

Notiuni introductive

Ce este analiza exploratorie a datelor?

Analiza exploratorie a datelor (*Exploratory Data Analysis* -EDA) este partea Statisticii care se ocupă cu trecerea în revistă, comunicarea și utilizarea datelor în cazul unui nivel scăzut de informație asupra lor.

Reamintim:

Pentru a studia o anumită caracteristică a unei populații, nu vom analiza toată populația, ci vom considera un anumit eșantion din ea.

Populația statistică este mulțimea elementelor, de aceeași natură, care au unele însușiri esențiale comune, însușiri caracteristice mulțimii privit ca un tot unitar.

Eșantionul (sample) este o submulțime a populației statistice considerate.

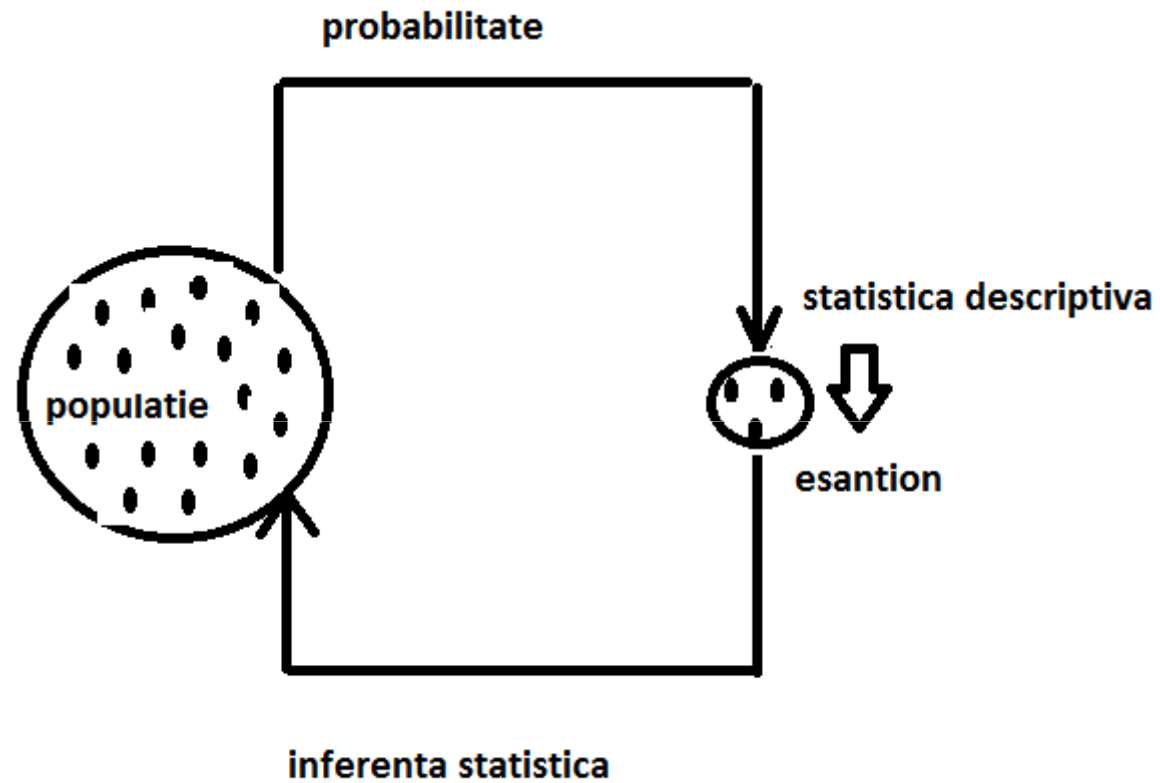
Reamintim:

Inferența statistică este totalitatea metodelor ce permit extragerea de concluzii demne de încredere pornind de la un eșantion statistic.

Prin metodele inferenței statistice deducem caracteristici necunoscute ale unei populații, pornind de la un eșantion obținut din acea populație. Caracteristicile eșantionului reflectă cu o anumită marjă de eroare caracteristicile întregii populații.

De reținut că inferența se aplică întregii populații, ca un tot unitar, și nu fiecărui element al populației respective.

Reamintim:



Spre deosebire de cazul clasic al testării ipotezelor, utilizat în Statistică pentru a verifica anumite supoziții apriorice

de exemplu: anumite corelații între diferite mărimi/variabile despre care există informații că ar fi cumva dependente

în cazul EDA se utilizează diferite tehnici pentru a identifica relații sistematice între anumite mărimi/variabile despre care nu există nicio informație prealabilă.

John Tukey, 1977

EDA a fost creată și numită astfel de statisticianul american John Tukey (Tukey J., *Exploratory Data Analysis*, Addison-Wesley, 1977).

Tehnicile computaționale EDA includ atât metode statistice elementare cât și altele avansate – *tehnici exploratorii multivariate* – create pentru a identifica anumite pattern-uri ascunse în mulțimi multivariate de date.

Tehnicile EDA sunt utilizate cu scopul de a:

- maximiza cunoașterea intimă a datelor;
- dezvălui structura de bază;
- extrage variabilele importante;
- detecta valorile extreme/exceptionale și anomaliile;
- identifica ipotezele fundamentale pentru a fi apoi testate;
- dezvolta modele suficient de simple;
- determina setarea optimă a parametrilor;
- sugera unele ipoteze privind cauzele fenomenelor observate;
- sugera tehnici statistice potrivite datelor disponibile;
- furniza cunoștințe pentru colectarea ulterioară de date pentru cercetare sau experimentare.

Odată utilizate tehnicile EDA se verifică rezultatele obținute.

Explorarea datelor este doar un prim stadiu de analiză a datelor și rezultatele obținute vor fi considerate cu titlu experimental atâta timp cât nu sunt validate alternativ.

De exemplu, rezultatul aplicării EDA sugerează un anumit model, atunci acesta trebuie validat aplicându-l la alt set de date, testându-i astfel calitatea predictivă.

In concluzie:

Analiza exploratorie a datelor- EDA poate fi considerată –
în principiu – o filosofie despre modul în care să fie „disecate”,
să fie „privite” și, în sfârșit, să fie interpretate datele.

Tipuri de date

Obiectele cu care lucrează statistica sunt reprezentate de *date*, adică acele caracteristici numerice sau nenumerică care descriu obiectele/subiecții unui studiu statistic

Datele obținute prin măsurare pot fi clasificate în funcție de tipul de informație conținut:

- *Datele numerice (cantitative)*
- *Datele categoriale (calitative)* sunt acele date care împart obiectele în diferite categorii.

Datele numerice sunt de două tipuri:
date ***discrete*** și date ***continue***.

Datele *discrete* apar atunci când este vorba de observații numerice întregi, privitoare la un anumit proces de numărare

Exemple: numărul de copii ai unei familii, pulsul, numărul de consultații pe an la care a fost supus un pacient, codul numeric, codul PIN.

Datele numerice *continue* se obțin de obicei în urma unor măsurători. Aceste date sunt de regulă exprimate prin numere reale, spre deosebire de cele discrete care sunt restricționate la numerele întregi.

Exemple: valoarea contului din bancă sau valoarea acțiunilor tranzacționate la Bursă; temperatura, presiunea atmosferică, viteza vântului; înălțimea, greutatea, tensiunea arterială, colesterolul unei anumite persoane, etc.

Din date continue se pot obține date discrete:

Exemplu: evaluarea venitului lunar

venit lunar $<$ 1000 lei,

1000lei $<$ venit lunar $<$ 2000 lei,

2000 lei $<$ venit lunar

Datele categoriale sunt de două tipuri: date ***nominale*** și date ***ordinale***

- Datele ***nominale*** sunt acele datele ce reprezintă mai multe categorii.

Exemple: grupa sanguină (A/B/AB/O), culoarea ochilor, specia de Iris (Iris Setosa, Virginica și Versicolour).

Datele nominale pot fi:

- binare (2 categorii)

Exemple: 0 sau 1; da / nu; adevărat / fals.

- enumerative (date discrete pentru care nu este definită o ordine).

Exemple: categoriile socio-profesionale sau culoarea ochilor).

Datele nominale pot fi și numerice
Exemplu: codurile poștale.

- Datele *ordinale* sunt date enumerative ordonate

Exemple: gradul fumatului (nefumător, fost fumător, fumător ,amator', fumător ,înraît'); ierarhizarea durerii (mică, medie, mare); răspunsurile la o anchetă de opinie (foarte mulțumit, mulțumit, nemulțumit, foarte nemulțumit) etc.

Observații

Datele numerice discrete sunt câteodată tratate ca date categoriale

Exemplu: numărul de copii născuți de o femeie, 0, 1, 2, 3, 4, împart mamele în categoriile corespunzătoare numărului de copii.

Invers, nu este corect să interpretăm datele categoriale ordonate ca date numerice:

Exemplu: la stadiile unor anumite boli, stadiul IV nu este de două ori mai grav decât stadiul II.

Este important ca în această situație să se ignore noțiunile de ordine sau de parametri numerici ca, de exemplu, media.

Alte tipuri de date

- **Rangul** reprezintă locul pe care îl ocupă un subiect într-o ierarhie
Exemple: competiție sportivă, examinare, preferința clienților pentru o anumită marcă.
- **Procentajul** descrie o anumită proporție între două cantități.
Exemple: procentajul celor cu studii superioare dintr-o populație, procentajul minorilor dintr-o populație.
- **Ratele și rapoartele** se referă la frecvența observată a unui anumit fenomen sau rapoartele dintre două mărimi, altele decât procentajele.
Exemple: mortalitatea la mia de locuitori, rata de apariție a unei boli pe arii geografice

- *Scorul* este folosit atunci când nu este posibilă o măsurătoare directă, dar trebuie cuantificată o anumită mărime.
Exemple: Scorul Apgar În neonatologie, gravitatea unei boli cuantificată ca ușoară, moderată, severă, foarte severă.

De cele mai multe ori datele sunt stocate sub forma unei matrice.
Prin convenție valorile variabilelor (atributelor/caracteristicilor) sunt stocate pe coloană în timp ce observațiile (obiectele) sunt stocate pe linii.

Considerând n vectori de dimensiune p , notația x_i^k utilizată de obicei se referă la a i -a variabilă (atribut/caracteristică observată) a vectorului \mathbf{x}_k (observația/ obiectul numărul k).

Se pot scrie într-un tabel cele n observații și attributele lor, astfel:

observația	atributul 1	..	atributul i	..	atributul p
observația 1	x_1^1		x_i^1		x_p^1
.....					
observația k	x_1^k		x_i^k		x_p^k
.....					
observația n	x_1^n		x_i^n		x_p^n

Reprezentarea datelor sub forma unei matrice \mathbf{X} , cu n linii și p coloane:

$$\mathbf{X} = \begin{pmatrix} x_1^1 & \dots & x_i^1 & \dots & x_p^1 \\ \dots & \dots & \dots & \dots & \dots \\ x_1^k & \dots & x_i^k & \dots & x_p^k \\ \dots & \dots & \dots & \dots & \dots \\ x_1^n & \dots & x_i^n & \dots & x_p^n \end{pmatrix}$$

1. Exemplu

Bază de date Iris, datorată lui R.A. Fisher, 1936, constă în 50 de exemplare din fiecare specie de iris (*Iris Setosa*, *Virginica* și *Versicolor*), fiecare exemplar fiind caracterizat de lungimea și lățimea sepalelor, respectiv a petalelor.

Prezentăm un tabel cu 7 flori, atributele lor (date *numerice continue*) și specia de iris /clasa căreia îi aparține fiecare floare (dată *nominală*). Menționăm că pentru atribute avem notațiile: PW lățimea petalei, PL lungimea petalei, SW lățimea sepalei, SL lungimea sepalei.

Specie de iris	PW	PL	SW	SL
Setosa	2	14	33	50
Virginica	24	56	31	67
Virginica	23	51	31	69
Setosa	2	10	36	46
Virginica	20	52	30	65
Versicolor	13	45	38	57
Versicolor	16	47	33	63

Prezentăm datele sub formă de matrice, codând datele nominale în felul următor:

0-Setosa, 1-Virginica, 2-versicolor.

$$\begin{pmatrix} 0 & 2 & 14 & 33 & 50 \\ 1 & 24 & 56 & 31 & 67 \\ 1 & 23 & 51 & 31 & 69 \\ 0 & 2 & 10 & 36 & 46 \\ 1 & 20 & 52 & 30 & 65 \\ 2 & 13 & 45 & 38 & 57 \\ 2 & 16 & 47 & 33 & 63 \end{pmatrix}$$

Pentru a accesa baza de date Iris a lui Fisher, in Matlab vom scrie:

```
>>load fisheriris
```

încărcând variabila **meas** (measures) o matrice numerică cu 150 linii și 4 coloane reprezentând attributele considerate (lungimea sepalei, lățimea sepalei, lungimea petalei, lățimea petalei) și variabila **species**.

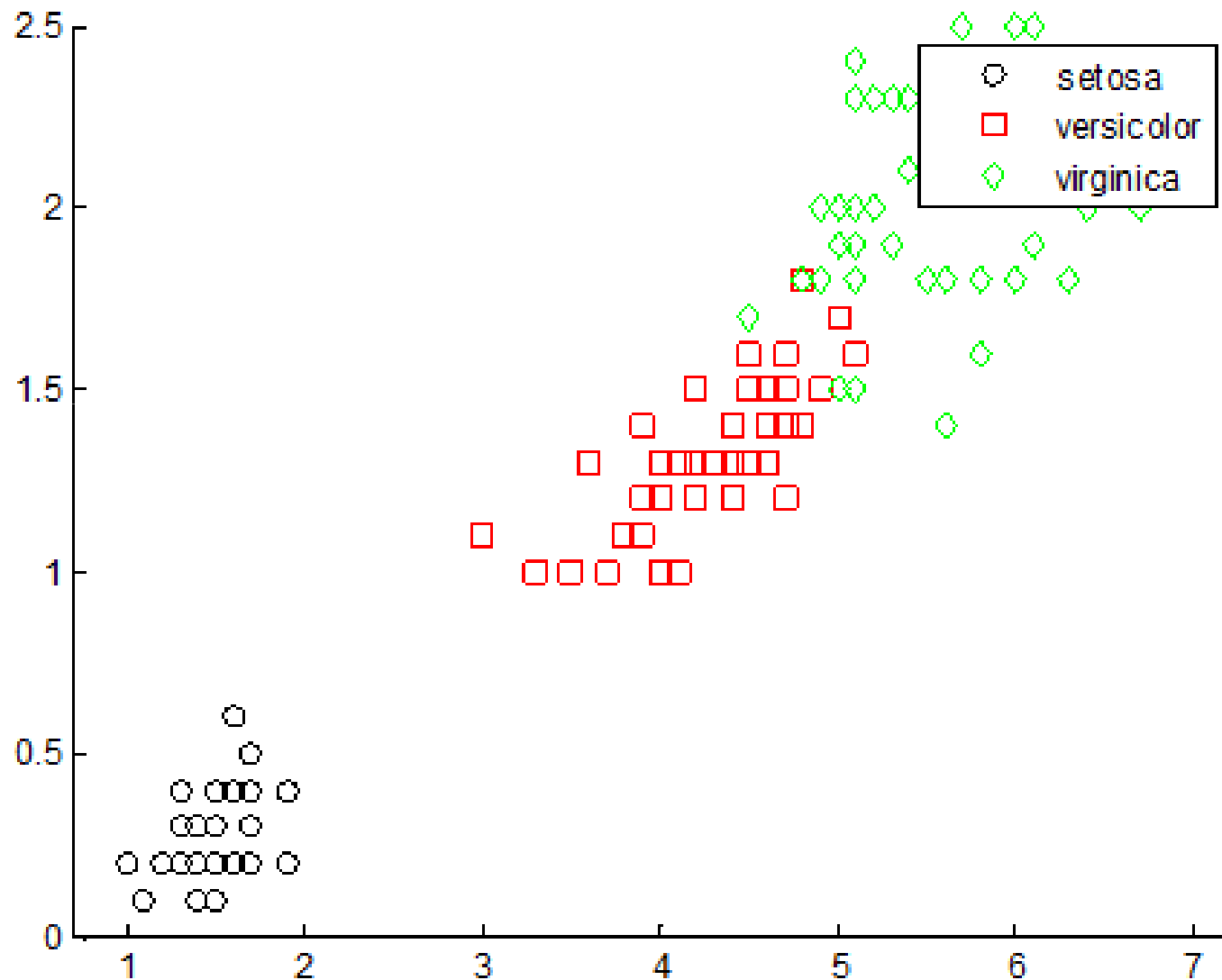
Vom prezenta o vizualizare a florilor de Iris, din baza de date a lui Fischer, luând în considerare dimensiunile petalelor.

Notând cu x_1 = lungimea în cm a petalei și x_2 = lățimea în cm a petalei, construim vectorul $\mathbf{x} = (x_1, x_2)$, corespunzător unui iris. Exemplarele de flori vor fi reprezentate în spațiul bidimensional al atributelor. Să reținem că lungimea, respectiv lățimea unei petale diferă și în cazul aceluiași tip de Iris.

Vom specifica atributele folosite `meas(:,3)`, `meas(:,4)`. Impunem ca toate florile, în funcție de specia căreia îi aparțin, să fie reprezentate prin cerc (setosa), pătrat (versicolor) și romb (virginica)


```
>> load fisheriris
```

```
>> gscatter(meas(:,3), meas(:,4), species, 'krg', 'osd')
```



2. Exemplu

O baza de date medicală conține datele unor pacienți, atributele acestora fiind: numărul de consultații pe an la care a fost supus pacientul - nr. cons (dată *numerică discretă*), sex-S (dată *nominală*), vârsta - V, indicele masei ponderale - IMC. Glicemia (glyc), colesterolul (Ch), trigliceride(trigl) (date *numerice continue*), gradul fumatului- gr F (dată *ordinală*) și cele două categorii cărora le aparțin: pacient hipertensiv- Ω_1 sau pacient sănătos - Ω_2 (dată *nominală*).

Prezentăm 3 linii din această bază, corespunzătoare a 3 pacienți.

nr cons	S	V	IMC	glyc	Ch	trigl.	gr F	clasa
9	F	60	30	92	286	349	fost	Ω_1
2	F	40	23,2	79	180	145	nefumător	Ω_2
5	M	55	25.1	128	230	210	amator	Ω_1

Vom face următoarele codificări:

- pentru feminin/F-0, pentru masculin/M-1;
 - gradul fumatului: nefumător - 0, fost fumător -1, fumător amator -2, fumător înrăit - 3;
 - pacient hipertensiv Ω_1 -1, pacient sănătos Ω_2 - 0
- obținând astfel următoarea matrice atașată datelor avute:

$$\begin{pmatrix} 0 & 9 & 60 & 30 & 92 & 286 & 349 & 1 & 1 \\ 0 & 2 & 40 & 23.2 & 79 & 180 & 145 & 0 & 0 \\ 1 & 5 & 55 & 25.1 & 128 & 230 & 210 & 2 & 1 \end{pmatrix}$$

Date cenzurate

Există situații în care o anumită dată există, dar din diferite motive nu poate fi bine precizată, situații în care considerăm că datele sunt *cenzurate*.

De exemplu în *Analiza supraviețuirii*, tehnică statistică clasică, care studiază dinamica timpului de supraviețuire după o anumită operație sau tratament, o parte din subiecții incluși în lotul de studiu decedează în perioada de observație, dar o altă parte supraviețuiesc acestei perioade sau se retrag benevol.

Alt exemplu când se efectuează anumite măsurători și aparatul de măsură nu poate înregistra valori înafara scalei sale.

Procesarea datelor în cursul unei analize statistice impune ca absolut necesară existența așa numitei *variabilități* a datelor.

Prin *variabilitate* se înțelege orice modificare care are loc într-o mulțime de date, indiferent de tipul acestora.

Nu se poate face analiză statistică pe variabile care sunt constante

Cu cât variabilitatea datelor este mai mare cu atât analiza statistică va da rezultate mai consistente.

Definirea datelor in context probabilist

Avem la dispoziție o anumita mulțime de obiecte/subiecți (o așa numită *populație statistică*) și suntem interesați de analiza principalelor lor caracteristici/atribute, care reprezintă date statistice.

Din punct de vedere matematic *data* este o aplicație definită pe mulțimea ce reprezintă populația și cu valori într-o anumită mulțime, ce depinde de data respectivă

Considerând un câmp de probabilitate (Ω, Σ, P) , unde Ω este chiar populația considerată, Σ este o σ -algebră de părți a lui Ω , data X a populației statistice Ω , dacă este numerică, este o variabilă aleatoare pe câmpul de probabilitate (Ω, Σ, P) , cunoscută și sub numele de variabilă statistică.,

În cazul în care X nu ia valori numerice, pe baza unor echivalări numerice ale acestor valori (codificări), putem considera X ca fiind o variabilă aleatoare.

3. Exemplu

Considerând ca populație Ω o clasă de elevi, considerăm o dată X , referitoare la această populație, reprezentată prin înălțimea elevilor.

Pentru elevul A ce are înălțimea de 1,70, avem $X(A) = 1.70$

Statistica descrittiva

Descrierea statistică constă într-o mulțime de diverse metode ce au ca scop rezumarea unui număr mare de observații privind datele, punând astfel în evidență principalele lor caracteristici.

Există două mari abordări a descrierii statistice a datelor:

- Determinarea unor parametri numerici, interesul fiind focalizat pe proprietățile lor matematice;
- Diverse reprezentări grafice simple ale datelor, a căror interpretare nu este dificilă, fiind de cele mai multe ori foarte sugestivă. Din punct de vedere strict informațional interpretare este totuși limitată.

Pentru a rezuma un mare număr de observații privind un set de obiecte/instanțe, punând în evidență caracteristicile principale ale atributelor lor, vom utiliza reprezentări numerice care cuprind principalele caracteristici statistice ale datelor respective. (media, mediana, deviația standard, modul etc.).

Este vorba de reprezentarea variabilității unor date. Această variabilitate poate fi una cu cauze cunoscute, o variabilitate ‚deterministă’ care este descrisă statistic pentru a o pune și mai bine în evidență și a o cuantifica precis, sau poate fi o variabilitate cu cauze doar bănuite sau chiar necunoscute – variabilitatea ‚aleatorie’ – și care, folosind statistica, se speră a fi clarificată cauzal.

Pentru a procesa datele în cursul unei analize statistice este absolut necesară existența așa numitei variabilități a datelor

Un obiect/o instanță corespunde din punct de vedere probabilist unei variabile aleatoare multidimensionale adică un vector aleator n -dimensional, subiect al analizei statistice multivariate.

Fiecare componentă a instanței, reprezentând un anumit atribut, este privită la rândul său ca o variabilă aleatoare.

