	3.26	3.11	3.00	2.91	2.85
0.975		6.55	5.10	4.47	4.12	3.89	3.73	3.61	3.51
0.99		9.33	6.93	5.95	5.41	5.06	4.82	4.54	4.50
0.95	13	4.67	3.81	3.41	3.18	3.03	2.92	2.83	2.77
0.975		6.41	4.97	4.35	4.00	3.77	3.60	3.48	3.39
0.99		9.07	6.70	5.74	5.21	4.86	4.62	4.44	4.30
0.95	14	4.60	3.74	3.34	3.11	2.96	2.85	2.76	2.70
0.975		6.30	4.86	4.24	3.89	3.66	3.50	3.38	3.29
0.99		8.86	6.51	5.56	5.04	4.70	4.46	4.28	4.14
0.95	15	4.54	3.68	3.29	3.06	2.90	2.79	2.71	2.64
0.975		6.20	4.76	4.15	3.80	3.58	3.41	3.29	3.20
0.99		8.68	6.36	5.42	4.89	4.56	4.32	4.14	4.00
0.95	16	4.49	3.63	3.24	3.01	2.85	2.74	2.66	2.59
0.975		6.12	4.69	4.08	3.73	3.50	3.34	3.22	3.12
0.99		8.53	6.23	5.29	4.77	4.44	4.20	4.03	3.89
0.95	17	4.45	3.59	3.20	2.96	2.81	2.70	2.61	2.55
0.975		6.04	4.62	4.01	3.66	3.44	3.28	3.16	3.06
0.99		8.40	6.11	5.18	4.67	4.34	4.10	3.93	3.79
0.95	18	4.41	3.55	3.16	2.93	2.77	2.66	2.58	2.51
0.975		5.98	4.56	3.95	3.61	3.38	3.22	3.10	3.01
0.99		8.29	7.01	5.09	4.58	4.25	4.01	3.84	3.71
0.95	19	4.38	3.52	3.13	2.90	2.74	2.63	2.54	2.48
0.975		5.92	4.51	3.90	3.56	3.33	3.17	3.05	2.96
0.99		8.18	5.93	5.01	4.50	4.17	3.94	3.77	3.63
0.95	20	4.35	3.49	3.10	2.87	2.71	2.60	2.51	2.45
0.975		5.87	4.46	3.86	3.51	3.29	3.13	3.01	2.91
0.99		8.10	5.85	4.94	4.43	4.10	3.87	3.70	3.56
0.95	21	4.32	3.47	3.07	2.84	2.68	2.57	2.49	2.42
0.975		5.83	4.42	3.82	3.48	3.25	3.09	2.97	2.87
0.99		8.02	5.78	4.87	4.37	4.04	3.84	3.64	3.51
0.95	22	4.30	3.44	3.05	2.82	2.66	2.55	2.46	2.40
0.975		5.79	4.38	3.78	3.44	3.22	3.05	2.93	2.84
0.99		7.95	5.72	4.82	4.31	3.99	3.76	3.59	3.45
0.95	23	4.28	3.42	3.03	2.80	2.64	2.53	2.44	2.37
0.975		5.75	4.35	3.75	3.41	3.18	3.02	2.90	2.81
0.99		7.88	5.66	4.76	4.26	3.94	3.71	3.54	3.41
0.95	24	4.26	3.40	3.01	2.78	2.62	2.51	2.42	2.36
0.975		5.72	4.32	3.72	3.38	3.15	2.99	2.87	2.78
0.99		7.82	5.61	4.72	4.22	3.90	3.67	3.50	3.36

Anexa IV (continuare)