Luăm în considerație în acest caz parametrii numerici (statistici) ce caracterizează o variabilă aleatoare, parametri deosebit de utili în descrierea dinamicii acesteia.

În cele ce urmează, vom considera că fiecare atribut x al unei instanțe reprezintă o valoare pe care o poate lua variabila aleatoare generatoare X , corespunzătoare aceluși atribut.

Reamintim că, pentru variabila aleatoare X , funcția $F_X: \mathbf{R} \rightarrow [0, 1]$, definită de:

$$F_X(x) = P\{X < x\} = P\{\omega \in \Omega; X(\omega) \in (-\infty, x)\},$$

se numește *funcția de repartiție (repartiția de probabilitate)*, unde (Ω, Σ, P) este un spațiu (câmp de probabilitate).

Funcția de repartiție a variabilei statistice X

În cazul unei analize statistice, unei variabile aleatoare X îi corespunde variabila statistică corespunzătoare, notată în mod firesc tot X . La fel ca în cazul unei variabile aleatoare, și în cazul unei variabile statistice putem defini noțiunea de funcție de repartiție.

Astfel, prin *funcția de repartiție* sau *funcția cumulativă (frecvența cumulată)* a variabilei statistice X (corespunzătoare variabilei aleatoare X), definită de seria statistică asociată (eșantionul corespondent) $\{x_i\}_{i=1, \dots, n}$, înțelegem aplicația $F: \mathbf{R} \rightarrow [0, 1]$, dată de:

$$F(x) = \frac{f_x}{n}, x \in \mathbf{R},$$

unde f_x reprezintă numărul observațiilor x_i strict mai mici decât x .

Cuantila

Definim *cuantila* (*quantile* - Kendall, 1940) de ordin α , $0 < \alpha < 1$, a variabilei statistice X , ca fiind numărul q_α cu proprietatea $F(q_\alpha) = \alpha$, adică $f_{q_\alpha} = \alpha \cdot n$, unde f_{q_α} reprezintă numărul observațiilor strict mai mici decât q_α .

Din punctul nostru de vedere, al explorării datelor, cuantilele implică practic divizarea în $1/\alpha$ submulțimi de mărimi egale a unei mulțimi ordonate de date, ele reprezentând de fapt tocmai frontierele (pragurile) dintre submulțimile consecutive.

În acest context, o *percentilă* (*percentile* -Galton, 1885) reprezintă oricare dintre cele 99 de valori care împart o mulțime ordonată de date în 100 de submulțimi de mărimi egale, consecutive.

De exemplu a 50-a percentilă împarte setul de date în 50% date dedesubtul său și 50% date deasupra sa.

Astfel:

- *decilele* (*deciles*) reprezintă cele 9 valori care împart o mulțime ordonată de date în 10 submulțimi consecutive de mărimi egale;
- *cuartilele* (*quartiles* - Galton, 1882) reprezintă cele trei valori notate Q_1 , Q_2 , Q_3 , care împart o mulțime ordonată de date în 4 submulțimi consecutive de mărimi egale. Cele mai utilizate cuantile sunt cuartilele Q_1 , Q_2 , și Q_3 . Astfel, prima cuartilă Q_1 (cuartila inferioară – a 25-a percentilă) are sub ea 25% din date și deasupra 75% din date; a doua cuartilă Q_2 are sub ea 50% din date și deasupra tot 50% din date; a treia cuartilă Q_3 (cuartila superioară – a 75-a percentilă) are sub ea 75% din date și deasupra 25% din date.

Calculul cuantilelor in Matlab- *Statistics Toolbox*

```
>> Y = quantile(X,p)
```

unde

Y returnează cuantilele corespunzătoare valorilor lui X,

p este scalarul sau vectorul valorilor probabilității și anume:

dacă dorim să calculăm

- percentilele, luăm $p=0.01:0.01:0.99$;
- decilele luăm $p=0.1:0.1:0.9$;
- quartilele, luăm $p=0.25:0.25:0.75$.

Dacă X este un vector, Y va avea dimensiunea lui p și componenta Y(i) va conține cuantila p(i).

Dacă X este o matrice, atunci linia a i-a conține cuantilele p(i) ale fiecărei coloane a lui X

4. Exemplu: decilele si quartilele corespunzatoare dimensiunilor sepalelor si petalelor florilor de iris

```
>> load fisheriris  
>> X=[meas(:,1) meas(:,2) meas(:,3) meas(:,4)];  
>> y10=quantile(X,0.1:0.1:0.9)
```

```
y10 =
```

4.8000	2.5000	1.4000	0.2000
5.0000	2.7000	1.5000	0.2000
5.2500	2.8000	1.7000	0.4000
5.6000	3.0000	3.9000	1.1500
5.8000	3.0000	4.3500	1.3000
6.1000	3.1000	4.6500	1.5000
6.3000	3.2000	5.0000	1.8000
6.5500	3.4000	5.3500	1.9000
6.9000	3.6500	5.8000	2.2000

```
>> y1=quantile(X,0.25:0.25:0.75)
```

```
y1 =
```

5.1000	2.8000	1.6000	0.3000
5.8000	3.0000	4.3500	1.3000
6.4000	3.3000	5.1000	1.8000

Valorile tipice corespunzătoare unei analize a datelor

- măsuri tipice ale tendinței centrale (localității): *modul (mode)*, *mediana (median)*, *media (mean, arithmetic mean, average)*, *media geometrică (geometric mean)* și *media armonică (harmonic mean)*;
- măsuri tipice ale împrăstierii (deviației): *dispersia (variance)* și *deviația standard/abaterea medie pătratică (standard deviation)*;
- măsuri tipice ale formei repartiției: *asimetria (skewness)* și *excesul (kurtosis)*.

Media

Media este cel mai comun parametru ce măsoară „tendința centrală” a unei serii statistice, reprezentând practic media aritmetică a tuturor observațiilor:

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

In Matlab, pentru calcularea mediei, utilizăm funcția

```
>>M=mean(X)
```

unde vectorul X secvența de date a cărei medie vrem să o calculăm.

Dacă X este o matrice `mean(X)` consideră coloanele lui X ca fiind vectori, returnând un vector linie ale cărui componente sunt mediile vectorilor coloană

5. Exemplu

Considerăm următoarea secvență de date reprezintă valorile înălțimii băieților dintr-o clasă de gimnaziu:

{ 1,40; 1,37; 1,57; 1,46; 1,49; 1,46; 1,39; 1,55; 1,29; 1,58; 1,50; 1,33 }

Să calculăm înălțimea medie a băieților din această clasă:

```
>> X1=[1.40, 1.37, 1.57, 1.46, 1.49, 1.46, 1.39, 1.55, 1.29, 1.58, 1.50, 1.33 ];
```

```
>>m1= mean(X1)
```

```
m1 =
```

```
1.4492
```

6. Exemplu

Vom da un exemplu de calcul al mediei pentru date discrete, considerând un tabel de distribuire (repartiție) a datelor, în care valorile x_i au frecvențele de apariție n_i :

Un medic de familie are în îngrijire 140 de copii. Prezentăm situația numărului de vizite ale copiilor la cabinet pe durata unui an:

Nr vizite	0	1	2	3	4	5	6	7
Nr copii	2	38	27	45	18	5	4	1

Să calculăm media numărului de vizite efectuate de un copil, într-un an.

```
>> a=[0 0];
>> b=[1];for i=2:38 b=[1 b];end
>> c=[2];for i=2:27 c=[2 c];end
>> d=[3];for i=2:45 d=[3 d];end
>> e=[4];for i=2:18 e=[4 e];end
>> f=[5];for i=2:5 f=[5 f];end
>> g=[6];for i=2:4 g=[6 g];end
>> h=[7];
>> X2=[a b c d e f g h]
>> m=mean(X2)
m =
    2.5357
```

Cu toate că media este o caracteristică foarte sugestivă a datelor pe care le reprezintă, existența valorilor extreme (*outliers*) îi pot perturba serios capacitatea de ilustrare a datelor.

Considerăm secvență de date din exemplul nr 5, la care adaugăm înălțimea unui alt băiat care suferă de o ușoară forma de nanism:

```
>> X=[1.40, 1.37, 1.57, 1.46, 1.49, 1.46, 1.39, 1.55, 1.29, 1.58, 1.50, 1.33 1.05];  
>>m= mean(X)  
m =  
    1.4185
```

Remarcăm faptul că media $m=1.4185$ diferă de cea calculată în exemplul nr 5, $m_1= 1.4492$, valori fiind suficient de distincte datorită influenței unei singure date.

Mediana

Pentru a evita asemenea situații, în astfel de cazuri se utilizează mediana în locul mediei. Astfel, *mediana* este definită ca numărul real care împarte în două efective egale seria statistică dată, observațiile fiind ordonate crescător, cu alte cuvinte mediana este chiar a doua cuartilă Q_2 . Formal, mediana este dată de:

$$P\{X \leq Q_2\} = P\{X > Q_2\} = 1/2 .$$

În Matlab, pentru calcularea medianei, utilizăm funcția

```
>>M=median(X)
```

unde vectorul X este secvența de date a cărei mediană vrem să o calculăm.

Dacă X este o matrice `median(X)` consideră coloanele lui X ca fiind vectori, returnând un vector linie ale cărui componente sunt valorile medianelor.

Folosind cuantilele avem:

```
>>M =quantile(X,0.50)
```


Exemple

In exemplul nr 5 avem:

$\text{median}(X1) = 1.4600$ și $\text{mean}(X1) = 1.4492$

$\text{median}(X) = 1.4600$ și $\text{mean}(X) = 1.4185$

In exemplul nr 6 avem:

$\text{median}(X2)=3$ și $m=\text{mean}(X2)=2.5357$

7. Exemplu

O fabrică produce bare de oțel de diametru 20mm. Datorită unor reclamații asupra calității produselor s-a ales un eșantion măsurându-se diametrul fiecărei bare din acest eșantion; rezultatele sunt prezentate în următoarea secvență de date:

```
{ 19.9 19.8 20.1 19.9 19.7 20.1 20 19.6 19.7 20.1 20.4 20 19.9 19.8 20.2 20  
19.8 19.6 29.9 20.3}
```

Vom calcula media, respectiv mediana valorilor diametrului barelor din acest eșantion

```
>> X3=[19.9 19.8 20.1 19.9 19.7 20.1 20 19.6 19.7 20.1 20.4 20 19.9 19.8 20.2  
20 19.8 19.6 29.9 20.3]
```

```
>> m=mean(X3)
```

```
m =  
    20.4400
```

```
>> M=median(X3)
```

```
M =  
    19.9500
```

Dacă efectivul seriei statistice este un număr impar $n = 2k + 1$, atunci mediana este a $(k + 1)$ -a valoare a seriei, iar dacă efectivul este un număr par $n = 2k$, atunci mediana se înlocuiește cu *intervalul median* dat de valorile a k -a și a $(k + 1)$ -a (mediana se poate considera astfel ca mijlocul acestui interval –media aritmetică a capetelor sale).

Mediana mai este folosită atunci când există posibilitatea ca unele valori extreme ale seriei statistice să fie cenzurate. Atunci când există observații care nu sunt suficient de exact precizate, nu putem folosi media, înlocuind-o prin mediană dacă **avem valori exacte pentru mai mult de jumătate din observații** (cazul măsurătorilor fizico-chimice, când există valori în afara scalei normale a aparatului de măsură).

Ambele măsuri sunt la fel de eficiente și, cu toate că media este mai frecvent folosită decât mediana, aceasta din urmă poate fi mai valoroasă în anumite circumstanțe.

În exemplul nr 5, al setului de date privind înălțimea unor elevi, să calculăm medianele corespunzătoare celor două cazuri (cu și fără valoarea extremă). sunt egale, ceea ce dovedește că mediana nu este influențată semnificativ de valori extreme.

Modul (modă)

Modul (modă) este o altă măsură a tendinței centrale care reprezintă cea mai frecventă valoare a secvenței, des utilizată în cazul datelor categoriale. Este foarte probabil ca modul să nu fie unic, deoarece pot apărea mai multe valori ale seriei statistice ce au aceeași frecvență de apariție. Este vorba de date *plurimodale*

Dacă datele sunt grupate în clase adică valorile aparțin unor intervale, vom numi *clasă modală* orice clasă corespunzând unui maxim al frecvenței.

În Matlab, pentru calcularea modului, utilizăm funcția

```
>> M=mode(X)
```

unde vectorul X este secvența de date al cărei mod vrem să-l calculăm (vrem să știm care date sunt cele mai frecvente). Dacă X este o matrice `mode(X)` consideră coloanele lui X ca fiind vectori, returnând un vector linie ale cărui componente sunt valorile modelor.

Exemple

În exemplul nr 5 avem:

```
>> mode(X1)
```

```
ans =
```

```
1.4600
```

```
>> mode(X)
```

```
ans =
```

```
1.4600
```

În exemplul nr 7 avem:

```
>> mode(X3)
```

```
ans =
```

```
19.8000
```

Media geometrica

Media geometrică este dată de formula: $\sqrt[n]{x_1 \cdot x_2 \cdot \dots \cdot x_n}$

Media geometrică este utilizată cu precădere în cazul măsurătorilor cu scală neliniară (e.g. în psihometrie, unde rata intensității unui stimul este adesea o funcție logaritmică în raport cu intensitatea, caz în care se folosește media geometrică și nu media aritmetică).

În Matlab, pentru calcularea mediei geometrice, utilizăm funcția

```
>>M=geomean(X)
```

unde vectorul X este secvența de date a cărei medie geometrică vrem să o calculăm.

Dacă X este matrice, **geomean**(X) este un vector linie, ale cărui componente sunt mediile geometrice ale coloanelor.

Media armonica

Media armonică, dată de formula:
$$\frac{n}{\sum_{i=1}^n \frac{1}{x_i}}$$

se utilizează câteodată în cazul determinării mediei frecvențelor.

În Matlab, pentru calcularea mediei armonice, utilizăm funcția

```
>>M=harmmean(X)
```

unde vectorul X este secvența de date a cărei medie armonică vrem să o calculăm.

Dacă X este matrice, `geomean(X)` este un vector linie, ale cărui componente sunt mediile armonice ale coloanelor.

Calculăm pentru secvențele de date din exemplul nr 5 (X și $X1$), respectiv pentru secvența din exemplul nr 7 ($X3$) valorile pentru medie, modă, mediană, medie geometrică și medie aritmetică.

Principalii parametri ai tendinței centrale (cu două zecimale)

Secvența de date	medie	mediana	modă	medie geometrică	media armonică
X	1.45	1.46	1.46	1.45	1.44
$X1$	1.42	1.46	1.46	1.41	1.40
$X3$	20.44	19.95	19.80	20.35	20.28

Din tabelul de mai sus se observă că în toate cele trei cazuri repartiția este suficient de „simetrică”, parametrii tendinței centrale fiind apropiați ca valoare.

Dispersia

O altă abordare privind analiza datelor este reprezentată de măsurarea împrăstierii (*dispersiei*) față de medie, adică **măsurarea distanței fiecărei valori a seriei statistice față de media eșantionului.**

Plecând de la cazul probabilist clasic al dispersiei, vom defini *dispersia* sau *varianța (variance)* - termen introdus de Fisher, 1918 - corespunzătoare unei serii statistice $\{x_i\}_{i=1, \dots, n}$, cu ajutorul formulei:

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

unde m este media cunoscută a variabilei statistice (a populației originare, cu alte cuvinte).

In Matlab, pentru calcularea dispersiei, utilizăm funcția

```
>>V=var(X)
```

unde vectorul X este secvența de date a cărei dispersie vrem să o calculăm.

Dacă X este matrice, $\text{var}(X)$ este un vector linie, ale cărui componente sunt dispersiile vectorilor coloană.

De obicei, considerăm că seria statistică cu care lucrăm nu reprezintă toată populația ci este doar un eșantion al ei și astfel media m nu este cunoscută ci putem calcula doar media eșantionului \bar{x} .

Vom folosi în locul formulei de mai sus o formulă de aproximație (o estimare) a dispersiei, înlocuind media m cu media seriei \bar{x} și împărțind prin $(n - 1)$ în loc de n , deci:

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

Remarcăm aici că, pentru serii statistice de dimensiuni mari, diferența dintre valoarea dată de formula de mai sus și formula clasică:

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

este neglijabilă.

Deviația standard

De multe ori, este preferabil ca în loc de dispersie să folosim o mărime care este măsurată cu aceeași unitate ca și seria statistică, și anume *deviația standard* (sau *abaterea medie pătratică*), dată de formula:

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

În Matlab, pentru calcularea deviației standard, utilizăm funcția

```
>> s=std(X)
```

unde vectorul X este secvența de date a cărei deviație standard vrem să o calculăm.

Dacă X este matrice, $\text{std}(X)$ este un vector linie, ale cărui componente sunt deviațiile standard a vectorilor coloană.

Interval de încredere

Deviația standard este folosită în statistica descriptivă mai ales pentru definirea unor intervale în care se găsesc marea majoritate a observațiilor. Astfel, în cazul unor repartiții rezonabil de simetrice, marea majoritate a observațiilor ce compun seria statistică (aproximativ 95% din ele) se găsesc în intervalul definit de:

medie $\pm 2 \times$ deviația standard,

numit *interval de încredere (confidence interval).*

Cele spuse mai sus nu mai au relevanță dacă nu avem repartiții relativ simetrice.

Exemple

Vom ilustra măsurile tipice împrăştierii și formei repartiției în cazul secvențelor de date din exemplul nr 5 (X și $X1$), respectiv pentru secvența din exemplul nr 7 ($X3$).

```
>> X=[1.40, 1.37, 1.57, 1.46, 1.49, 1.46, 1.39, 1.55, 1.29, 1.58, 1.50, 1.33 ];  
>> m=mean(X); v=var(X);s=std(X); [v, s, m-2*s, m+2*s]  
ans =  
    0.0089    0.0944    1.2603    1.6380
```

Principali parametri ai împrăştierii (dispersia, deviația standard, intervalul de încredere pentru medie):

Dispersia	Deviația standard	Intervalul de încredere
0.0089	0.0944	(1.2603, 1.6380)

```
>> X1=[1.4,1.37, 1.57, 1.46, 1.49, 1.46, 1.39, 1.55, 1.29, 1.58, 1.50, 1.33, 1.05 ];
>> m1=mean(X1); v1=var(X1);s1=std(X1); [v1, s1, m1-2*s1, m1+2*s1]
ans =
    0.0204    0.1429    1.1326    1.7043
```

Principalii parametri ai împrăştierii (dispersia, deviația standard, intervalul de încredere pentru medie):

Dispersia	Deviația standard	Intervalul de încredere
0.0204	0.1429	(1.1326 1.7043)


```

>> X3=[19.9 19.8 20.1 19.9 19.7 20.1 20 19.6 19.7 20.1 20.4 20 19.9 19.8 20.2
20 19.8 19.6 29.9 20.3];
>> m3=mean(X3); v3=var(X3);s3=std(X3); [v3, s3, m3-2*s3, m3+2*s3]
ans =
    5.0057    2.2373   15.9653   24.9147

```

Principalii parametri ai împrăstierii (dispersia, deviația standard, intervalul de încredere pentru medie):

⊕	Dispersia	Deviația standard	Intervalul de încredere
	5.0057	2.2373	(15.9653 24.9147)

□

Remarcă

În afara parametrilor statistici amintiți mai sus, se mai utilizează câteodată și:

- *Domeniul (range)* valorilor, reprezentat de diferența între maximum și minimum valorilor datelor;
- *Domeniul inter-cuartile (interquartiles range)*, definit de diferența $Q_3 - Q_1$;
- Media abaterilor absolute de la medie:

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

- Mediana abaterilor absolute de la medie:

$$MAD(x) = \text{mediana } \{|x_1 - \bar{x}|, |x_2 - \bar{x}|, \dots, |x_n - \bar{x}|\}$$

Reprezentarea grafică
a unei mulțimi de date

Reprezentarea grafică elementară a datelor înseamnă conversia datelor în format vizual sau tabular simplu, astfel încât să se poată analiza și raporta rapid caracteristicile datelor, precum și relațiile existente între atribute.

Reprezentările grafice aferente unui set de date (atribute ale unei instanțe în cazul de față) depind de natura atributelor: calitative sau cantitative.

Reprezentarea circulară (pie) a datelor calitative

Datele calitative pot fi reprezentate grafic cu ajutorul diferitelor diagrame formate din bastoane verticale sau orizontale, cercuri, etc., bi- sau tri-dimensionale, plecând de la partiția setului de date pe care atributul o induce.

În cazul utilizării *reprezentării circulare* (numită și '*pie*' ~ *plăcintă* în engleză), întreaga mulțime de obiecte este reprezentată de cerc, fiecare atribut al obiectelor este reprezentat printr-un sector circular a cărui suprafață este proporțională cu numărul obiectelor având atributul respectiv (sau cu procentajul corespunzător).

8. Exemplu

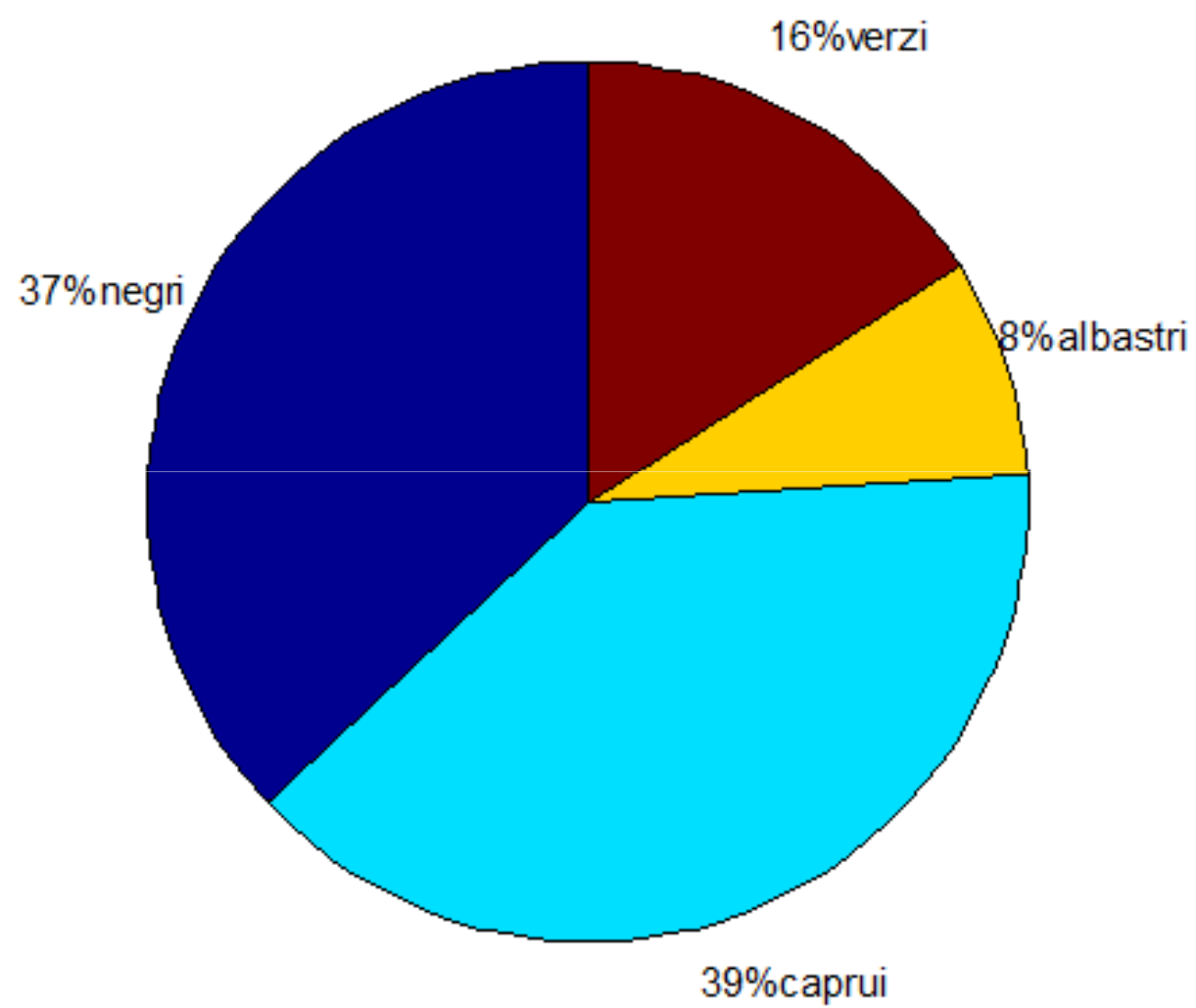
În Matlab `pie(X)` realizează o reprezentare circulară bidimensională, pe baza datelor din vectorul `X`, vector ce are drept componente procentajele corespunzătoare unor atribute. Fiecare atribut va fi reprezentat de un sector circular.

Suntem interesați de culoarea ochilor indivizilor - atribut calitativ, ce aparțin unei anumite populații. Presupunem că mulțimea culorilor ochilor este $Culoare_ochi = \{negru, albastru, verde, căprui\}$ și că din studiul efectuat a rezultat că 37% din populație are ochii negri, 39% are ochii căprui, 8% are ochii albaștri și 16% are ochii verzi, așadar $X=[0.37 \ 0.39 \ 0.08 \ 0.16]$.

```
>> X=[0.37 0.39 0.08 0.16];
```

```
>>pie(X,{'37%negri','39%caprui', ' 8%albaștri', ' 16%verzi'})
```

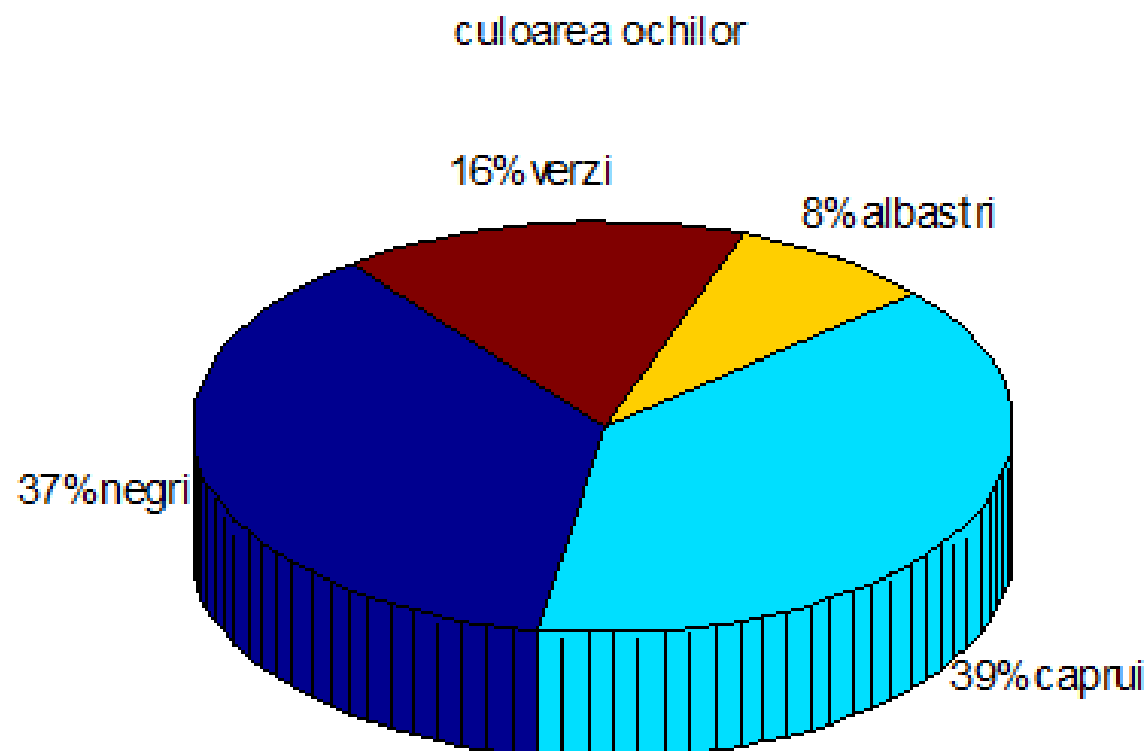
culoarea ochilor



În Matlab `pie3(X)` realizează o reprezentare circulară tridimensională, reprezentare ce nu furnizează informație suplimentară fața de cazul bidimensional;

```
>>X=[0.37 0.39 0.08 0.16];
```

```
>>pie3(X,{'37%negri','39%caprui', ' 8%albastri', ' 16%verzi'})
```



9. Exemplu

Vom folosi funcția `pie` pentru a vizualiza contribuția ce o aduc trei produse în totalul vânzării. Coloanele matricei `X` conțin vânzările anuale pentru fiecare produs de-a lungul a 5 ani:

```
>>X = [19.3 22.1 51.6;  
       34.2 70.3 82.4;  
       61.4 82.9 90.8;  
       50.5 54.9 59.1;  
       29.4 36.3 47.0];
```

Vom calcula vânzările totale pentru fiecare produs în cei 5 ani utilizând funcția `sum`.

```
>>x = sum(X);  
     194.8000 266.5000 330.9000
```

Folosind argumentul **explode**, putem pune în evidență sectorul circular care are cea mai mare contribuție.

explode este un vector ale cărui componente nenule corespund unui anumit sector circular.

Pentru început vom crea un vector de componente nule:

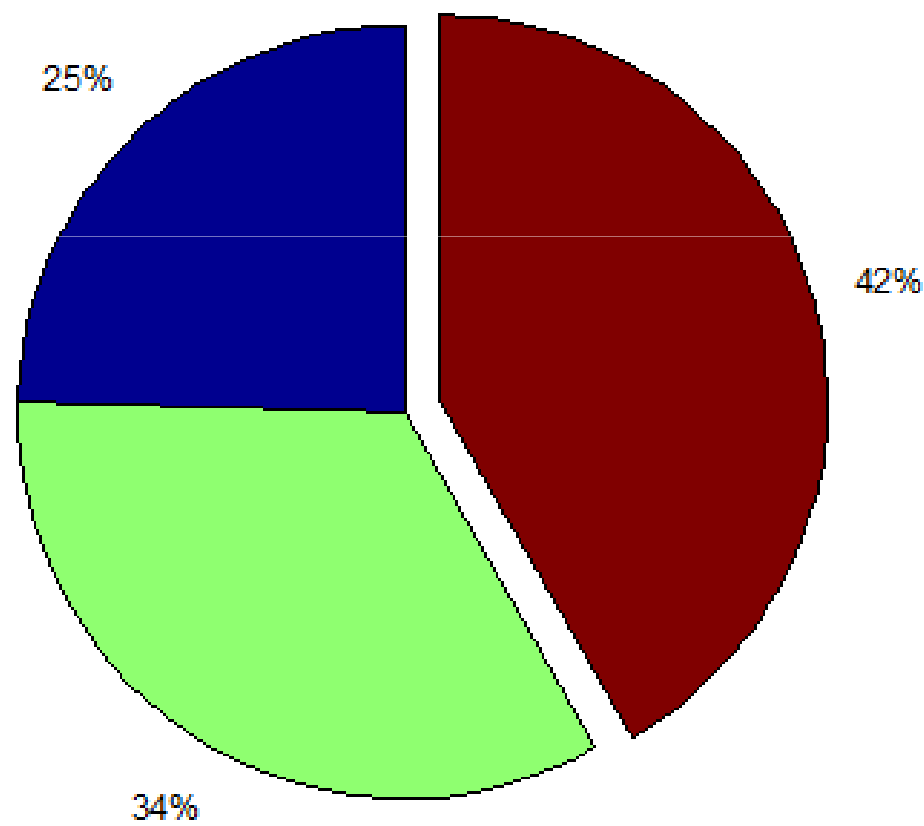
```
>> explode = zeros(size(x))  
explode =  
    0    0    0
```

Vom găsi produsul care își aduce cea mai mare contribuție și vom atribui componentei corespunzătoare a argumentului **explode** valoarea 1:

```
>> [c,offset] = max(x)  
c =  
 330.9000  
offset =  
    3  
>> explode(offset) = 1  
explode =  
    0    0    1
```

Pentru a crea reprezentarea circulară vom scrie:

```
>>h = pie(x,explode);
```



Reprezentarea prin diagrame cu bastoane (bar) a datelor calitative

Cealaltă modalitate de reprezentarea a unui atribut calitativ – cel al *diagramelor cu bastoane (bar graph)* – se raportează la o vizualizare într-un sistem de axe: pe axa absciselor apar attributele iar pe axa ordonatelor apare numărul obiectelor cu atributul respectiv (sau procentajul corespunzător).

În Matlab pentru a vizualiza bidimensional diagramele cu bastoane folosim funcția

```
>>bar(Y)
```

unde Y este vectorul ale cărui componente reprezintă numărul obiectelor cu atributul respectiv sau procentajul corespunzător
Funcția

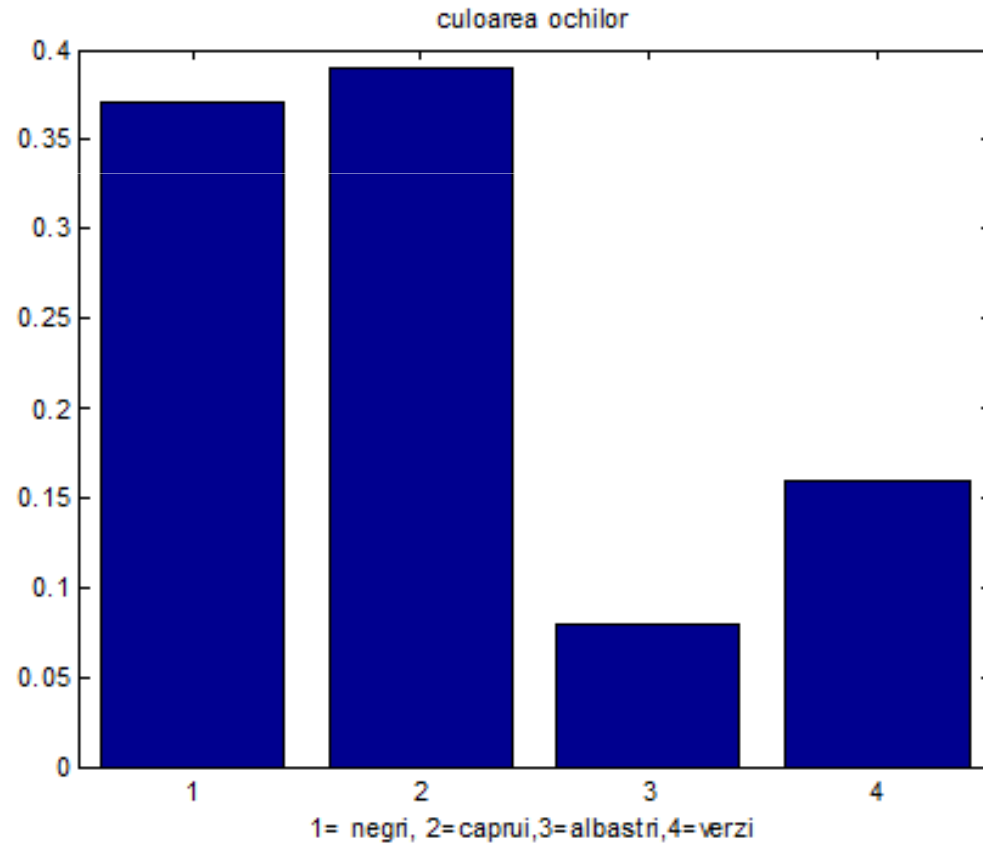
```
>>bar3(Y)
```

desenează tridimensional fiecare element ca un paralelipiped, cu elementele distribuite pe axa Oy

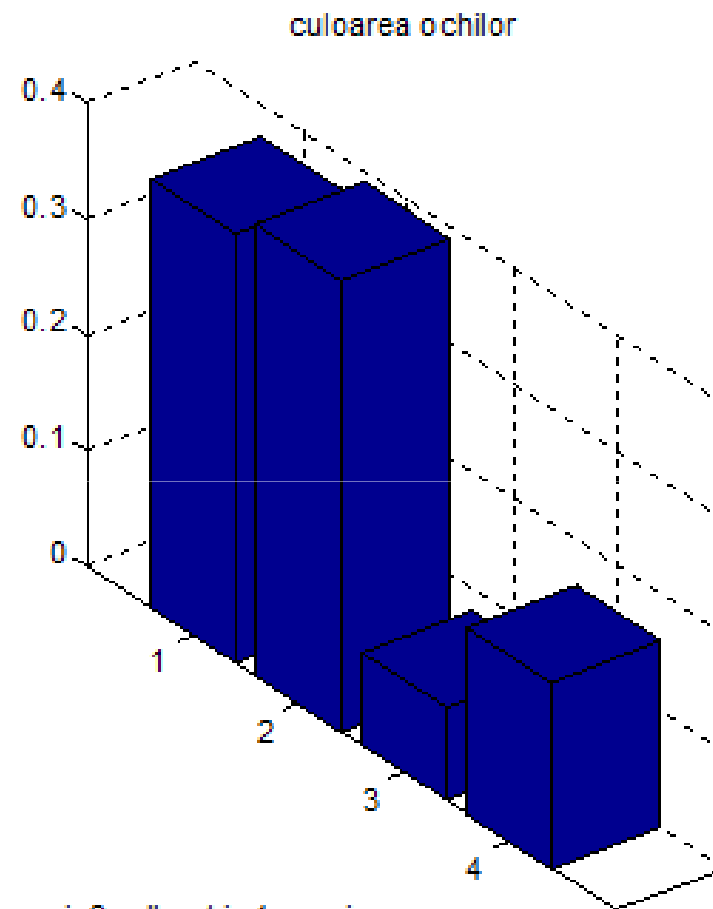
Exemplu

Reluăm exemplul nr 8, lucrând cu reprezentarea prin bastoane a frecvenței de apariție a diferitelor tipuri de culori ale ochilor în populația respectivă.

```
>> Y=[0.37 0.39 0.08 0.16];bar(Y)
```



```
>> Y=[0.37 0.39 0.08 0.16];  
>> bar3(Y)
```



1=negri, 2=caprui, 3=albastri, 4=verzi

Remarcă:

În acest context se mai folosește pentru aceste diagrame și termenul de histogramă, dar, reprezentarea prin histograme se referă la reprezentarea grafică a frecvențelor tabulate – cu referire directă la datele (atributele) cantitative.

Matematic vorbind, *histograma* reprezintă aplicația ce produce numărul de observații (valori ale datelor) ce aparțin unor anumite intervale (echivalent, frecvența observațiilor aparținând intervalelor de valori).

Vizual, prin histogramă (bidimensională) se înțelege o reprezentare grafică a acestei aplicații, adică a repartiției numărului de valori (frecvențelor) unui anumit atribut numeric în care fiecare baston (coloană) reprezintă un anumit interval de valori ale atributului iar înălțimea sa este proporțională cu numărul (frecvența) valorilor din intervalul respectiv. Termenul a fost introdus de Pearson -1895.

Reprezentarea grafică în cazul atributelor numerice, cantitative.

Histograme

În cazul unor date numerice, întâlnim cele două moduri de reprezentare grafică prin *histograme*, corespunzătoare felului datelor: discrete sau continue.

În cazul datelor *discrete*, reprezentarea grafică este asemănătoare cazului datelor calitative, cu toate că există diferența conceptuală subliniată anterior. Dacă considerăm diagramele cu bastoane, lungimea acestora are o semnificație numerică precisă. Concret pe axa absciselor sunt reprezentate valorile variabilei (datei) discrete, în timp ce pe axa ordonatelor este reprezentată frecvența relativă a apariției fiecărei valori.

Problema se complică atunci când este vorba de date numerice *continue*. Aici, pentru trasarea unei histograme, este necesară împărțirea (gruparea) datelor numerice în anumite clase (de regulă intervale), cărora să le corespundă pe cealaltă axă frecvența relativă de apariție (sau numărul de observații), corespunzătoare fiecărei clase. În general, pe axa absciselor sunt reprezentate clasele /intervale de valori, iar pe cea a ordonatelor este reprezentat numărul corespunzător de observații, înălțimea dreptunghiului fiind egală cu numărul de elemente din clasa respectivă

În Matlab putem utiliza funcțiile:

hist(x) care desenează o histogramă cu 10 intervale pentru elementele vectorului **x**;

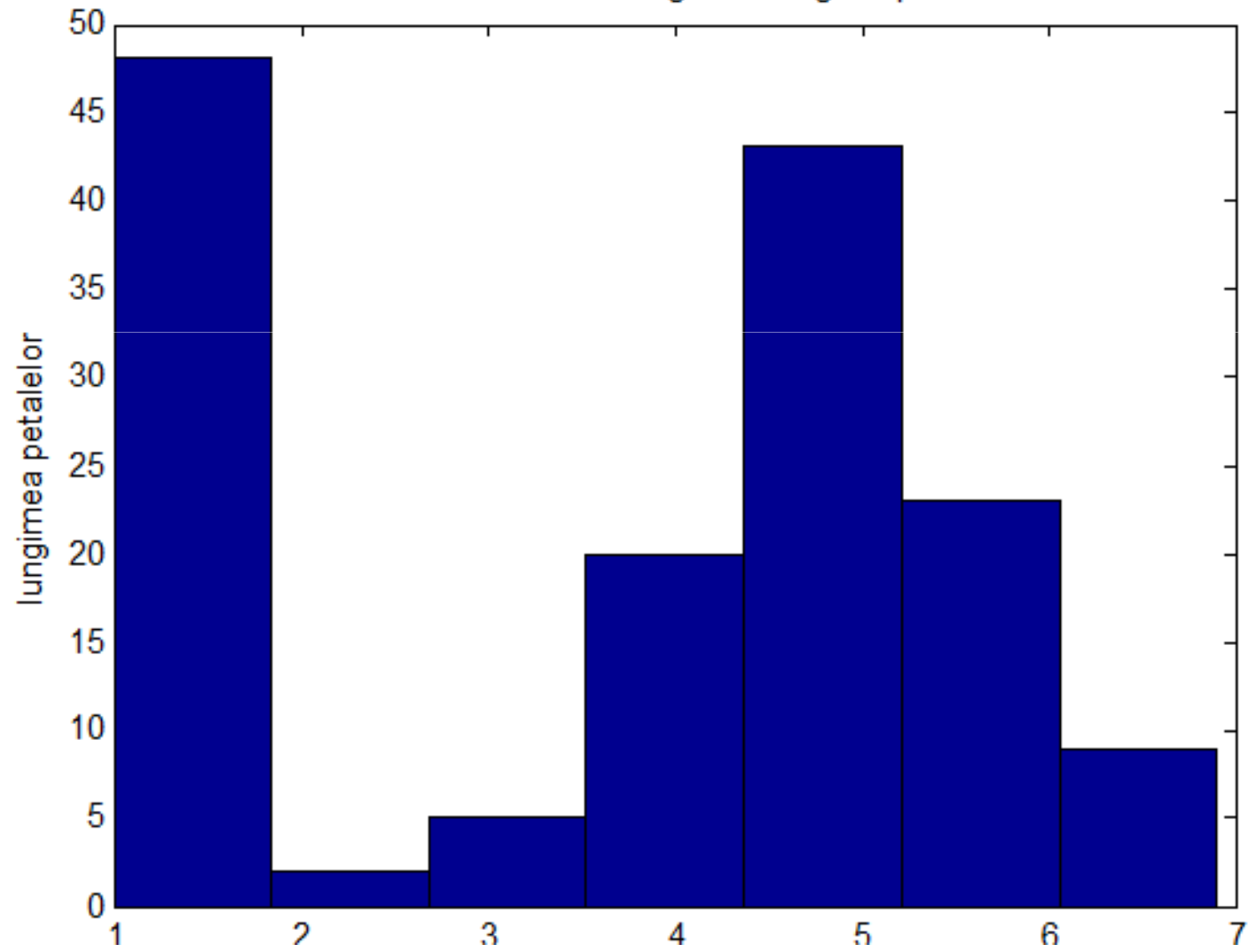
hist(x,n) care desenează o histogramă cu n intervale pentru elementele vectorului **x**;

10. Exemplu

Vom desena histograma pentru vectorul ce reprezintă lungimea petalelor florilor de iris din baza de iriși a lui Fisher. Vom calcula valorile extreme ale acestor lungimi și în funcție de aceasta vom decide numărul de intervale.

```
>> load fisheriris
>> x=meas(:,3);
>> min(x)
ans =
    1
>> max(x)
ans =
    6.9000
>> hist(x 7)
```