γ	m	9	10	11	12	13	14	15	16
	n	$f_{m,n;\gamma}$							
0.95	2	19.4	19.4	19.4	19.4	19.4	19.4	19.4	19.4
0.975		39.4	39.4	39.4	39.4	39.4	39.4	39.4	39.4
0.99		99.4	99.4	99.4	99.4	99.4	99.4	99.4	99.4
0.95	3	8.81	8.79	8.76	8.74	8.73	8.71	8.70	8.69
0.975		14.5	14.4	14.4	14.3	14.3	14.3	14.3	14.2
0.99		27.3	27.1	27.1	27.1	27.0	26.9	26.9	26.8
0.95	4	6.00	5.96	5.94	5.91	5.89	5.87	5.86	5.84
0.975		8.90	8.84	8.79	8.75	8.72	8.69	8.66	8.64
0.99		14.7	14.5	14.4	14.4	14.3	14.2	14.2	14.2
0.95	5	4.77	4.74	4.70	4.68	4.66	4.64	4.62	4.60
0.975		6.68	6.62	6.57	6.52	6.49	6.46	6.43	6.41
0.99		10.2	10.1	9.96	9.89	9.82	9.77	9.72	9.68
0.95	6	4.10	4.06	4.03	4.00	3.98	3.96	3.94	3.92
0.975		5.52	5.46	5.41	5.37	5.33	5.30	5.27	5.25
0.99		7.98	7.89	7.79	7.72	7.66	7.60	7.56	7.52
0.95	7	3.68	3.64	3.60	3.57	3.55	3.53	3.51	3.49
0.975		4.82	4.76	4.71	4.67	4.63	4.60	4.57	4.54
0.99		6.72	6.62	6.54	6.47	6.41	6.36	6.31	6.27
0.95	8	3.39	3.35	3.31	3.28	3.26	3.24	3.22	3.20
0.975		4.36	4.30	4.24	4.20	4.16	4.13	4.10	4.08
0.99		5.91	5.81	5.73	5.67	5.61	5.56	5.52	5.48
0.95	9	3.18	3.14	3.10	3.07	3.05	3.03	3.01	2.99
0.975		4.03	3.96	3.91	3.87	3.83	3.80	3.77	3.74
0.99		5.35	5.26	5.18	5.11	5.05	5.00	4.96	4.92
0.95	10	3.02	2.98	2.94	2.91	2.89	2.86	2.85	2.83
0.975		3.78	3.72	3.66	3.62	3.58	3.55	3.52	3.50
0.99		4.94	4.85	4.77	4.71	4.65	4.60	4.56	4.52
0.95	11	2.90	2.85	2.82	2.79	2.76	2.74	2.72	2.70
0.975		3.59	3.53	3.47	3.43	3.39	3.36	3.33	3.30
0.99		4.63	4.54	4.46	4.40	4.34	4.29	4.25	4.21
0.95	12	2.80	2.75	2.72	2.69	2.66	2.64	2.62	2.60
0.975		3.44	3.37	3.32	3.28	3.24	3.21	3.18	3.15
0.99		4.39	4.30	4.22	4.16	4.10	4.05	4.01	3.97
0.95	13	2.71	2.67	2.63	2.60	2.58	2.55	2.53	2.51
0.975		3.31	3.25	3.20	3.15	3.12	3.08	3.05	3.03
0.99		4.19	4.10	4.02	3.96	3.91	3.86	3.82	3.78
0.95	14	2.65	2.60	2.57	2.53	2.51	2.48	2.46	2.44
0.975		3.21	3.15	3.09	3.05	3.01	2.98	2.95	2.92
0.99		4.03	3.94	3.86	3.80	3.75	3.70	3.66	3.62

Anexa IV (continuare)

	m	9	10	11	12	13	14	15	16
γ	n	$f_{m,n;\gamma}$							
0.95	15	2.59	2.54	2.51	2.48	2.45	2.42	2.40	2.38
0.975		3.12	3.06	3.01	2.96	2.92	2.89	2.86	2.84
0.99		3.89	3.80	3.73	3.67	3.61	3.56	3.52	3.49
0.95	16	2.54	2.49	2.46	2.42	2.40	2.37	2.35	2.33
0.975		3.05	2.99	2.93	2.89	2.85	2.82	2.79	2.76
0.99		3.78	3.69	3.62	3.55	3.50	3.45	3.41	3.37
0.95	17	2.49	2.45	2.41	2.38	2.35	2.33	2.29	2.27
0.975		2.98	2.92	2.87	2.82	2.79	2.75	2.72	2.70
0.99		3.68	3.59	3.52	3.46	3.40	3.35	3.31	3.27
0.95	18	2.46	2.41	2.37	2.34	2.31	2.29	2.27	2.25
0.975		2.93	2.87	2.81	2.77	2.73	2.70	2.67	2.64
0.99		3.60	3.51	3.43	3.37	3.32	3.27	3.23	3.19
0.95	19	2.42	2.38	2.34	2.31	2.28	2.26	2.23	2.21
0.975		2.88	2.82	2.76	2.72	2.68	2.65	2.62	2.59
0.99		3.52	3.43	3.36	3.30	3.24	3.19	3.15	3.12
0.95	20	2.39	2.35	2.31	2.28	2.25	2.22	2.20	2.19
0.975		2.84	2.77	2.72	2.68	2.64	2.60	2.57	2.55
0.99		3.46	3.37	3.29	3.23	3.18	3.13	3.09	3.05
0.95	21	2.37	2.32	2.28	2.25	2.22	2.20	2.18	2.16
0.975		2.80	2.73	2.68	2.64	2.60	2.56	2.53	2.51
0.99		3.40	3.31	3.24	3.17	3.12	3.07	3.03	2.99
0.95	22	2.34	2.30	2.26	2.23	2.20	2.17	2.15	2.13
0.975		2.76	2.70	2.65	2.60	2.56	2.53	2.50	2.47
0.99		3.35	3.26	3.18	3.12	3.07	3.02	2.98	2.94
0.95	23	2.32	2.27	2.23	2.20	2.18	2.15	2.13	2.11
0.975		2.73	2.67	2.62	2.57	2.53	2.50	2.47	2.44
0.99		3.30	3.21	3.14	3.07	3.02	2.97	2.93	2.89
0.95	24	2.30	2.25	2.21	2.18	2.15	2.13	2.11	2.09
0.975		2.70	2.64	2.59	2.54	2.50	2.47	2.44	2.41
0.99		3.26	3.17	3.09	3.03	2.98	2.93	2.89	2.85