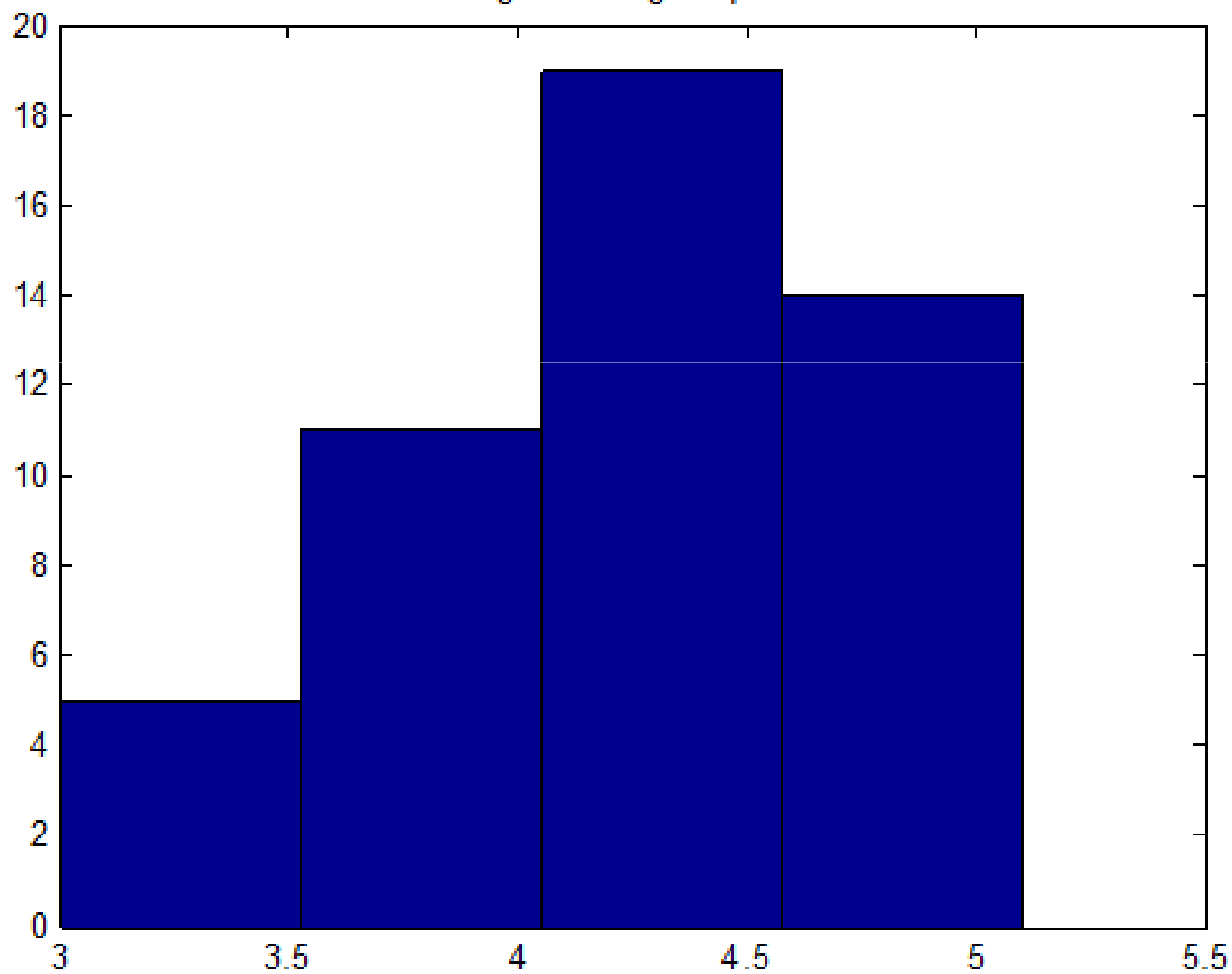
Baza de date Iris- histograma lungimii petalelor



Vom desena histograma lungimii petalelor de iris versicolor:

```
>> load fisheriris
>> versicolor_indices = strcmp('versicolor',species);
>> versicolor = meas(versicolor_indices,:);
>> versicolorT = versicolor(1:49,:);
>> x=versicolorT(:,3);
>> min(x)
ans =
    3
>> max(x)
ans =
    5.1000
>> hist(x,4)
```

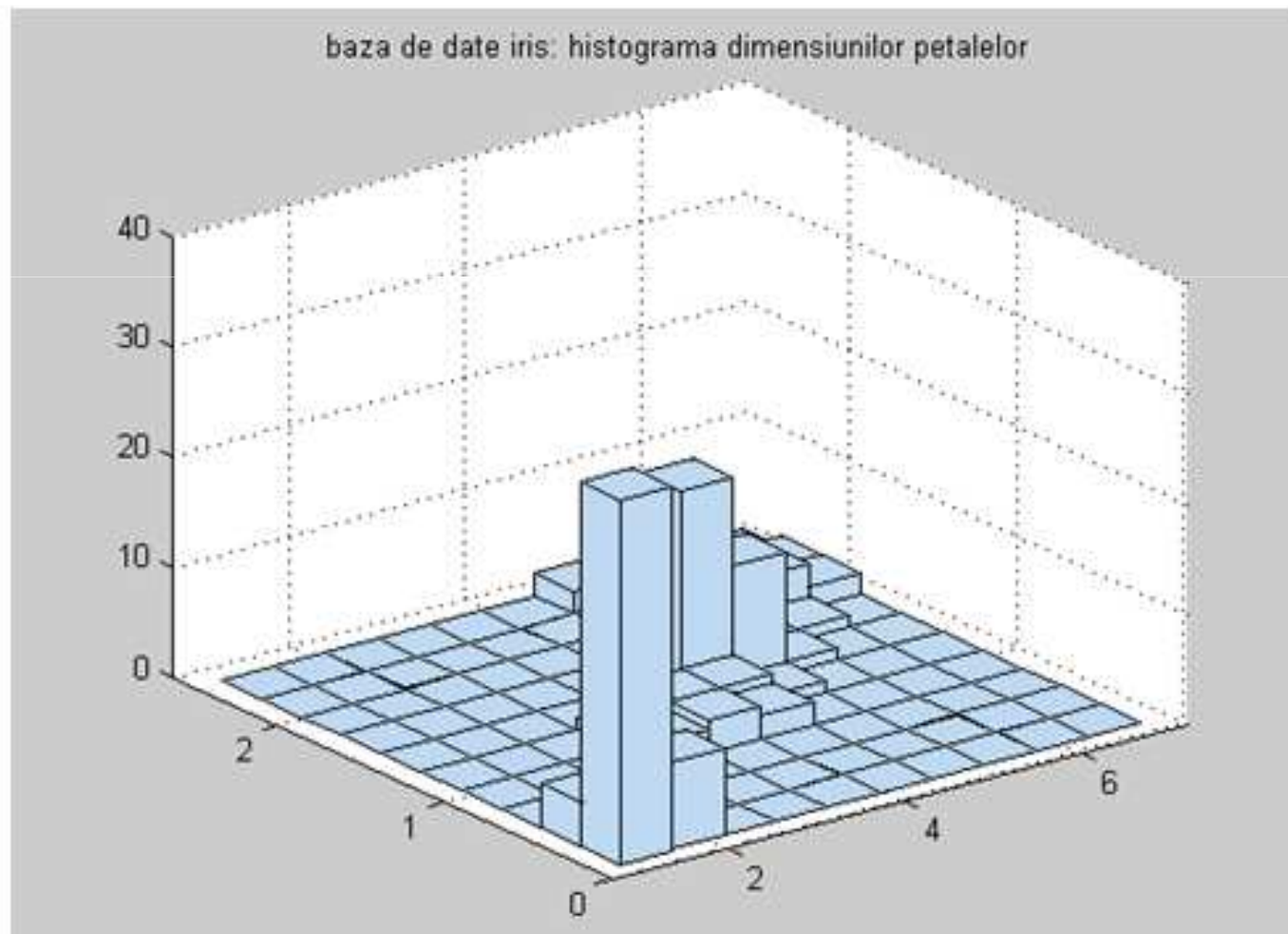
baza de date Iris- histograma lungimii petalelor de iris versicolor



`hist3(X)` depozitează elementele unei matrici de tip $m \times 2$ rețea de tip 10×10 , formată din pătrate egale și desenează o histogramă.
Fiecare coloană a matricii X corespunde unei dimensiuni a rețelei
`hist3(X, [n1,n2])` desenează histograma având în spațiu xOy o rețea de tip $n_1 \times n_2$.

Desenăm histogramele tridimensionale pentru dimensiunea petalelor în cazul bazei de date Iris:

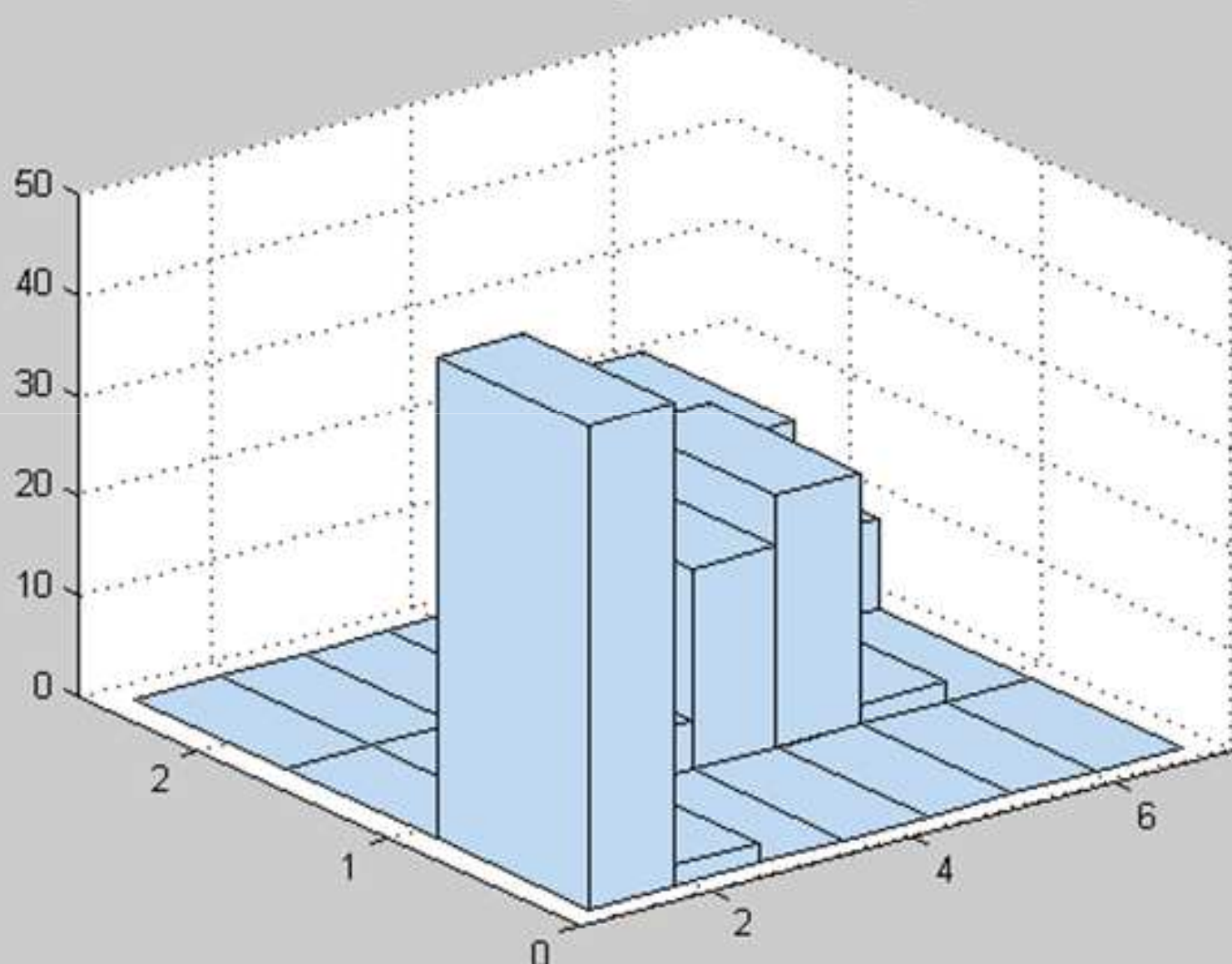
```
>> load fisheriris  
>> x=meas(:,3:4);  
>> hist3(x)
```




```
>> load fisheriris
>> x=meas(:,3:4);
>> min(meas(:,3))
ans =
    1
>> max(meas(:,3))
ans =
    6.9000
>> min(meas(:,4))
ans =

    0.1000
>> max(meas(:,4))
ans =
    2.5000
>> hist3(x,[7 3])
```

baza de date iris: histograma dimensiunii petalelor



Examinarea repartițiilor variabilelor

Asimetria(skewness)

În cazul în care repartiția datelor nu este simetrică este necesară analiza abaterii de la simetrie.

Definim *asimetria (skewness – Pearson, 1895)* ca fiind măsura deviației repartiției date de la simetrie. Formula de calcul a asimetriei este:

$$Asimetria = \frac{n \cdot \sum_{i=1}^n (x_i - \bar{x})^3}{(n-1)(n-2) \cdot \sigma^3}.$$

Dacă asimetria este diferită de zero, atunci repartiția este asimetrică.

În repartițiile asimetrice media și mediana sunt diferite.

11. Exemplu: Calculul asimetriei in Matlab

Pentru vectori `skewness(x)` este asimetria elementelor lui `x`.
Pentru matrice `skewness(X)` este un vector linie ale cărui elemente sunt asimetriile fiecărei coloane.

O parte din elevii unei clase au dat bacalaureatul la literatura română, matematică și informatică. Notele sunt prezentate sub forma unei matrice, `X`, prima coloană reprezentând notele de la română, a doua pe cele de la matematică iar a treia reprezintă notele de la informatică. Să studiem simetria acestor rezultate, calculând valoarea `skewness(X)`:

```
>> X=[7 7 6; 9 5 7; 8 4 5; 9 6 8;10 8 9; 6 9 8;9 4 5; 8 5 4; 9 8 7;10 10 8];  
>> skewness(X)  
ans =  
-0.6873  0.1862 -0.2951
```

Dacă repartiția este asimetrică, atunci ea va avea o “coadă” fie în stânga, fie în dreapta.

Dacă această “coadă” este la dreapta, spunem că repartiția are o asimetrie pozitivă, iar în celălalt caz, al “cozii” la stânga, asimetria este negativă.

12. Exemplu

Notele unei subgrupe de studenți la un examen mai dificil sunt
{ 2 2 3 3 3 4 5 5 9 10 }

```
>> x=[2, 2, 3, 3, 3, 4, 5, 5, 9, 10];  
>> m=mean(x); M=median(x); [m M]  
ans =  
    4.6000    3.5000
```

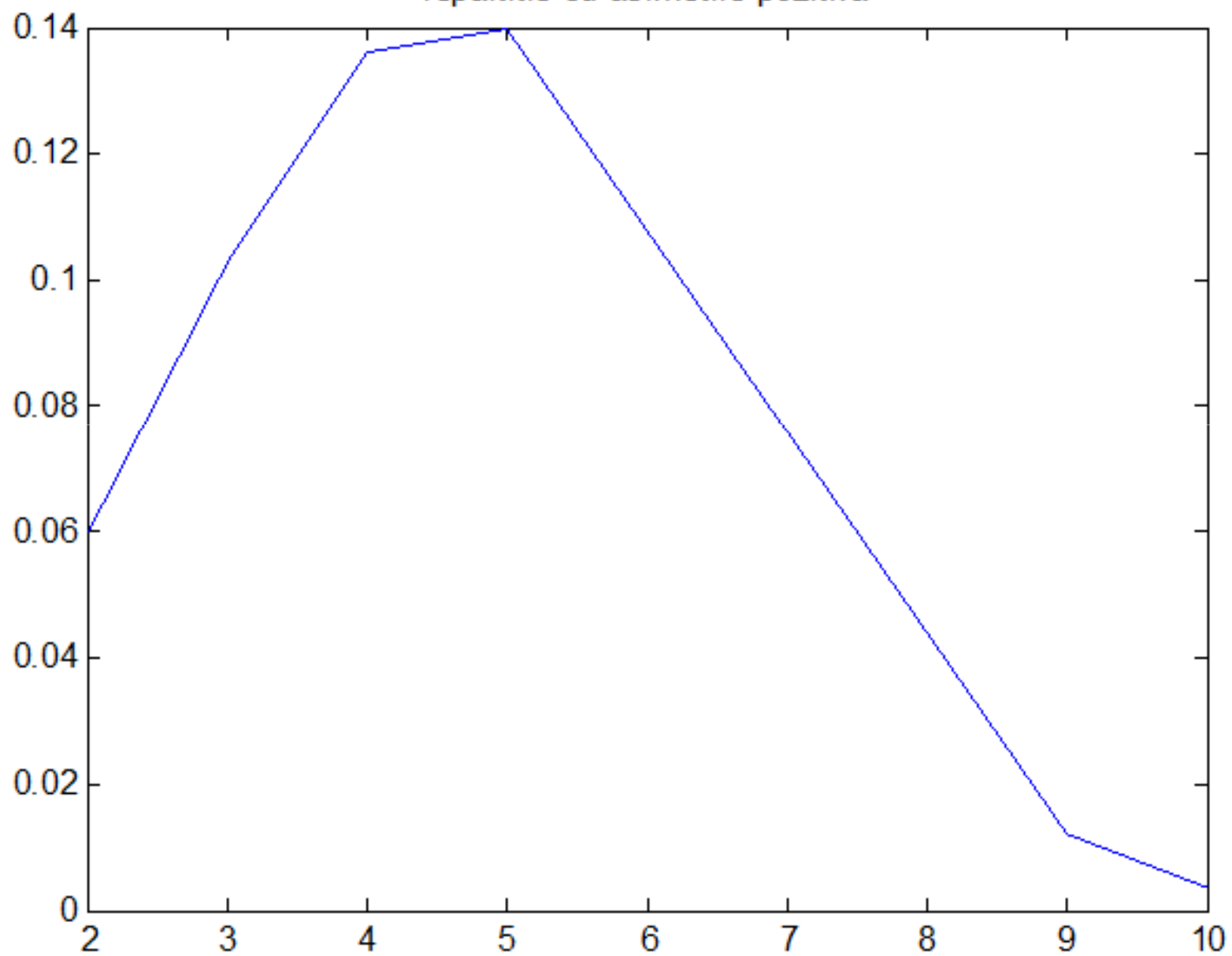
Repartiția este asimetrică.

```
>> skewness (x)  
ans =  
    1.0446
```

Repartiția are o asimetrie pozitivă

```
>> s=std(x);  
>> f=1/(s*sqrt(2*pi))*exp(-((x-m)./s).^2);plot(x,f)
```

repartitie cu asimetrie pozitiva



13. Exemplu

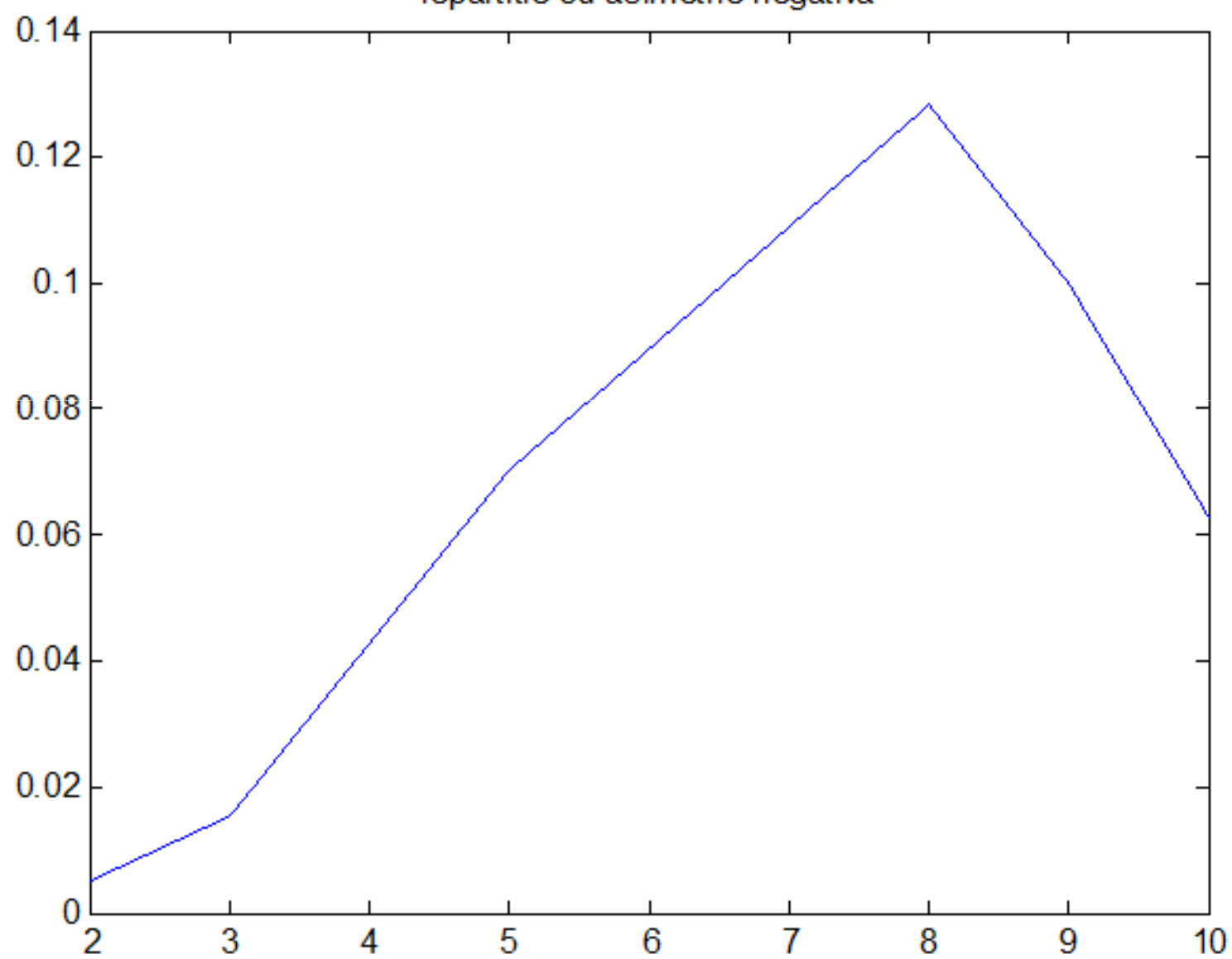
Notele aceleiași subgrupe de studenți la un examen mai ușor sunt:
{2 3 5 8 8 9 9 10 10 10 }

```
>>y=[2 3 5 8 8 9 9 10 10 10];  
>> m1=mean(y);M1=median(y); [m1 M1]  
ans =  
    7.4000    8.5000
```

```
>> skewness(y)  
ans =  
   -0.8559
```

```
>> s1=std(y);  
>> g=1/(s1*sqrt(2*pi))*exp(-((y-m1)./s1).^2);plot(y,g)
```

repartitie cu asimetrie negativa



Repartiția normală (gaussiană) este perfect simetrică, reprezentând modelul generic de simetrie. Reamintim că funcția de repartiție, cunoscută în acest caz sub numele de funcția gaussiană este definită de:

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \cdot e^{-\frac{(x-\mu)^2}{\sigma^2}},$$

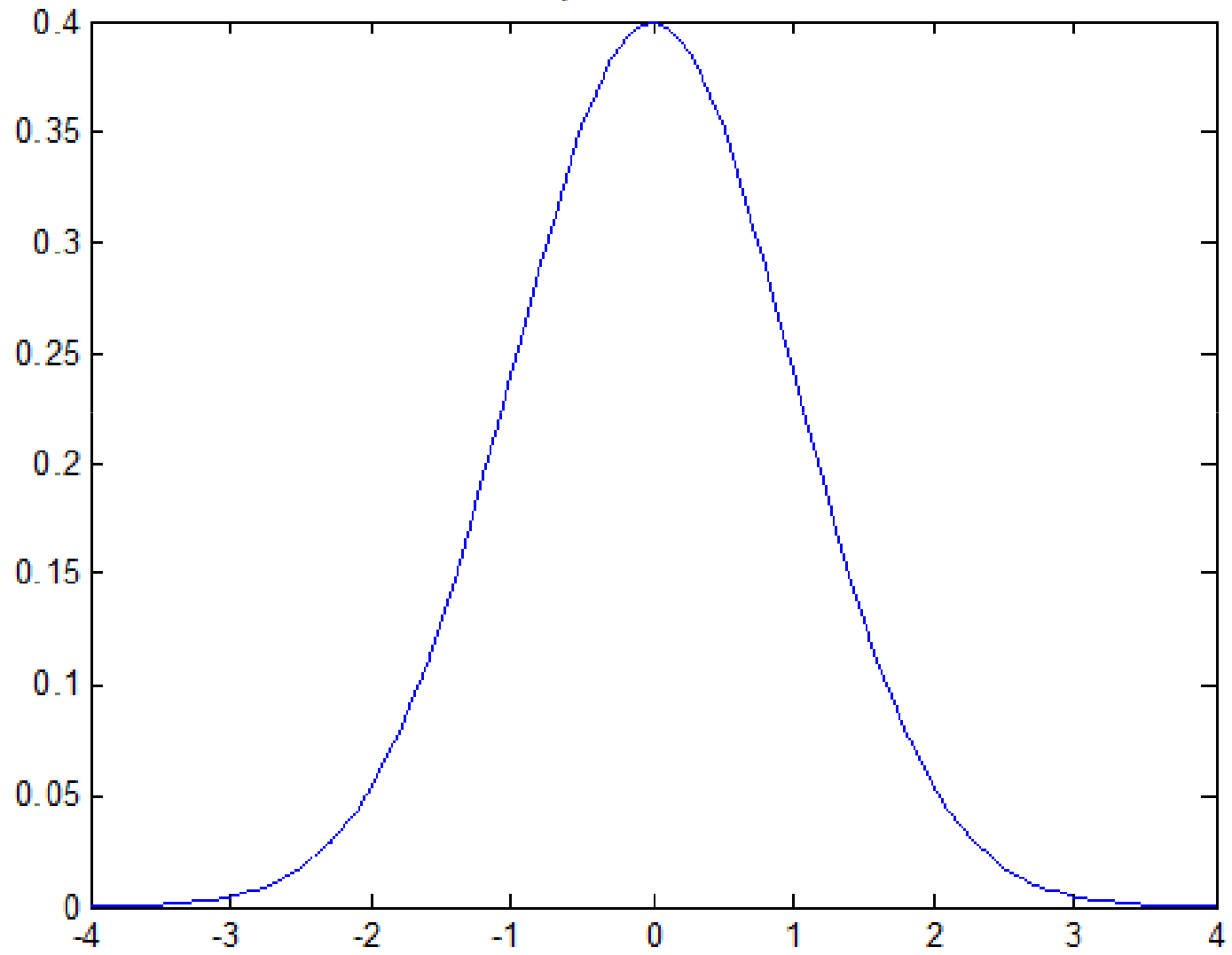
unde μ este media, σ^2 este dispersia și σ deviația standard.

Repartiția cu $\mu = 0$ și $\sigma = 1$ este repartiția normală standard

Figura următoare reprezintă „clopotul lui Gauss”, adică cea mai „simetrică” repartiție.

```
>>x=-4:.1:4; plot(x,1/sqrt(2*pi)*exp(-x.^2/2))
```

clopotul lui Gauss



Excesul (kurtosis)

Cealaltă măsură tipică formei repartiției, *excesul* (*kurtosis* Pearson, 1905), măsoară „ascuțimea” unei repartiții, în sensul cât este aceasta de dispusă de a avea valori extreme (outliers).

Dacă excesul este diferit de zero, atunci repartiția este ori mai „plată” ori mai „ascuțită” decât repartiția normală, care este chiar zero.

Excesul se calculează utilizând formula următoare:

$$Exces = \frac{n \cdot (n + 1) \cdot \sum_1^n (x_i - \bar{x})^4 - 3 \cdot (n - 1) \left[\sum_1^n (x_i - \bar{x})^2 \right]^2}{(n - 1) \cdot (n - 2) \cdot (n - 3) \cdot \sigma^4}$$

Această formulă este folosită pentru calculele în Excel.

În Matlab se utilizează formula

$$K = \frac{\sum_{i=1}^n (x_i - \bar{x})^4}{\sigma^4}, \text{ unde } \bar{x} \text{ este media iar } \sigma \text{ este deviatia standard.}$$

Avem egalitatea $K = \text{Exces} + 3$

Funcția `kurtosis(x)` returnează în cazul vectorului x , excesul K al lui x , iar în cazul matricei X , un vector linie, ale cărui componente sunt valorile Kurtosis-ului pentru fiecare coloană în parte.

Exemple

Reluăm exemplul nr 11 (cu notele de la bacalaureat) calculând asimetria și excesul:

```
X=[7 7 6; 9 5 7; 8 4 5; 9 6 8;10 8 9; 6 9 8;9 4 5; 8 5 4; 9 8 7;10 10 8];
```

```
>> skewness(X)
```

```
ans =
```

```
    -0.6873    0.1862   -0.2951
```

```
>> kurtosis(X)
```

```
ans =
```

```
    2.5981    1.7103    1.8363
```

Reluăm exemplele nr 12 (cu notele la un examen dificil) și nr 13 (cu notele la un examen ușor) calculând asimetria și excesul.

```
>> x=[2, 2, 3, 3, 3, 4, 5, 5, 9, 10];
```

```
>> skewness(x)
```

```
ans =
```

```
1.0446
```

```
>> kurtosis(x)
```

```
ans =
```

```
1.7996
```

```
>> y=[2 3 5 8 8 9 9 10 10 10];
```

```
>> skewness(y)
```

```
ans =
```

```
-0.8559
```

```
>> kurtosis(y)
```

```
ans =
```

```
2.1793
```


Datele date întâlnite frecvent în realitate sunt caracterizate de repartiții non-normale (negaussiene), care sunt mai mult sau mai puțin îndepărtate de repartiția normală.

Vom prezenta, patru asemenea repartiții, cele mai des întâlnite.

Repartiția exponențială

Repartiția *exponențială* $X \sim \text{Exp}(\lambda)$, descrie timpii dintre evenimentele unui proces Poisson, (proces în care evenimentele apar continuu și independent cu o rată medie constantă).

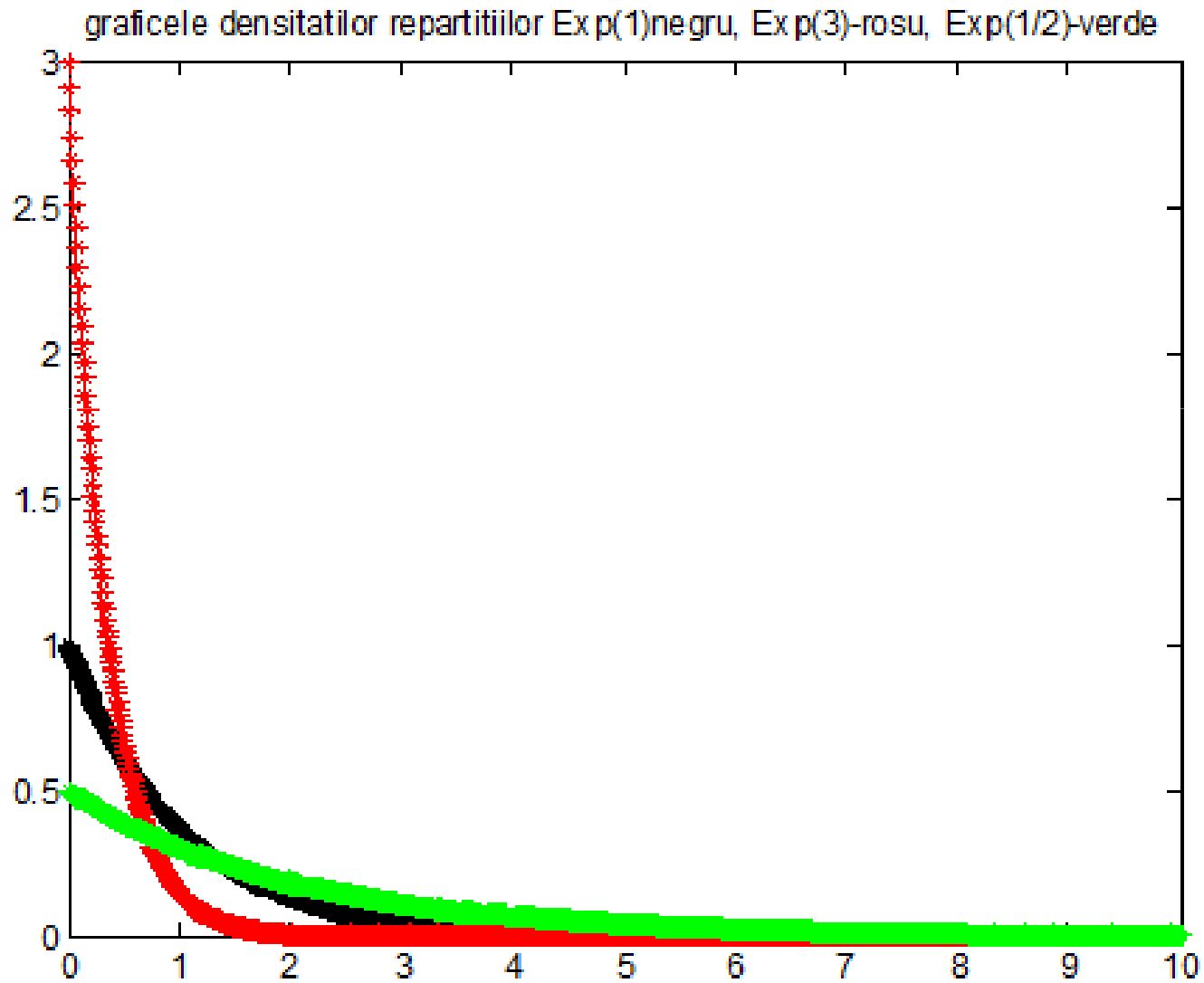
Repartiția *exponențială* de parametru $\lambda > 0$ este caracterizată

de densitatea $f_X(x) = \begin{cases} 0, & x < 0 \\ \lambda e^{-\lambda x}, & x \geq 0 \end{cases}$.

Media corespunzătoare este $E(X) = \frac{1}{\lambda}$.

Dispersia corespunzătoare este $D^2(X) = \frac{1}{\lambda^2}$.

```
>> x=0:.01:10;plot(x,exp(-x),'k',x,3.*exp(-3.*x),'r',x,1./2.*exp(-x./2),'g')
```

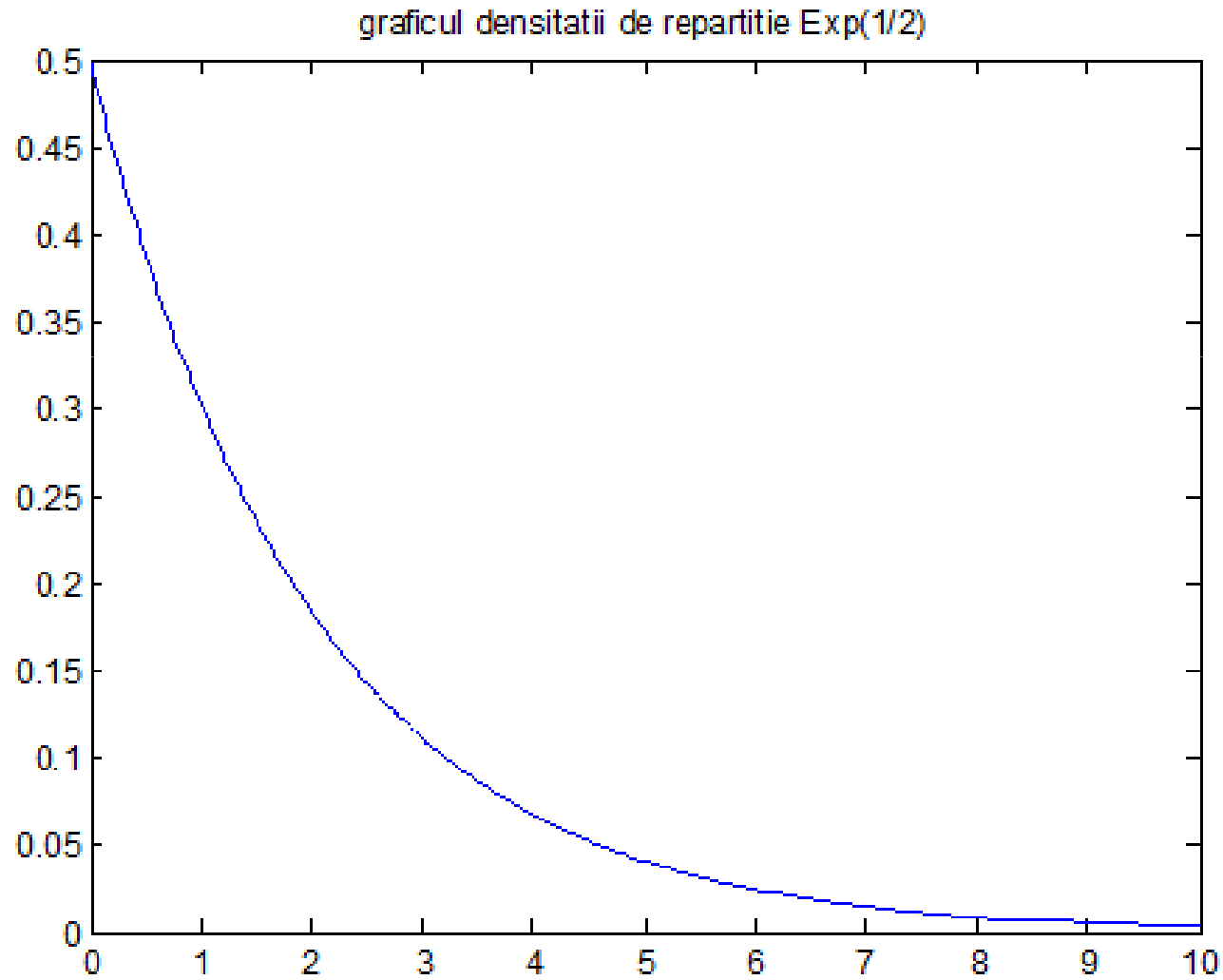


Graficul densității de repartiție (*probability density function* – pdf) în cazul repartiției exponențiale poate fi desenat în Matlab folosind sintaxa:

$Y = \text{exppdf}(X, \mu)$

care returnează valorile în X ale densității de repartiție exponențială de parametru pozitiv μ , unde μ reprezintă media lui X

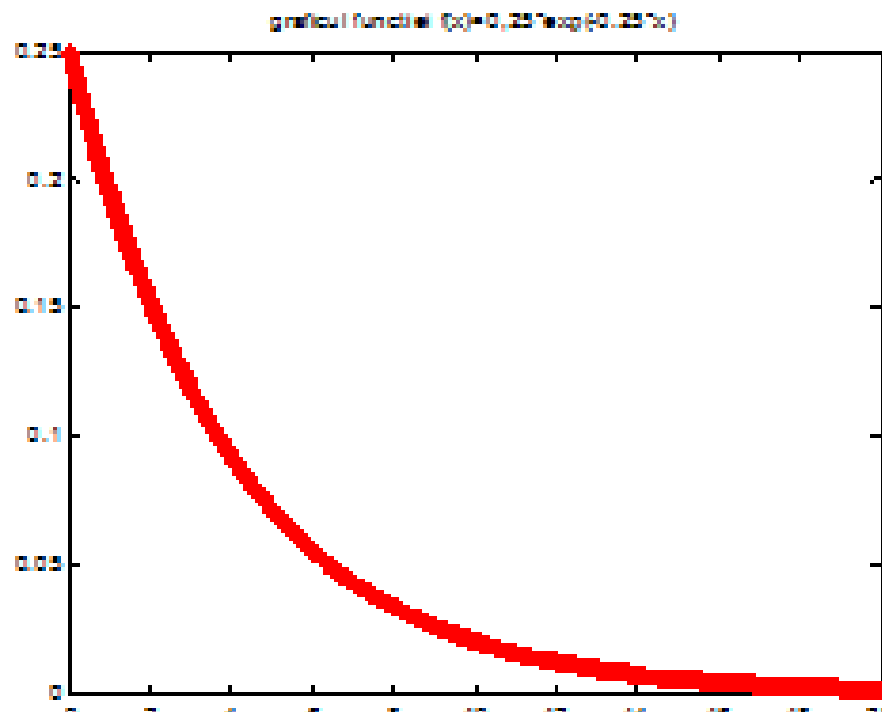
```
>> X = (0:0.02:10);  
>> Y = exppdf(X,2);  
>> plot(X,Y)
```



14. Exemplu

X reprezintă timpul, în minute, pe care un funcționar de la ghișeu de informații al unei bănci îl alocă fiecărui client. Se știe că acești timpi au o repartiție exponențială, iar timpul mediu alocat unui client este de 4 minute. Astfel

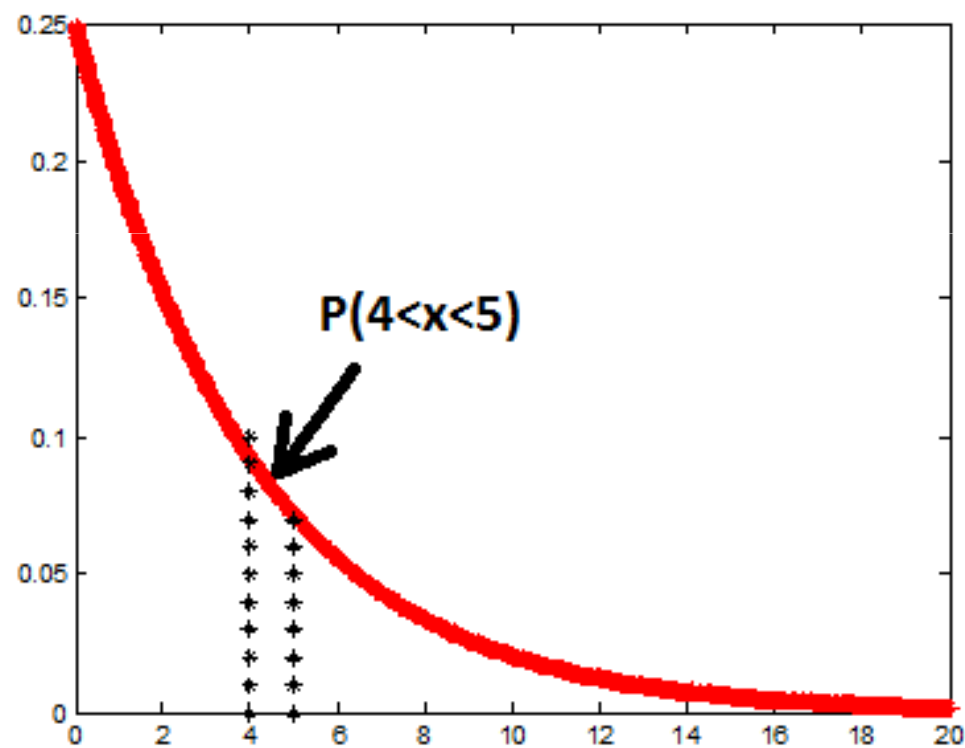
$\lambda = \frac{1}{\mu}$, unde $\mu = 4$. Funcția de repartiție este $f(x) = \frac{1}{4} e^{-\frac{1}{4}x}$ pentru $x \geq 0$



Să calculăm probabilitatea ca funcționarul să aloce unui client între 4 și 5 minute, adică $P(4 < x < 5)$

Știm că $P(x < a) = 1 - e^{-\lambda a}$ și astfel

$$P(x < 4) = 1 - e^{-0.25 \cdot 4} = 0.6321 \text{ și } P(x < 5) = 1 - e^{-0.25 \cdot 5} = 0.7135$$



Ordonăm crescător timpii alocați fiecărui client. Care este media de timp alocată fiecărui client din prima jumătate? Această înseamnă să găsim a 50-a percentilă, rezolvând ecuația $P(x < k) = 1 - e^{-0.25 \cdot k} = 0.5$

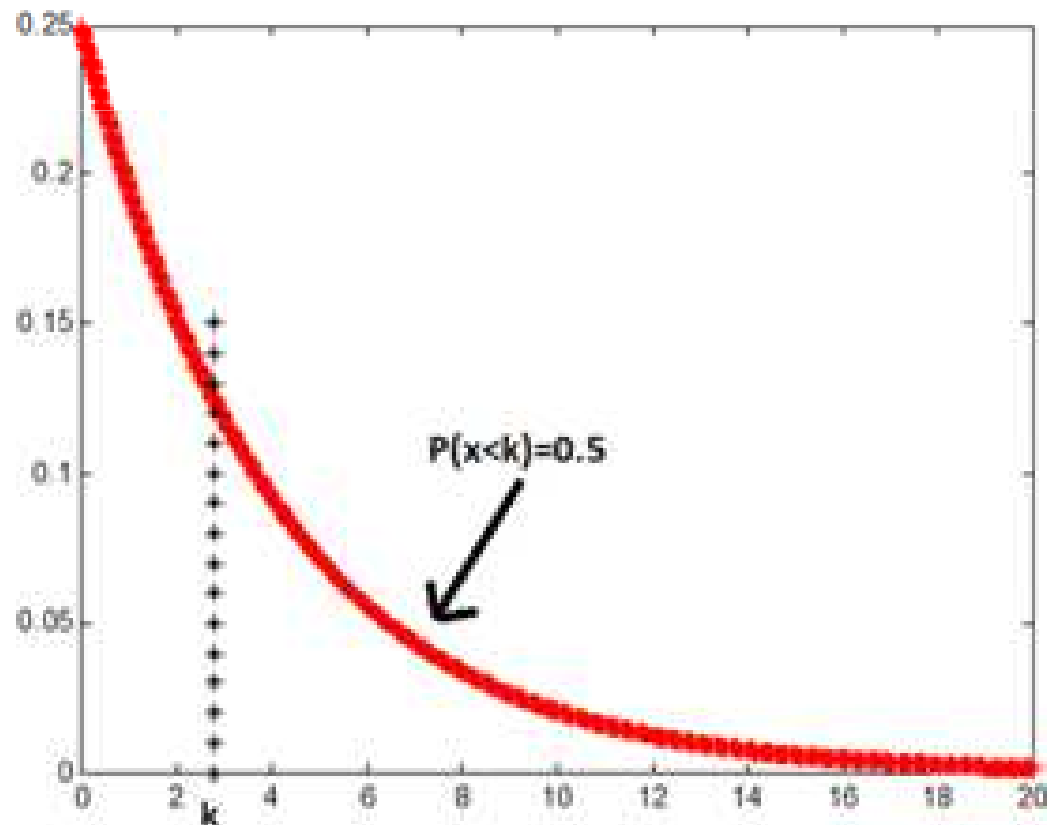
```
>> syms k
```

```
>> solve(0.5-exp(-0.25*k))
```

```
ans =
```

```
4*log(2)
```

Avem: $k=2.7726$



Care este mai mare media sau mediană?

Media= 4 (enunțul problemei);

Mediana =2.7726 (50-a percentilă).

Generarea de numere aleatoare exponențial repartizate - Matlab

`R = exprnd(mu)` generează aleator numere exponențial repartizate; `mu` poate fi un vector sau o matrice. Dimensiunea lui `R` este egală cu dimensiunea lui `mu`

```
>> R = exprnd(1:8)
```

```
R =
```

```
1.2840 6.1508 6.9951 0.7770 1.8205 6.8912 0.3574 26.9469
```

```
>> R = exprnd(1:8)
```

```
R =
```

```
0.8238 1.9270 0.8016 0.9166 8.3866 4.2830 5.6586 3.4918
```

```
>> A=[1 3; 2 4; 7 5];
```

```
>> R = exprnd(A)
```

```
R =
```

```
1.8476 2.1685
```

```
0.0597 0.8912
```

```
0.3064 9.7636
```

```
>> R = exprnd(A)
```

```
R =
```

```
0.8633 0.1241
```

```
0.1761 1.6880
```

```
4.0005 10.0004
```

Repartiția gamma

Repartiția *gamma* $X \sim \text{GAM}(\lambda, k)$, de parametri $\lambda > 0$ și $k > 0$, are densitatea dată de:

$$f_X(x) = \frac{\lambda^k x^{k-1} e^{-\lambda x}}{\Gamma(k)}, \quad x \geq 0 \quad \text{și} \quad f_X(x) = 0, \quad x < 0,$$

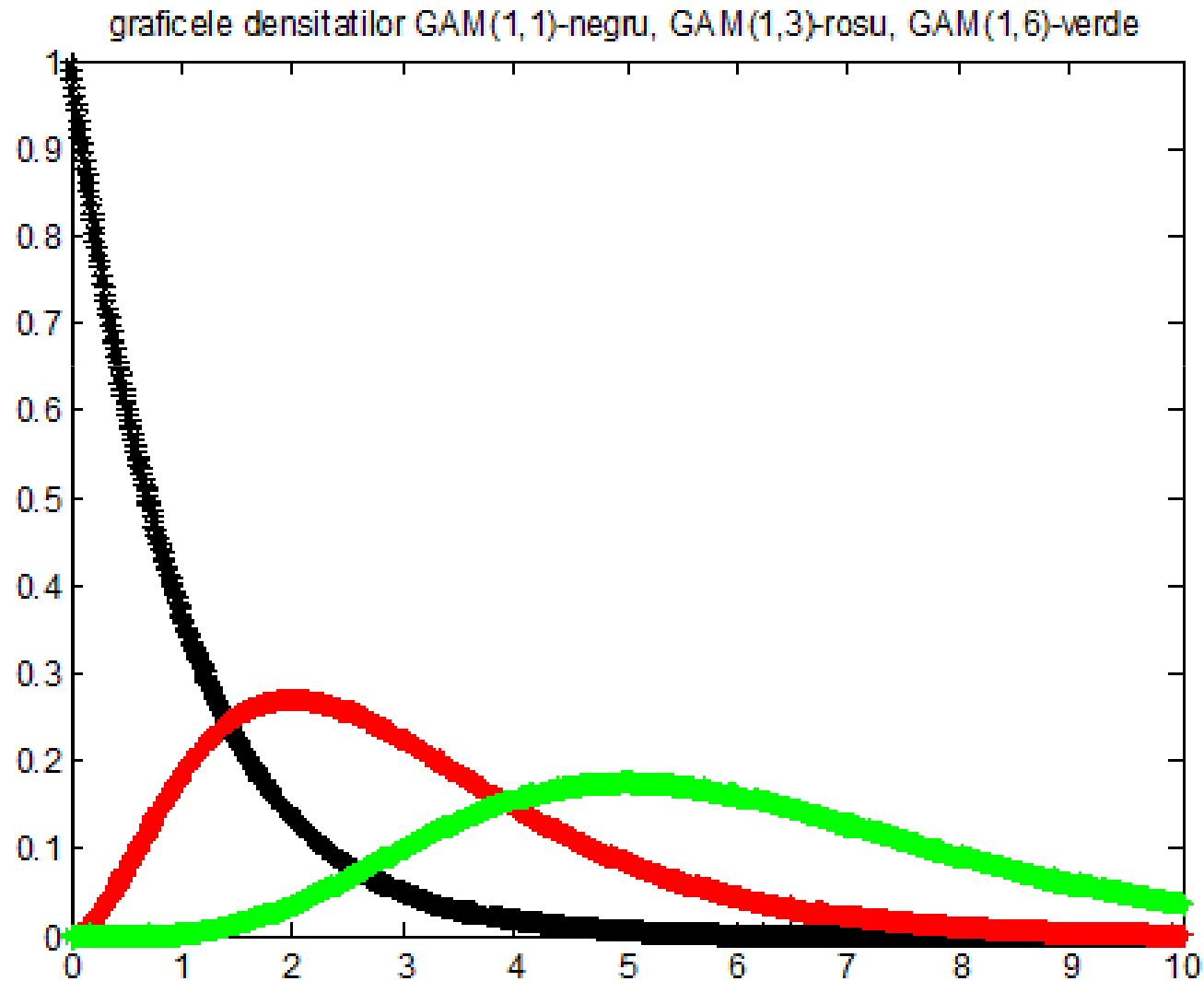
unde $\Gamma(k) = \int_0^{\infty} e^{-x} x^{k-1} dx$ este funcția *gamma*.

Media corespunzătoare este $E(X) = \frac{k}{\lambda}$.

Dispersia corespunzătoare este $D^2(X) = \frac{k}{\lambda^2}$

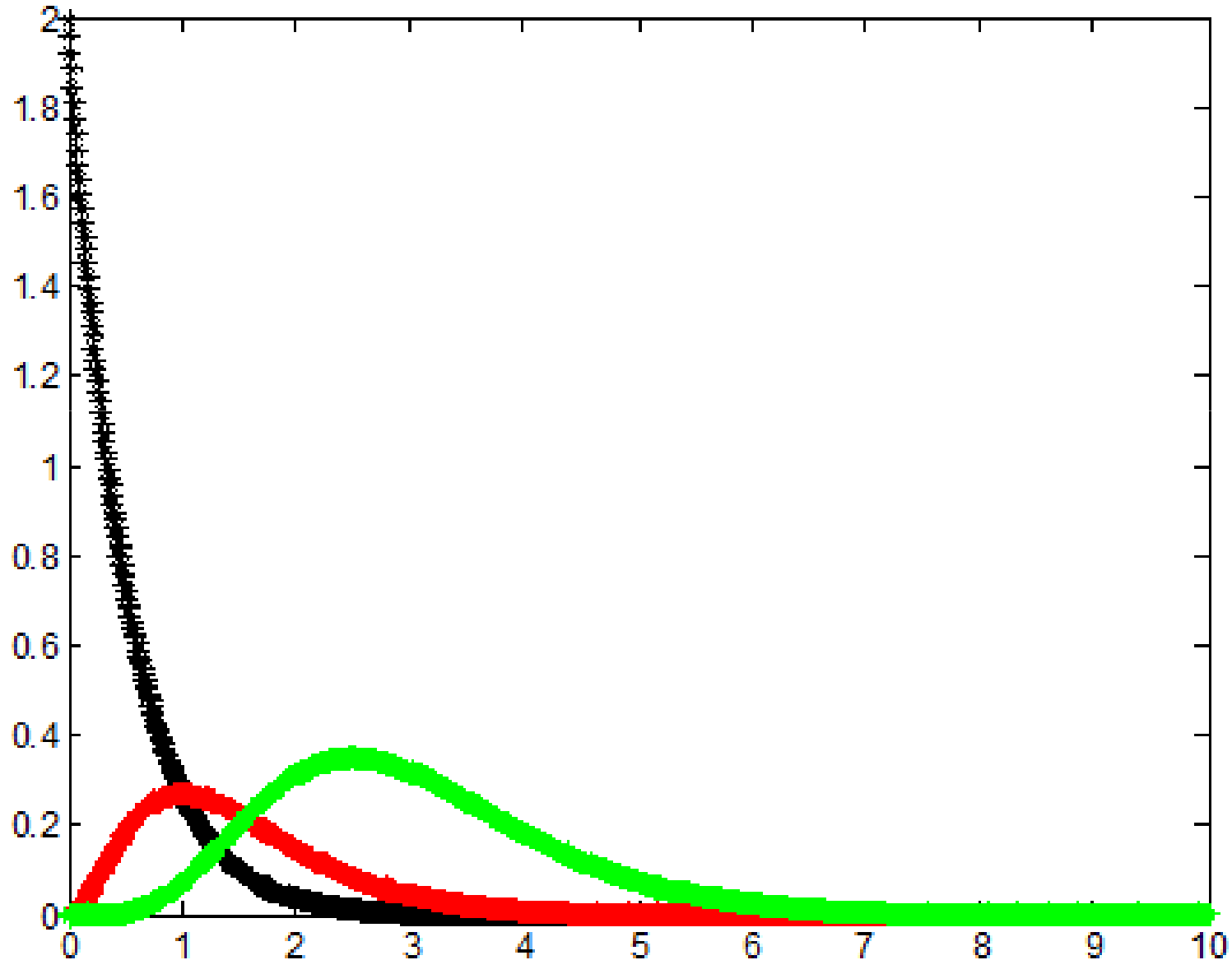
$$\text{GAM}(\lambda, 1) = \lambda \cdot e^{-\lambda x} = \text{Exp}(\lambda)$$

```
>> syms x
>> k=1; gam(1)=int(exp(-x),0,inf); for k=2:6 gam(k)=int(x^(k-1)*exp(-x),0, inf);end
>> x=0:.01:10;plot(x,exp(-x)/gam(1),'k*',x,(x.^2).*exp(-x)/gam(3),'r*',
x,(x.^5).*exp(-x)/gam(6),'g*')
```



```
>> x=0:.01:10;plot(x,2*exp(-2*x)/gam(1),'k',x,2^2*(x.^2).*exp(-2*x)/gam(3),'r',  
x,2^6*(x.^5).*exp(-2*x)/gam(6),'g')
```

graficele densitatilor repartitiilor GAM(2,1)-negru, GAM(2,3)-rosu,GAM(2,6)-verde



Graficul densității de repartiție (*probability density function* – pdf) în cazul repartiției gamma poate fi desenat folosind sintaxa

`Y = gampdf(X,k,lamda)`

care returnează valorile în X ale densității repartiției $GAM(\lambda, k)$.

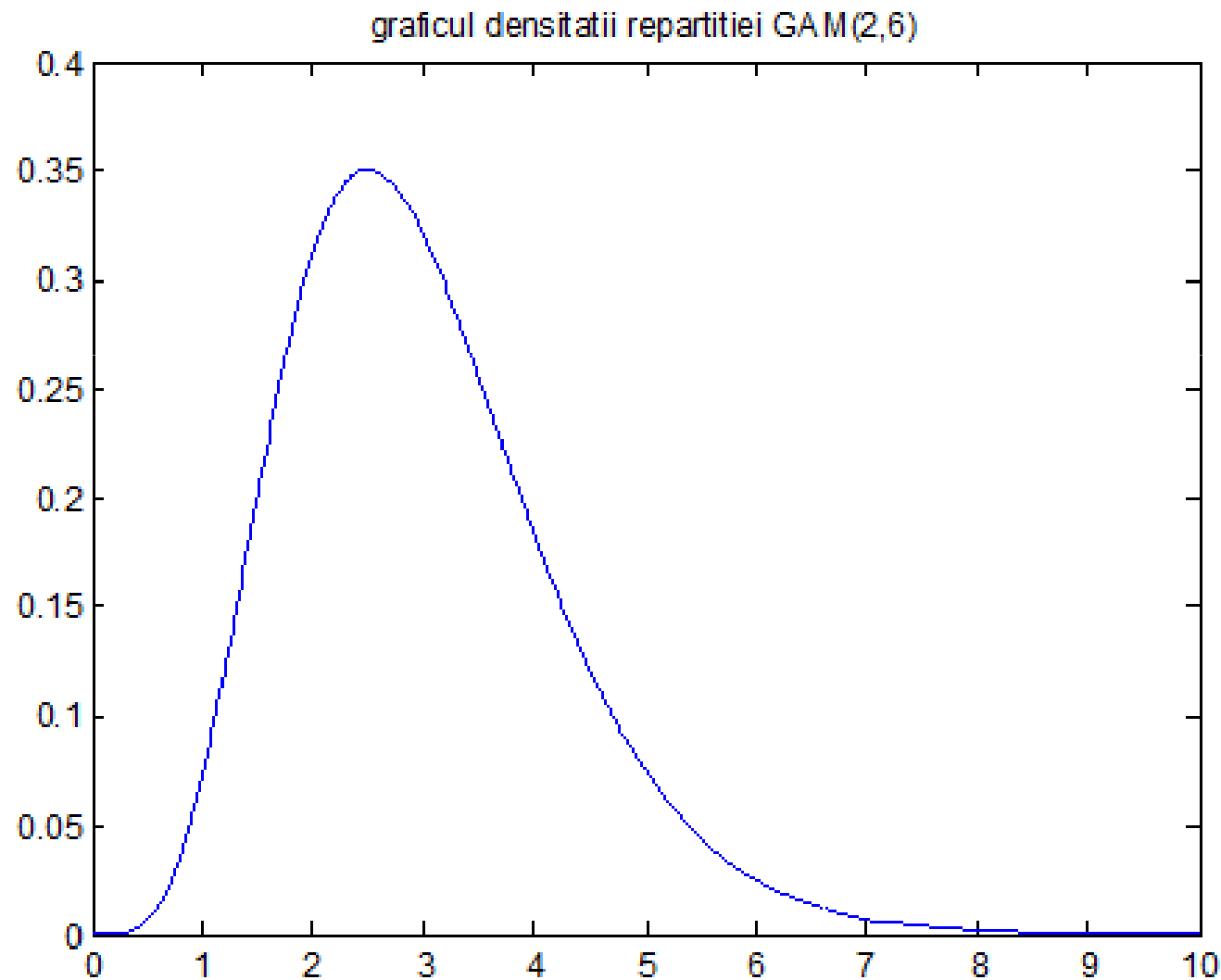
Este obligatoriu de reținut că în Matlab funcția de repartiție gamma are formula:

$$f_X(x) = \frac{x^{k-1} e^{-\frac{x}{\lambda}}}{\lambda^k \Gamma(k)}$$

care este o altă exprimare pentru formula anterioară prin înlocuirea lui λ cu $\frac{1}{\lambda}$.

Să desenăm densitatea repartiției GAM(2,6):

```
>> X=0:.01:10;Y=gampdf(X,6,1/2);plot(X,Y)
```



Dacă X_1, X_2, \dots, X_k sunt variabile aleatoare exponențiale, independente și identic repartizate de parametru λ , atunci variabila sumă

$$X_1 + X_2 + \dots + X_k$$

este o variabilă gamma de parametrii λ și k .

Generarea de numere aleatoare gamma repartizate - Matlab

`R = gamrnd(A,B,m,n)` generează aleator numere cu gamma repartizate de parametri `A` și `B`, scalarii `m` și `n` reprezentând numărul de linii, respectiv coloane a lui `R`.

```
>> R = gamrnd(1,1,1,7)
```

```
R =
```

```
2.8643 1.0182 2.0584 0.6936 2.0296 0.7393 0.2257
```

```
>> R = gamrnd(1,1,1,7)
```

```
R =
```

```
5.0164 1.4830 1.6873 0.0521 1.1500 0.5191 0.9342
```

```
>> R = gamrnd(2,5,1,9)
```

```
R =
```

```
23.3573 13.6873 6.8571 2.8459 9.0288 38.6431 9.6020 1.0473
```

```
1.2877
```

```
>> R = gamrnd(2,5,1,9)
```

```
R =
```

```
3.8150 5.2894 5.0279 14.0750 7.1420 4.0038 3.1517 9.1600
```

```
2.7154
```

Repartiția Weibull

Repartiția *Weibull*, $X \sim \text{WEI}(\alpha, \beta)$, de parametri $\alpha > 0$ și $\beta > 0$ are densitatea dată de:

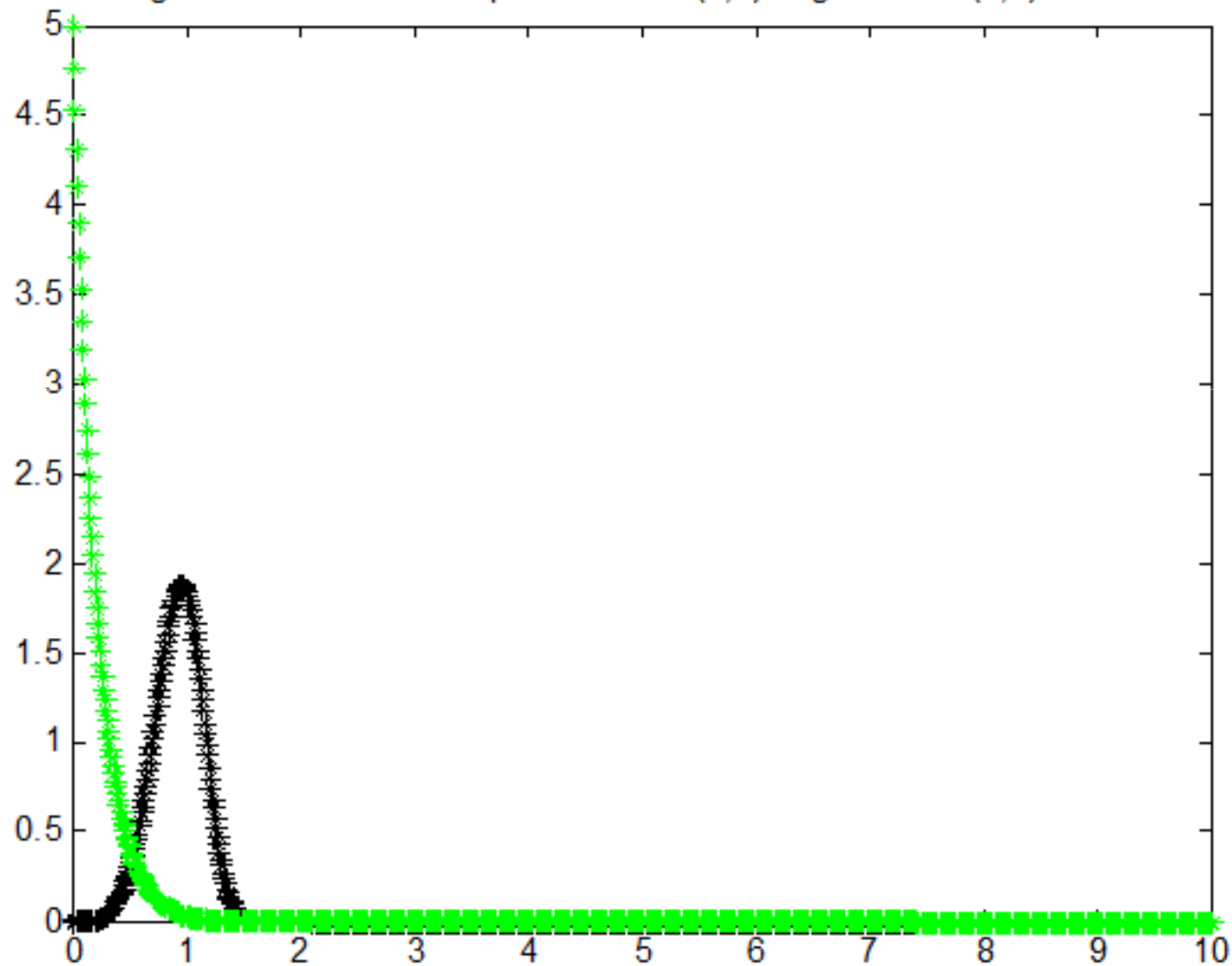
$$f_X(x) = \alpha\beta(\alpha x)^{\beta-1} e^{-(\alpha x)^\beta}, x \geq 0 \text{ și } f_X(x) = 0, x < 0,$$

Este un instrument în modelarea rezistenței materialor și în modelarea timpilor de supraviețuire.