$$K(x) = \sum_{k=-\infty}^{+\infty} (-1)^k e^{-2k^2 x^2}$$
[illegible]

Bibliografie

- [1] Abramowitz, M., Stegun, I. A., *Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables*, tenth ed., Dover Publications, Inc., New York, 1972.
- [2] Blaga, P., *Calculul probabilităților. Culegere de probleme*, Universitatea "Babeș-Bolyai", Cluj-Napoca, 1984.
- [3] ———, *Metode statistice în modelarea cu calculatorul. Lucrări de laborator*, Universitatea "Babeș-Bolyai", Cluj-Napoca, 1993.
- [4] ———, *Calculul probabilităților și statistică matematică. Vol. II. Curs și culegere de probleme*, Universitatea "Babeș-Bolyai", Cluj-Napoca, 1994.
- [5] ———, *Statistică matematică. Lucrări de laborator*, Universitatea "Babeș-Bolyai", Cluj-Napoca, 1999.
- [6] ———, *Statistică matematică*, Universitatea "Babeș-Bolyai", Cluj-Napoca, 2000.
- [7] ———, *Statistică matematică (ediția II)*, Universitatea "Babeș-Bolyai", Cluj-Napoca, 2001.
- [8] Blaga, P., Lupaș, A., Mureșan A. S., *Matematici aplicate, Vol.I-II*, Promedia Plus, Cluj-Napoca, 1999.
- [9] Blaga, P., Mureșan A. S., *Matematici aplicate în economie, Vol.I-II*, Transilvania Press, Cluj-Napoca, 1996.
- [10] Blaga, P., Rădulescu, M., *Calculul probabilităților*, Universitatea "Babeș-Bolyai", Cluj-Napoca, 1987.
- [11] Ciucu, G., *Elemente de teoria probabilităților și statistică matematică*, Editura didactică și pedagogică, București, 1963.

- [12] Ciucu, G., Craiu, V., *Probleme de statistică matematică*, Editura didactică și pedagogică, București, 1968.
- [13] ———, *Teoria estimăției și verificarea ipotezelor statistice*, Editura didactică și pedagogică, București, 1968.
- [14] ———, *Introducere în teoria probabilităților și statistică matematică*, Editura didactică și pedagogică, București, 1971.
- [15] ———, *Inferență statistică*, Editura didactică și pedagogică, București, 1974.
- [16] Ciucu, G., Craiu, V., Ștefănescu, A., *Statistică matematică și cercetări operaționale, Vol. I*, Editura didactică și pedagogică, București, 1974.
- [17] Ciucu, G., Tudor, C., *Probabilități și procese stochastice. Vol. I*, Editura Academiei, București, 1978.
- [18] Craiu, V., *Verificarea ipotezelor statistice*, Editura științifică și enciclopedică, București, 1972.
- [19] Craiu, V., Enache, R., Bâscă, O., *Teste de concordanță cu programe în FORTRAN*, Editura științifică și enciclopedică, București, 1974.
- [20] Cramér, H., *Mathematical methods of statistics*, Princeton University Press, Princeton, 1951.
- [21] Dacunha–Castelle, D., Duflo, M., *Exercices de probabilités et statistiques. I. Problèmes à temps fixe. Tome I*, Masson, Paris, 1990.
- [22] ———, *Probabilités et statistiques. I. Problèmes à temps fixe. Tome I*, Masson, Paris, 1990.
- [23] Deák, I., *Random number generators and simulation*, Akadémiai Kiadó, Budapest, 1990.
- [24] Dumitrescu, M., Florea, D., Tudor, C., *Probleme de teoria probabilităților și statistică matematică*, Editura Tehnică, București, 1985.
- [25] Enders, W., *Applied econometric. Time series*, John Wiley & Sons, 1995.
- [26] Ermakov, S. M., *Metoda Monte Carlo și probleme înrudite*, Editura Tehnică, București, 1976.
- [27] Feller, W., *An introduction to probability theory and its applications. Vol. I*, John Wiley, New York, 1957.