Prezentăm graficele densităților de repartiție pentru diferite valori ale parametrilor α și β

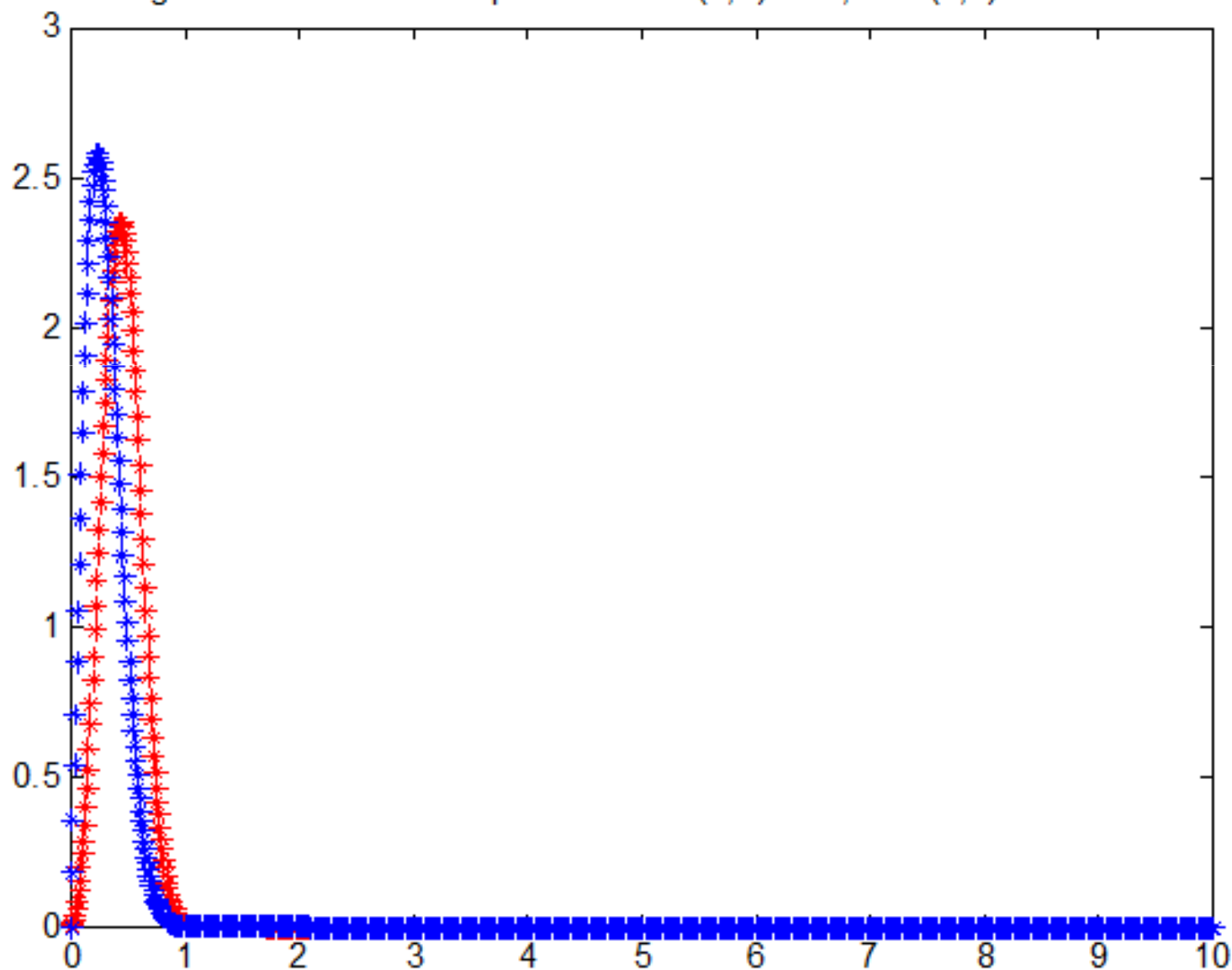
```
>>x=0:.01:10;plot(x,(5*x.^4).*exp(-x.^5),'k*',x,5*exp(-5*x),'g*')
```

graficele densitatilor repartitiilor WEI(1,5)-negru si WEI(5,1)-verde



```
>>x=0:.01:10;plot(x,6*(2*x).^2.*exp((-2*x).^3),'r',x,6*(3*x).^2.*exp(-(3*x).^2),'b')
```

graficele densitatilor repartitiilor WEI(2,3)-rosu, WEI(3,2)-albastru

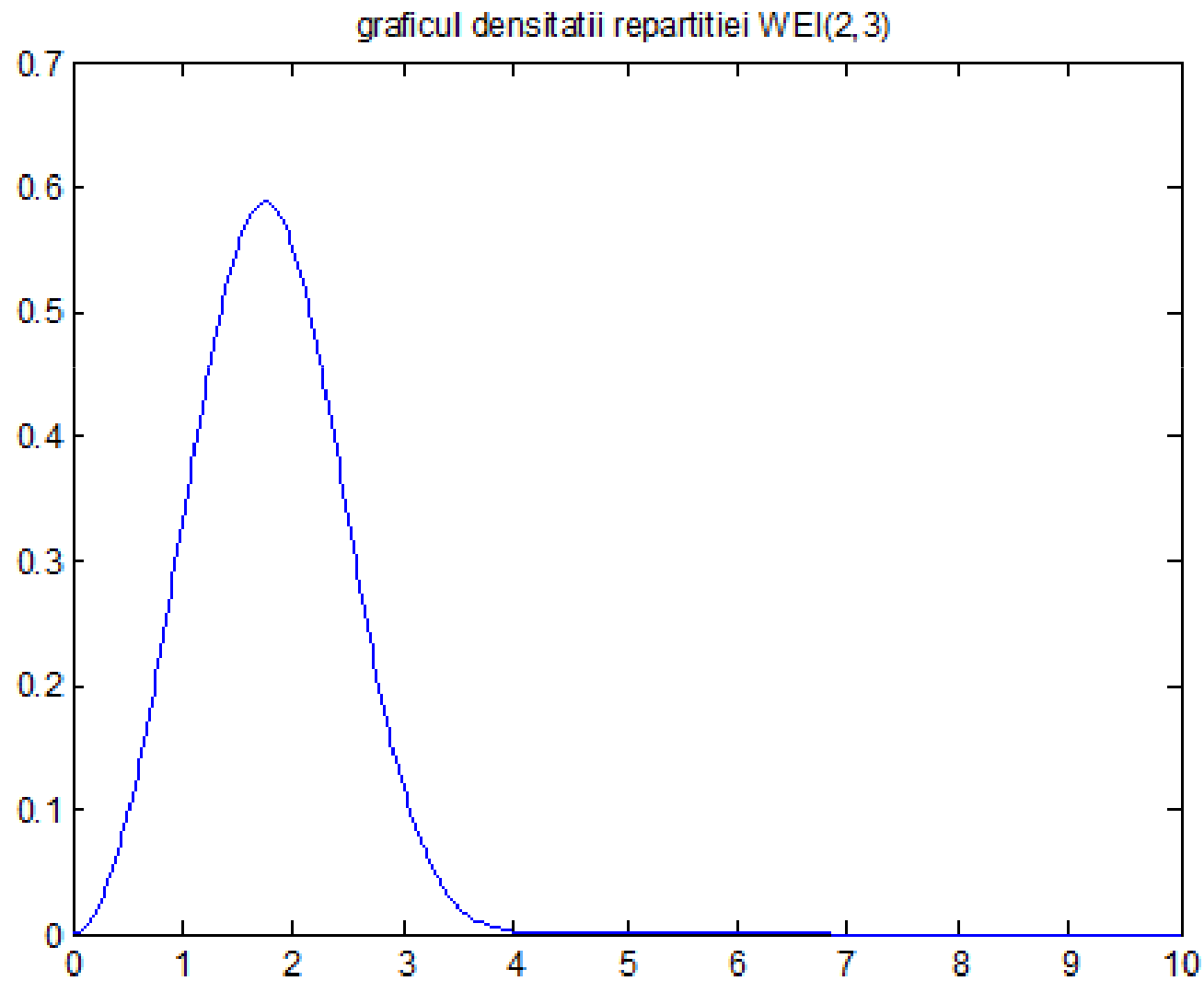


Graficul densității de repartiție (*probability density function* – pdf) în cazul repartiției Weibull poate fi desenat folosind sintaxa

$Y = \text{wblpdf}(X, A, B)$

care returnează valorile în X ale densității de repartiție Weibull de parametri A și B .

```
>> X = (0:0.02:10);  
>> Y = wblpdf(X,2,3);  
>> plot(X,Y)
```



Generarea de numere aleatoare Weibull repartizate - Matlab

$R = wblrnd(A,B,m,n)$ generează aleator numere Weibull distribuite de parametri A și B , scalarii m și n reprezentând numărul de linii, respectiv coloane a lui R .

```
>> R = wblrnd(1,5,1,7)
```

```
R =
```

```
0.7283 0.6296 1.1559 0.6186 0.8555 1.1841 1.0503
```

```
>> R = wblrnd(1,5,1,7)
```

```
R =
```

```
0.9039 0.5340 0.5136 1.1306 0.4954 0.5349 0.9371
```

```
>> R = wblrnd(2,5,1,8)
```

```
R =
```

```
1.4812 2.2864 1.9421 1.2301 1.4944 1.0576 1.6830 2.5444
```

```
>> R = wblrnd(2,5,1,8)
```

```
R =
```

```
1.3924 1.1692 1.6546 1.5476 1.5687 1.9737 1.6833 2.2407
```

Repartiția normală

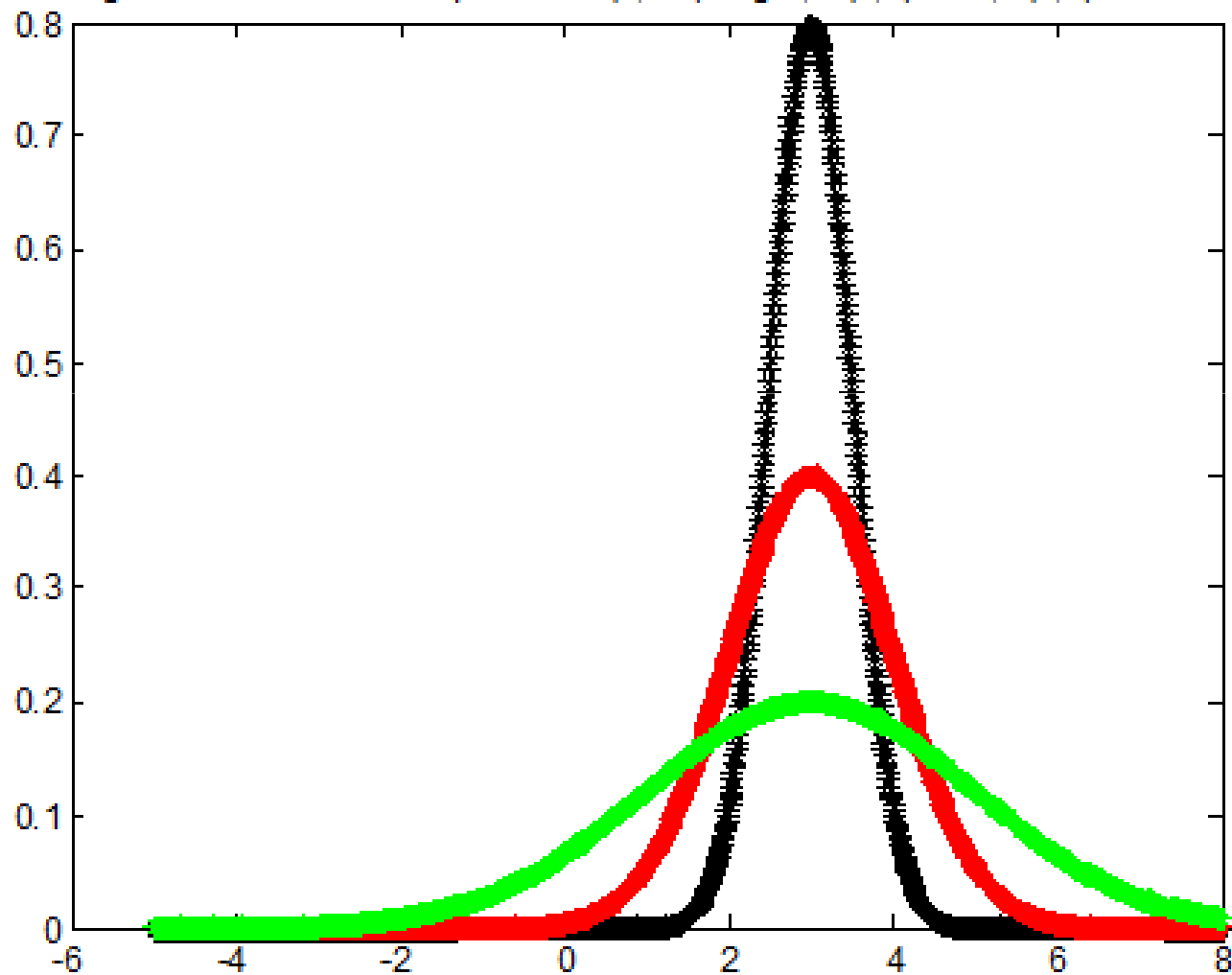
Repartiția normală (legea lui Gauss) $X \sim N(\mu, \sigma^2)$ de parametri μ și σ^2

are densitatea de repartiție $f(x) = \frac{1}{\sigma\sqrt{2\pi}} \cdot e^{-\frac{(x-\mu)^2}{\sigma^2}}$, unde μ este media,

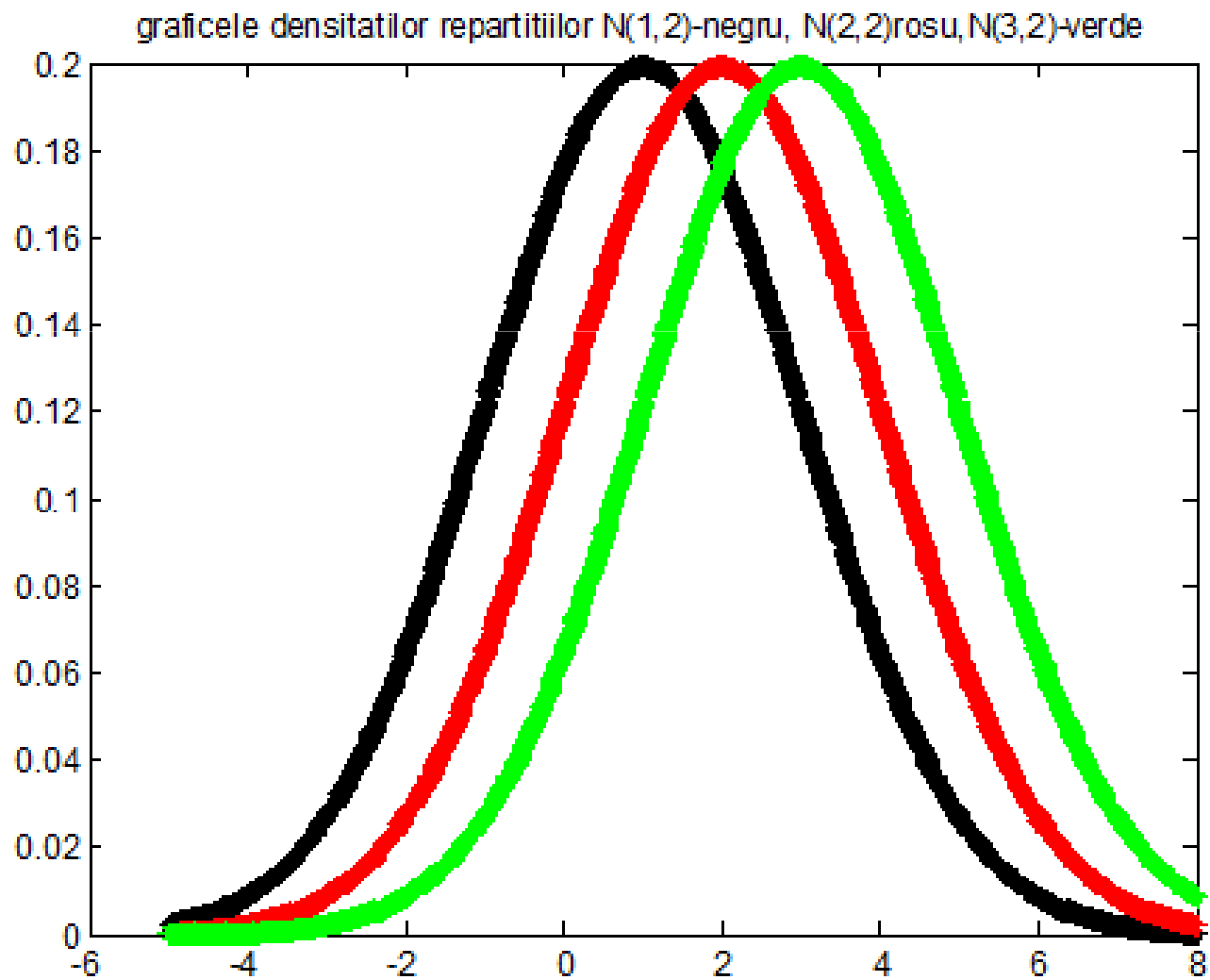
σ^2 este dispersia și σ deviația standard.


```
>>x=-5:.01:8;plot(x,1/(0.5*sqrt(2*pi))*exp(-(x-3).^2/(2*0.5^2)),'k',  
x,1/sqrt(2*pi)*exp(-(x-3).^2/2),'r',x,1/(2*sqrt(2*pi))*exp(-(x-3).^2/(2*2^2)),'g')
```

graficele densitatiiilor repartiilor $N(3, 0.5)$ -negru, $N(3, 1)$ -rosu, $N(3, 2)$ -verde



```
>>x=-5:.01:8;plot(x,1/(2*sqrt(2*pi))*exp(-(x-1).^2/(2*2^2)), 'k',  
x,1/(2*sqrt(2*pi))*exp(-(x-2).^2/(2*2^2)), 'r',  
x,1/(2*sqrt(2*pi))*exp(-(x-3).^2/(2*2^2)), 'g')
```

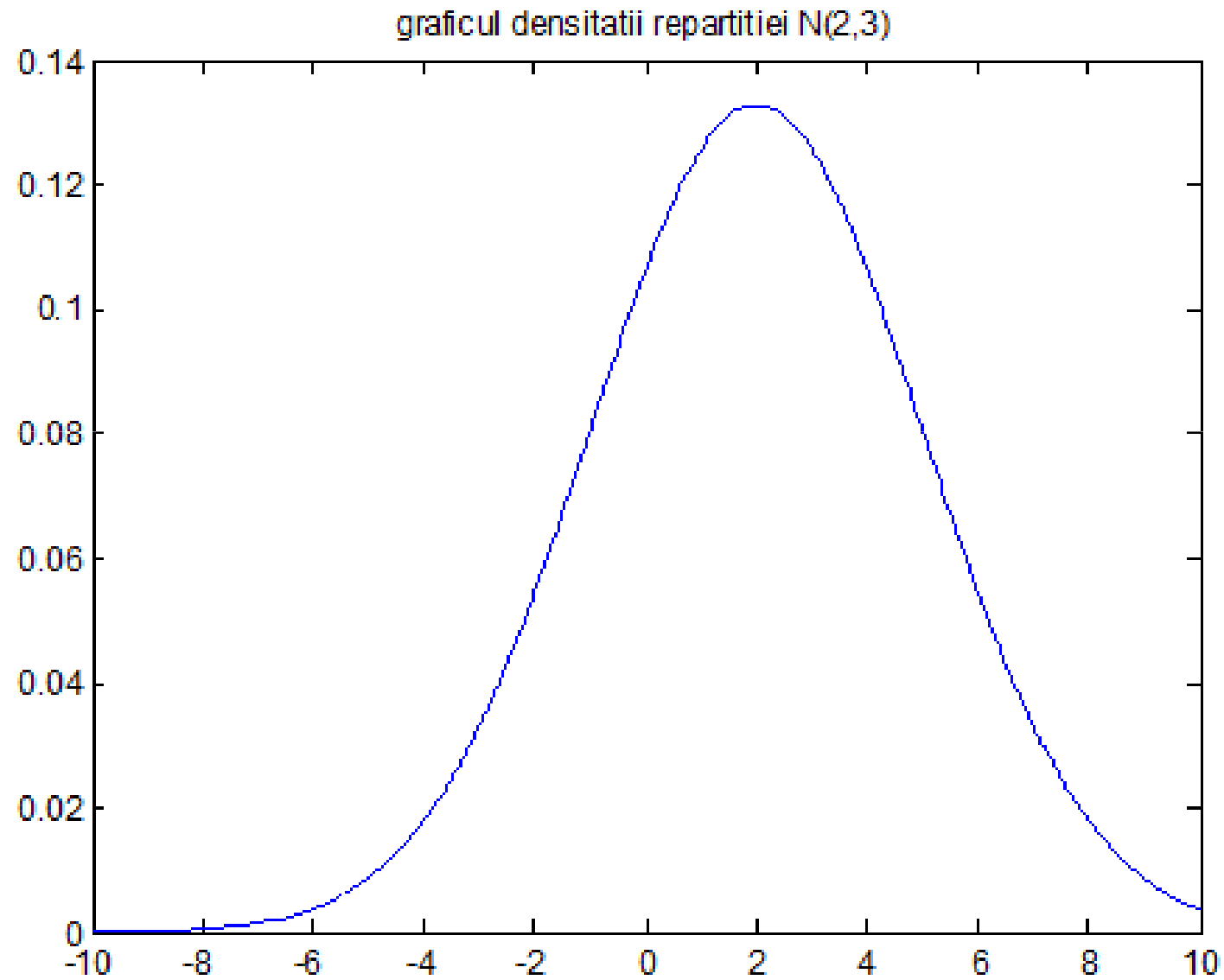


Pentru a desena în Matlab graficul densității repartiției normale se folosește sintaxa

```
Y = normpdf(X,mu,sigma)
```

care calculează valorile în X a densității de repartiție (pdf) pentru distribuția normală de medie μ și deviație standard σ .

```
>> X=-10:0.01:10;Y=normpdf(X,2,3);plot(X,Y)
```



Generarea de numere aleatoare normal repartizate - Matlab

`R = normrnd(mu,sigma,m,n)` generează aleator numere normal distribuite de parametri `mu` (media) și `sigma` (dispersia), scalarii `m` și `n` reprezentând numărul de linii, respectiv coloane a lui `R`.

```
>> R = normrnd(2,3,1,8)
```

```
R =
```

```
 4.6652 -1.4412 -1.2066 -0.4285 -6.8329  6.3151  2.9756 -0.2648
```

```
>> R = normrnd(2,3,1,8)
```

```
R =
```

```
 6.1109 -3.1345  1.6933  1.2757  2.9576  2.9386 -0.5946  1.9098
```

```
>> R = normrnd(3,5,1,9)
```

```
R =
```

```
 2.1756  6.1385  8.4663  8.5464 -1.3183  3.3868 -3.0706 -2.5675
```

```
2.9658
```

```
>> R = normrnd(3,5,1,9)
```

```
R =
```

```
10.6632 -0.8483  4.8569  1.8721  8.5868 -2.4453  3.1628  5.7626
```

```
8.5031
```

Graficul probabilității normale- Matlab

Graficul probabilității normale este o tehnică care testează vizual dacă un set de date are sau nu o repartiție normală.