- [28] ———, *An introduction to probability theory and its applications. Vol. II*, John Wiley, New York, 1966.
- [29] Gnedenko, B. V., *The theory of probability*, Mir Publishers, Moscow, 1976.
- [30] Good, Ph. I., *Resampling Methods. A Practical Guide to Data Analysis*, Birkhäuser, Boston • Basel • Berlin, 1999.
- [31] Griffiths, D. F., *An Introduction to Matlab, Version 2.1*, University of Dundee, 1997.
- [32] Grinstead, Ch. M., Snell, J. L., *Introduction in probability, Second edition*, American Mathematical Society, 1997.
- [33] Gurskiĭ, E. I., *Culegere de probleme pentru teoria probabilităților și statistica matematică (l. rusă)*, Edit. Înv. superior Minsk, Minsk, 1984.
- [34] Hoel, P. J., *Introduction to mathematical statistics, Fourth edition*, John Wiley & Sons, New York–London–Sydney–Toronto, 1971.
- [35] ———, *Statistique mathématique, Tome II*, Armand Colin, Paris, 1991.
- [36] Hoffmann–Jørgensen, J., *Probability with a view toward statistics, Vol. I-II*, Chapman & Hall, New York–London, 1994.
- [37] Iosifescu, M., Mihoc, Gh., Theodorescu, R., *Teoria probabilităților și statistica matematică*, Editura Tehnică, București, 1966.
- [38] Iosifescu, M., Moineanu, C., Trebici, V., Ursianu, E., *Mică enciclopedie de statistică*, Editura științifică și enciclopedică, București, 1985.
- [39] Jansson, B., *Random number generators*, Victor Petterson Bokindustri Aktiebolag, Stockholm, 1966.
- [40] Johnson, N. L., Kotz, S., *Distributions in statistics. Discrete distributions*, Houghton Mifflin Company, Boston, 1969.
- [41] Kendall, M. G., Stuart, A., *The advanced theory of statistics. Vol. 1. Distribution theory, Second edition*, Charles Griffin & Company Limited, London, 1961.
- [42] ———, *The advanced theory of statistics. Vol. 2. Inference and relationship*, Charles Griffin & Company Limited, London, 1963.
- [43] ———, *The advanced theory of statistics. Vol. 3. Design and analysis, and time-series*, Charles Griffin & Company Limited, London, 1966.

- [44] Knuth, D. E., *The art of computer programming, Vol. II. Seminumerical algorithms*, Addison-Wesley, Mass., 1969.
- [45] Lebart, L., Morineau, M. G., Fénelon, J.-P., *Traitement des données statistiques. méthodes et programmes*, Dunod, Paris, 1982.
- [46] Lecoutre, J.-P., Tassi, Ph., *Statistique non paramétrique et robustesse*, Economica, Paris, 1987.
- [47] Lehmann, E. L., *Testing statistical hypotheses, Second edition*, Springer, New York–Berlin, 1997.
- [48] Malița, M., Zidăroiu, C., *Incertitudine și decizie*, Editura științifică și enciclopedică, București, 1980.
- [49] Mihoc, Gh., Ciucu, G., Craiu, V., *Teoria probabilităților și statistică matematică*, Editura didactică și pedagogică, București, 1970.
- [50] Mihoc, Gh., Ciucu, G., Muja, A., *Modele matematice ale așteptării*, Editura Academiei, București, 1973.
- [51] Mihoc, Gh., Craiu, V., *Tratat de statistică matematică. Vol. I. Selecție și estimare*, Editura Academiei, București, 1976.
- [52] ———, *Tratat de statistică matematică. Vol. II. Verificarea ipotezelor statistice*, Editura Academiei, București, 1977.
- [53] Mihoc, Gh., Firescu, D., *Statistică matematică*, Editura didactică și pedagogică, București, 1966.
- [54] Mihoc, I., *Calculul probabilităților și statistică matematică. Part. I, II*, lito. Univ. “Babeș–Bolyai”, Cluj–Napoca, 1994, 1995.
- [55] Nakamura, S., *Numerical Analysis and Graphic Visualization with MATLAB*, Prentice Hall PTR, Upper Saddle River, New Jersey, 1996.
- [56] Oancea, E., Rădulescu, M., *Calculul probabilităților și statistică matematică*, lito. Univ. “Babeș–Bolyai”, Cluj–Napoca, 1974.
- [57] Ogden, R. T., *Essential wavelets for statistical applications and data analysis*, Birkhäuser, Boston • Basel • Berlin, 1997.
- [58] Onicescu, O., *Numere și sisteme aleatoare*, Editura Academiei, București, 1962.

- [59] Onicescu, O., Botez, M. C., *Incertitudine și modelare economică (econometrie informațională)*, Editura științifică și enciclopedică, București, 1985.
- [60] Parzen, E., *Modern probability theory and its applications*, John Wiley, New York–London, 1960.
- [61] Press, W. H., Flannery, B. P., Teukolsky, S. A., Vetterling, W. T., *Numerical Recipes. The Art of Scientific Computing*, Cambridge University Press, Cambridge, 1986.
- [62] Rancu, N., Tövissi, L., *Statistica matematică cu aplicații în producție*, Editura Academiei, București, 1963.
- [63] Rao, M. M., *Probability theory with applications*, Academic Press, New York, 1984.
- [64] Rényi, A., *Probability theory*, Akadémiai Kiadó, Budapest, 1970.
- [65] Rumşiski, L. Z., *Prelucrarea matematică a datelor experimentale*, Editura Tehnică, București, 1974.
- [66] Saporta, G., *Probabilités, analyse des données et statistique*, Éditions Technip, Paris, 1990.
- [67] Schervish, M. J., *Theory of statistics*, Springer, New York–Berlin, 1995.
- [68] Shiryaev, A. N., *Probability*, Springer, New York–Berlin, 1995.
- [69] Sigmon, K., *Matlab Primer, Third edition*, University of Florida, 1993.
- [70] Snedecor, G. W., Cochran, W. G., *Statistical methods*, Iowa State University Press, 1989.
- [71] Stapleton, J. H., *Linear statistical models*, John Wiley & Sons, New York–Chichester–Brisbane, 1995.
- [72] Stark, H., Woods, J. W., *Probability, random processes, and estimation theory for engineers*, Prentice-Hall, New Jersey, 1986.
- [73] Stoyanov, J., Mirazchiiski, I., Ignatov, Z., Tanushev, M., *Exercise manual in probability theory*, Kluwer Academic Publishers, Dordrecht–Boston–London, 1988.
- [74] Săcuiu, I., Zorilescu, D., *Numere aleatoare. Aplicații în economie, industrie și studiul fenomenelor naturale*, Editura Academiei, București, 1978.

-
- [75] Sveshnikov, A. A., *Problems in probability theory, mathematical statistics and theory of random functions*, W. B. Saunders Company, Philadelphia–London–Toronto, 1968.
- [76] Tassi, Ph., *Méthodes statistiques*, 2^e édition, Economica, Paris, 1989.
- [77] Văduva, I., *Analiză dispersională*, Editura Tehnică, București, 1970.
- [78] ———, *Modele de simulare cu calculatorul*, Editura Tehnică, București, 1977.
- [79] Wahba, G., *Spline models for observational data*, Society for Industrial and Applied Mathematics, Philadelphia, 1990.
- [80] Yule, G. U., Kendall, M. G., *Introducere în teoria statisticii*, Editura științifică, București, 1969.

Index

Abatere
 cuartilică, 152
 medie absolută, 152
 medie pătratică, 152
 standard, 107, 152
abs, 9
acos, 9
acosh, 9
acot, 9
acoth, 9
acsc, 9
acsch, 9
all, 15
Amplitudinea, 152
 clasei, 120
and, 14
angle, 9
ans, 12
any, 14
asec, 9
asech, 9
Asimetrie, 103, 154
asin, 9
asinh, 9
atan, 9
atan2, 9
atanh, 9
Atribut, 114
axis, 23

bar, 31, 65, 134, 138
bar3, 58
bar3h, 58
barh, 31
beta, 80
betafit, 282
betalike, 280
bino, 68
binofit, 282
boxplot, 164

Câmp de evenimente, 62
Câmp de probabilitate, 62
calendar, 6
Cantitate de informație, 232
Caracteristică, 114
 calitativă (atribut), 114
 cantitativă, 114
 continuă, 115
 discretă, 115
caseread, 127
casewrite, 127
cd, 4, 10
cdf, 65, 66, 74, 107
cdfplot, 223
ceil, 9
Centile, 104, 151
Centru de greutate, 174
chi2, 83
cla, 23
Clase, 120
clear, 11
clf, 23
clock, 6

- Coeficienții
 - lui Fisher, 154
 - lui Pearson, 154
- Coeficient de asimetrie
 - intercuartil, 154
- Coeficient de corelație, 95
 - al lui Pearson, 167
 - al rangurilor, 187
 - de selecție, 207
- Coeficient de regresie, 174
- Coeficient de variație, 154
 - intercuartil, 154
- Coeficientul
 - lui Kendall, 189
 - lui Spearman, 187
- Coeficientul de concordanță
 - al lui Kendall, 194
- Colectivitate, 114
- colormap, 27
- combnk, 10
- Concatenare, 9
- conj, 9
- contour, 54
- contourf, 55
- Convergență
 - în probabilitate, 108
 - în repartiție, 108
- Corecțiile
 - Sheppard, 158
- Corelație, 95, 165
- Corp borelian, 62
- corrcoef, 180
- cos, 9
- cosh, 9
- cot, 9
- coth, 9
- cov, 179
- Covarianță, 95
- crosstab, 126
- csc, 9
- csch, 9
- Cuantilă, 104
- Cuartila, 104, 150
 - inferioară, 104, 150
 - superioară, 104, 150
- cumprod, 157
- cumsum, 157
- Curbă de regresie, 103, 170
- date, 6
- Date de selecție, 198
- dbclear, 16
- dbcont, 16
- dbdown, 16
- dbmex, 17
- dbquit, 17
- dbstack, 16
- dbstatus, 16
- dbstep, 16
- dbstop, 16
- dbtype, 16
- dbup, 17
- Decile, 104, 151
- delete, 17
- demo, 4
- Densitate de probabilitate, 73, 74
 - marginală, 89
- det, 10
- diag, 10
- Diagramă
 - cumulativă ascendentă, 129
 - integrală cumulativă, 132
 - prin batoane (bare), 129
- diary, 11, 43
- diff, 157
- dir, 4
- disp, 10
- Dispersie, 95, 152
 - condiționată, 101, 170

- de selecție, 206
- Distribuție, 64
 - statistică, 123
- Distribuție marginală, 88
- disttool, 107
- Drepte de regresie, 173
- Eșantion, 197
- echo, 48
- edit, 16
- Eficiență, 242
- eps, 6
- eq, 14
- erf, 76
- Eroare
 - de speța întâi, 291
 - de speța a doua, 291
- Estimație, 236
 - absolut corectă, 237
 - consistentă, 236
 - corectă, 239
 - cu χ^2 minim, 256
 - de verosimilitate maximă, 252
 - nedeplasată, 236
- Estimator, 236
 - absolut corect, 237
 - consistent, 236
 - corect, 239
 - cu χ^2 minim, 256
 - de verosimilitate maximă, 252
 - nedeplasat, 236
 - optimal, 247
- Eveniment
 - cert, sigur, 62
- Exces, 103
- exit, 2
- exp, 9, 79
- Experiment, 62
- expfit, 282
- eye, 6
- ezcontour, 56
- ezcontourf, 57
- ezmesh, 57
- ezmeshc, 57
- ezplot, 35
- ezplot3, 51
- ezpolar, 25
- ezsurf, 57
- ezsurf, 57
- f, 85
- factor, 10
- factorial, 10
- feval, 47
- Fișiere script, 15, 43
- figure, 23
- file, 17
- fill, 35
- fix, 9
- floor, 9
- flops, 48
- format, 5
- Formula lui Daniels, 193
- fplot, 33
- fprintf, 13
- Frecvență
 - absolută, 120, 124
 - cumulată
 - ascendentă, 123
 - descendentă, 123
 - relativă, 123
- Funcția eroare, 76
- Funcția Laplace, 76, 77
- Funcția lui Euler
 - de speța întâi, 80
 - de speța a doua, 78
- Funcție caracteristică, 103
- Funcție de estimație, 236
 - absolut corectă, 237
 - corectă, 239

- eficientă, 242
- nedeplasată, 236
- Funcție de probabilitate, 65
- Funcție de repartiție, 63–65, 74
 - condiționată, 89
 - de selecție, 216
 - marginală, 88
- Funcție de selecție, 198
- Funcție de supraviețuire, 93
- Funcție de verosimilitate, 229
- Funcție hazard, 93
- Funcție nucleu, 131
 - Gaussian, 131
 - parabolic (Epanechnikov), 131
- Funcție wavelet, 29
 - mamă, 29
- Funcții Matlab, 43
- Funcțiile lui Haar, 29

- gam, 78
- gamfit, 282
- gamlike, 280
- gcd, 10
- ge, 14
- geo, 73
- geomean, 149
- ginput, 13
- gline, 181
- global, 44
- gplotmatrix, 146
- grid, 24
- grpstats, 163
- Grupare, 120
- gscatter, 144
- gt, 14
- gtext, 24

- harmmean, 149
- help, 2, 44
- helpdesk, 4

- helpwin, 3
- hist, 138, 139
- histc, 139, 140
- histfit, 278
- Histograma, 129
- hold, 22
- hyge, 69

- i, 6
- icdf, 104, 107
- imag, 9
- Indicatorul
 - de verosimilitate, 257
 - lui Berkson, 257
 - lui Kullbach, 257
 - lui Neyman, 257
 - lui Pearson, 257
- Indice
 - de dilatare, 29
 - de translatăre, 29
- inf, 6
- input, 13
- Instrucțiunea
 - break, 42
 - de atribuire, 12
 - error, 43
 - for, 40
 - if, 37
 - pause, 42
 - return, 42
 - switch, 38
 - try...catch, 41
 - while, 39
- Instrucțiuni
 - de citire, 13
 - de scriere, 13
- Interval
 - intercuartilic, 152
- Interval de încredere, 185, 187, 257
- inv, 10

- Ipoteză statistică, 285
 - alternativă, 285
 - compusă, 285
 - nulă, 285
 - simplă, 285
- iqr, 155
- j, 6
- kstest, 360
- kstest2, 366
- kurtosis, 155
- lcm, 10
- ldivide, 8
- le, 14
- Lege de probabilitate
 - de tip continuu, 74
 - clasică, 74
 - statistică, 81
 - de tip discret, 64
 - clasică, 67
- Legea
 - χ^2 , 83
 - necentrată, 85
 - arcsin, 80
 - Bernoulli, 67
 - beta, 80
 - binomială, 68
 - binomială negativă, 72
 - Cauchy, 81
 - evenimentelor rare, 71
 - exponențială, 71, 79
 - F (Fisher–Snedecor), 85
 - necentrată, 86
 - gamma, 78
 - geometrică, 73
 - hipergeometrică, 69
 - lognormală, 78
 - normală, 76
 - normală multidimensională, 90
 - Pascal, 72
 - Poisson, 70
 - Rayleigh, 80
 - t (Student), 81
 - necentrată, 82
 - uniformă, 74
 - discretă, 68
 - Weibull, 80
- Legea numerelor mari, 109
- legend, 24
- length, 10
- lillietest, 362
- load, 11
- log, 9
- log10, 9
- log2, 9
- logn, 78
- lookfor, 3, 44
- lsline, 181
- lt, 14
- mad, 155
- magic, 6
- Matrice vidă, 12
- Matricea covarianțelor, 95
- Matricea informației lui Fisher, 232
- max, 156
- mean, 149
- median, 150
- Mediana, 104, 149
- Medie
 - aritmetică, 148
 - armonică, 149
 - geometrică, 148
- Medie de selecție, 199
- mesh, 52
- meshgrid, 52
- Metoda ferestrei mobile, 130
- Metode de selecție
 - nerepetate, 198

- repetate(bernoulliene), 198
- min, 156
- minus, 8
- mkdir, 4, 10
- mldivide, 8
- mle, 278
- Mod, 151
- Moment, 151, 166
 - centrat, 103, 152, 167
 - centrat de selecție, 203
 - de selecție, 201
 - inițial, 103
- moment, 155
- mpower, 8
- mrdivide, 9
- mtimes, 8
- mvnrnd, 117
- mvtrnd, 117

- NaN, 6
- nanmax, 158
- nanmean, 158
- nanmedian, 158
- nanmin, 158
- nanstd, 158
- nansum, 158
- nargin, 44
- nargout, 44
- nbin, 72
- ncf, 86
- nchoosek, 10
- nct, 82
- ncx2, 85
- ndims, 10
- ne, 14
- Nivel de semnificație, 286
- nlinfit, 185
- nlintool, 186
- Nor statistic, 133
- norm, 76
- normfit, 282
- normlike, 280
- normspec, 87
- not, 14
- Numere
 - aleatoare, 116
 - pseudoaleatoare, 117
- Observare
 - curentă, 115
 - parțială, 115
 - periodică, 115
 - totală, 115
- ones, 6
- or, 14

- pdf, 65, 66, 74, 107
- perms, 10
- pi, 6
- pie, 147
- pie3, 147
- plot, 17, 25, 28, 65, 74, 138
- plot3, 49
- plotmatrix, 145
- plus, 8
- poiss, 70
- poissfit, 282
- polar, 25
- poly, 183
- polyfit, 182
- polytool, 184
- polyval, 183
- polyvalm, 183
- Populație, 114
- power, 8
- prctile, 151
- primes, 10
- Probă, 62
- Probabilitate, 62
- Probabilitate de încredere, 185

- Proceduri Matlab, 43
prod, 157
Puterea unui test, 291

quit, 2

rand, 118
randn, 118
random, 117
randperm, 119
randtool, 147
range, 155
rank, 10
Raport de corelație, 171
rayl, 80
raylfir, 282
rdivide, 8
real, 9
realmax, 6
realmin, 6
refcurve, 182
refline, 181
Regiune critică, 286
Regula lui Sturges, 120
rem, 9
Repartiție, 64
reshape, 7
Risc
 de speța întâi, 291
 de speța a doua, 291
roots, 183
round, 9

save, 11, 43
scatter, 142
scatter3, 143
sec, 9
sech, 9
Selecție, 197
 cu probabilități egale, 198
 sistematică, 197
 stratificată, 198
 tipică, 198
shading, 53
sign, 9
sin, 9
sinh, 9
Sistemul ecuațiilor
 de verosimilitate maximă, 253
size, 10
skewness, 155
Sondaj, 197
sort, 156
sortrows, 156
Spațiul probelor, 62
sqrt, 9
stairs, 28, 65, 134
Statistică, 198
 completă, 249
 suficientă, 229
std, 155
Subfuncții Matlab, 45
subplot, 32
sum, 157
surf, 52

t, 81
Tabel de contingență, 124, 349
Tabel de corelație, 124
Tabel statistic
 nesistematizat, 119
 sistematizat, 119, 120
tabulate, 125
tan, 9
tanh, 9
tblread, 128
tblwrite, 128
Teoremă limită
 centrală, 110
 Moivre–Laplace, 110
Test, 285

- cel mai puternic, 297
- nedeplasat, 297
- neparametric, 285
- parametric, 285
- Testul
 - F , 318, 321
 - bilateral, 318, 321
 - unilateral dreapta, 318, 321
 - unilateral stânga, 318, 321
 - T , 304, 325, 329
 - bilateral, 304, 310, 332
 - unilateral dreapta, 304, 332
 - unilateral stânga, 304, 332
 - Z , 286, 323
 - bilateral, 286
 - unilateral dreapta, 286
 - unilateral stânga, 286
 - χ^2 , 312, 334, 346, 349
 - bilateral, 312
 - neparametric, 338
 - parametric, 340, 342
 - unilateral dreapta, 312
 - unilateral stânga, 312
 - Fisher–Snedecor, 318
 - bilateral, 318
 - unilateral dreapta, 318
 - unilateral stânga, 318
 - Kolmogorov, 353
 - bilateral, 353
 - unilateral dreapta, 353
 - unilateral stânga, 353
 - Kolmogorov–Smirnov, 363
 - bilateral, 363
 - unilateral dreapta, 364
 - unilateral stânga, 364
 - raportului verosimilităților, 308
 - Student, 304
 - bilateral, 304, 332
 - unilateral dreapta, 304, 332
 - unilateral stânga, 304, 332
 - text, 24
 - tic, 49
 - times, 8
 - title, 24
 - toc, 49
 - trace, 10
 - transpose, 10
 - trimmean, 149
 - ttest, 307
 - ttest2, 327
 - unid, 68
 - unif, 74
 - unifit, 282
 - Valoare medie, 94, 95
 - condiționată, 99, 170
 - Valoarea funcției de selecție, 198
 - var, 155
 - Variație
 - intercuartilică, 152
 - Variabilă
 - nominală, 114
 - ordinală, 114
 - Variabilă aleatoare, 62
 - (absolut) continuă, 73
 - de tip continuu, 63
 - de tip discret, 63
 - simplă, 63
 - Variabile aleatoare
 - independente, 64
 - Variabile de selecție, 198
 - Varianță, 95
 - condiționată, 101
 - Vector aleator, 63
 - (absolut) continuu, 74
 - version, 4
 - Volumul
 - colectivității, 114

weib, 80
weibfit, 282
weiblike, 280
what, 17
which, 17
who, 11
whos, 11
workspace, 43

xlabel, 24
xor, 14

ylabel, 24

zeros, 6
zoom, 25
zscore, 290
ztest, 290