`h = normplot(X)` afișează un grafic al datelor din X .

Dacă X este o matrice, `normplot` afișează câte un grafic pentru fiecare coloană a lui X .

Graficul prezintă datele eșantionului notate cu '+'. Peste acest grafic este suprapusă o linie ce unește prima și a treia cuartilă a lui X , linie ce este extrapolată la tot eșantionul.

Scopul acestei reprezentări este de a vizualiza dacă datele din X sunt normal repartizate. Dacă da, graficul va fi liniar. Alt tip de repartiție face ca datele să fie pe o anumită curbă.

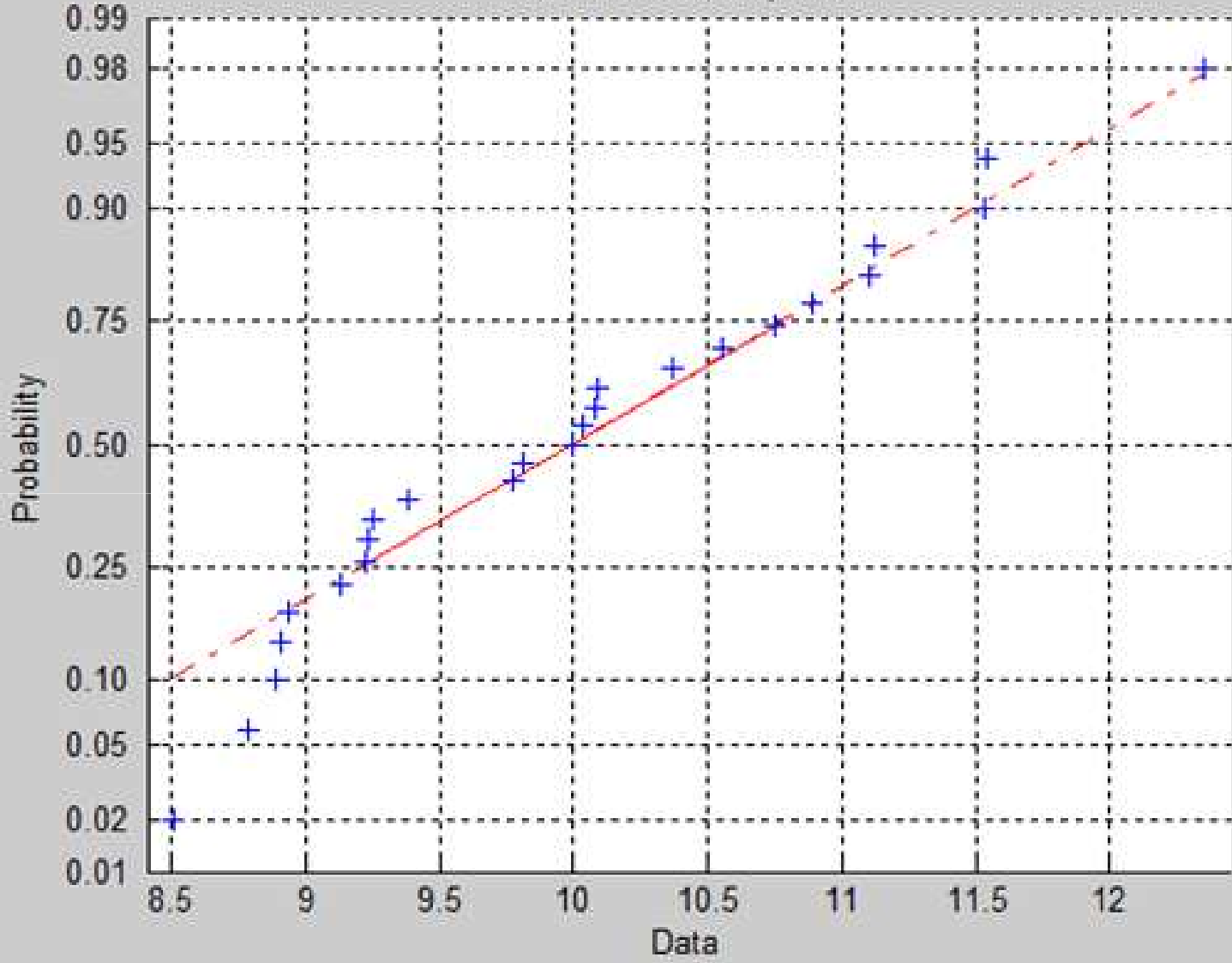
Dacă avem date cenzurate se utilizează `probplot`

15. Exemplu

— Să generăm un eșantion de 25 de numere aleatoare, normal repartizate $N(25,1)$ și să desenăm graficul probabilității normale, corespunzător acestui eșantion

```
>>x = normrnd(10,1,1,25);  
>>normplot(x)
```

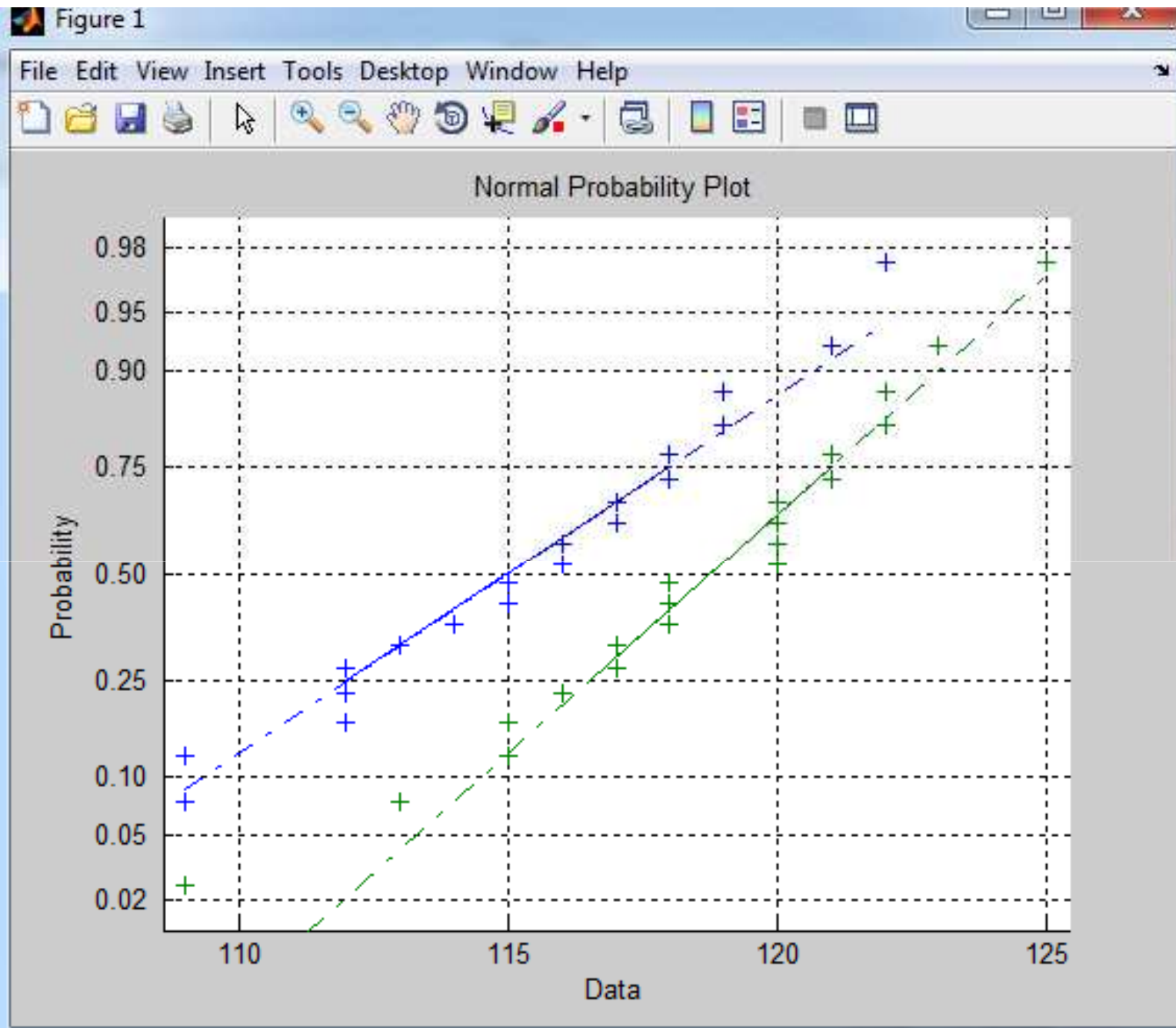
Normal Probability Plot



16. Exemplu

în Matlab bază de date `gas.mat` conține două eșantioane cu prețurile galonului de benzină într-un număr de stații, aleator alese din statul Massachusetts, în 1993. Primul eșantion `price1`, conține 20 de observații aleatoare din aceeași zi de ianuarie, stațiile fiind pe tot cuprinsul statului, în timp ce al doilea eșantion `price2`, conține 20 de observații aleatoare din aceeași zi, o lună mai târziu, stațiile fiind pe tot cuprinsul statului.

```
>> normplot(x)
>> load gas
>> prices = [price1 price2]
>> normplot(prices)
```



Repartiția log-normală

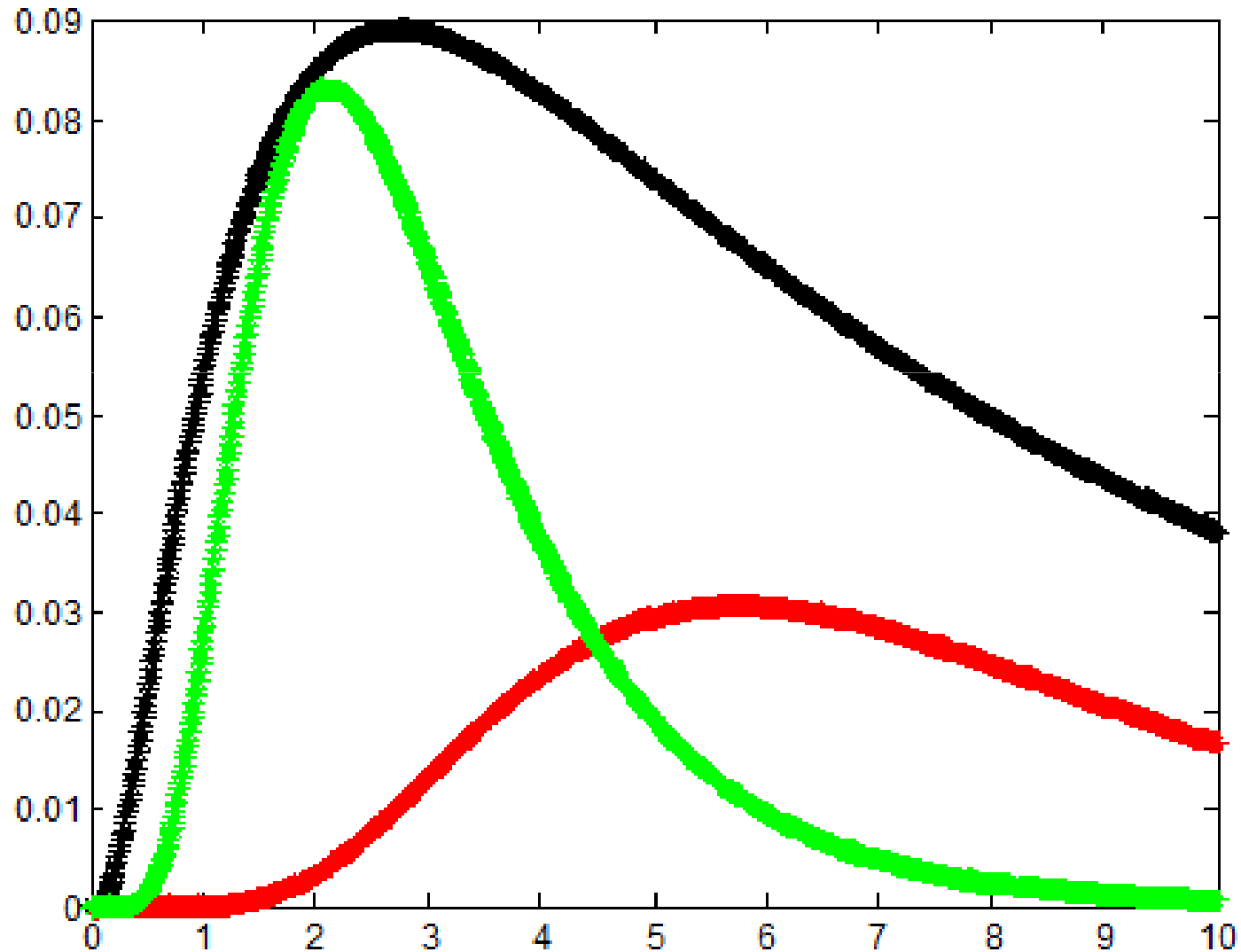
Repartiția *log-normală* $X \sim \text{LOGN}(\mu, \sigma^2)$: variabila $X \geq 0$ este *log-normal repartizată* dacă variabila $\ln(X)$ să fie normal repartizată.

Densitatea de repartiție este:

$$f_X(x) = \frac{1}{\sqrt{2\pi\sigma x}} \exp\left[-\frac{(\ln x - \mu)^2}{2\sigma^2}\right], \quad x > 0,$$

```
>> x=0:.01:10; plot(x,(1./(sqrt(2*pi)*x)).*exp(-(log(x)-2).^2/2),'k',  
x,(1./(sqrt(2*pi)*2*x)).*exp(-(log(x)-2).^2/2^2),'r',  
x,(1./(sqrt(2*pi)*2*x)).*exp(-(log(x)-1).^2/2^2),'g')
```

graficele densitatilor repartitiilor LOGN(2,1)-negru, LOGN(2,4)-rosu, LOGN(1,4)-verde

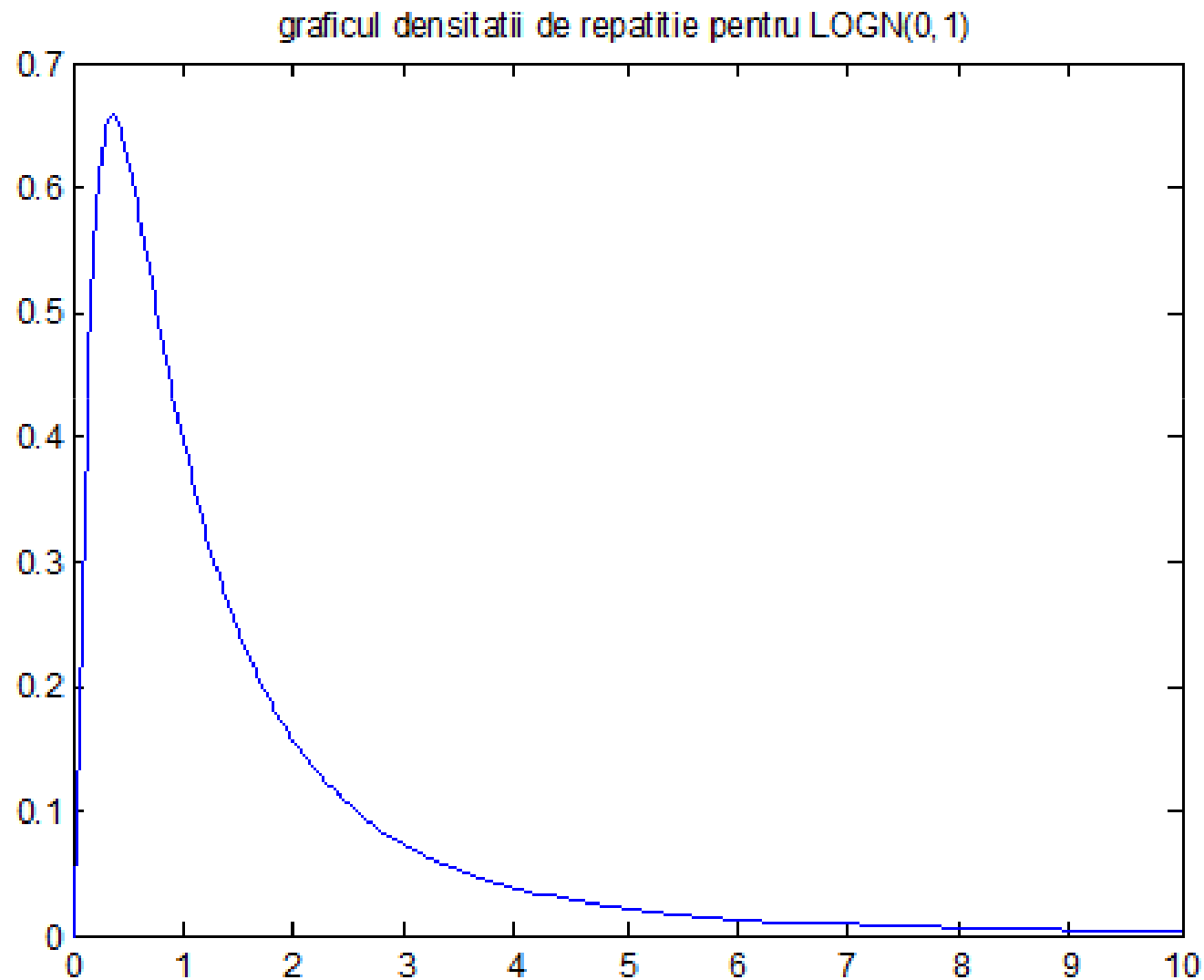


Graficul densității de repartiție (*probability density function* – pdf) în cazul repartiției lognormale poate fi desenat folosind sintaxa

`Y = lognpdf(X, mu, sigma)`

care returnează valorile în X ale densității de repartiție lognormale de parametri μ și σ . Reamintim că μ și σ reprezintă media și deviația standard a repartiției normale asociate X .

```
>> x = (0:0.02:10);  
>> y = lognpdf(x,0,1);  
>> plot(x,y)
```



Prezentăm formulele de calcul ale mediei și dispersiei variabilei lognormale X :

$$m = e^{\frac{\mu + \sigma^2}{2}} ;$$

$$D = e^{(2\mu + \sigma^2)(e^{\sigma^2} - 1)}$$

O repartiție lognormală de medie m și dispersie are ca parametri:

$$\mu = \ln \frac{m^2}{\sqrt{D + m^2}}$$

$$\sigma = \sqrt{\ln \left(\frac{v}{m^2} + 1 \right)}$$

17. Exemplu

Pentru a calcula media și dispersia unei variabile lognormale în Matlab folosim funcția `lognstat`.

Să generăm 15 numere aleatoare lognormal repartizate având media 1 și dispersia 2:

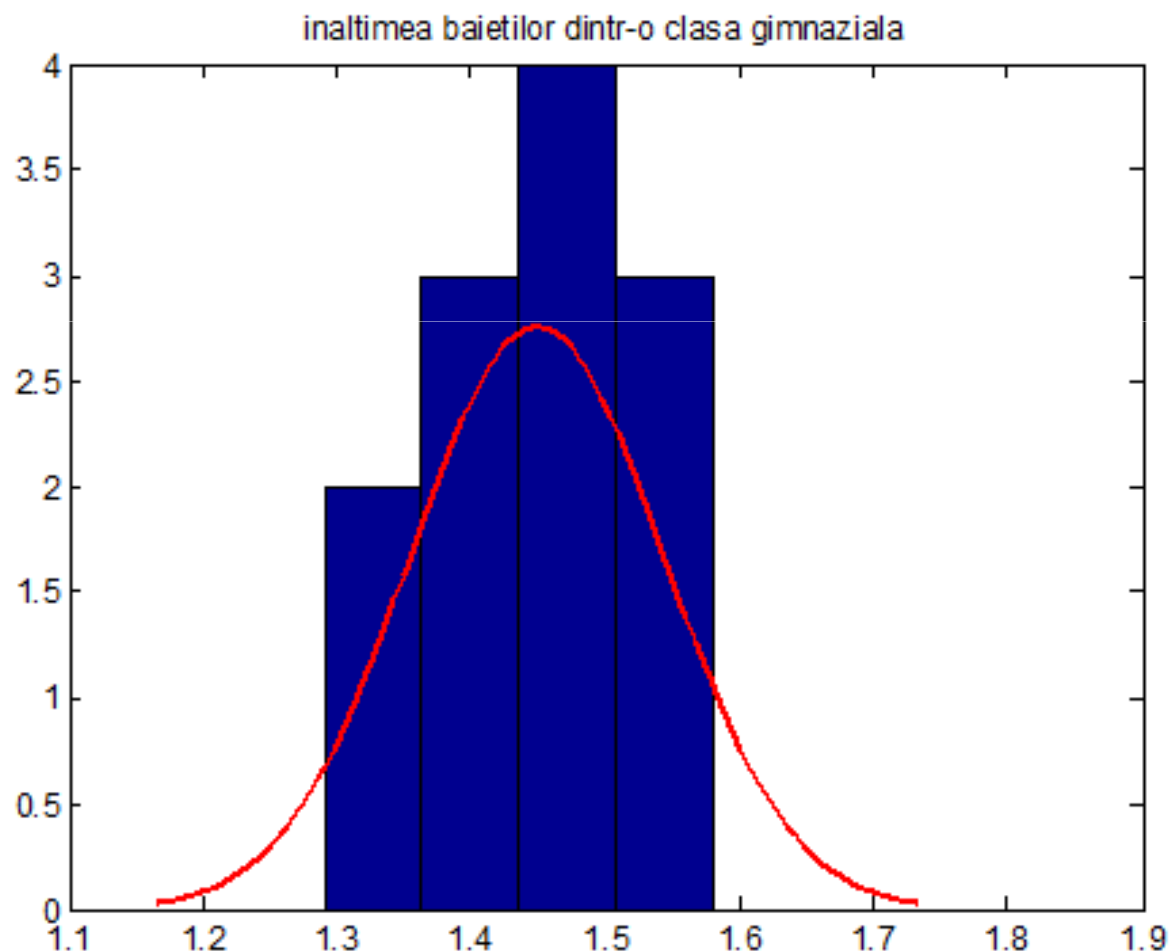
```
>> m = 1;v = 2;mu = log((m^2)/sqrt(v+m^2))
mu =
    -0.5493
>> sigma = sqrt(log(v/(m^2)+1))
sigma =
    1.0481

>> X = lognrnd(mu, sigma, 1, 15)
X =
Columns 1 through 9
    0.5187    0.4483    0.8068    0.8014    0.2332    0.5594    0.4857    1.1147
    1.8159
Columns 10 through 15
    1.8467    0.2335    0.6261    0.1617    0.1797    0.5732
>> MX = mean(X);VX = var(X);[MX VX]
ans =
    0.6937    0.2809
```


histfit

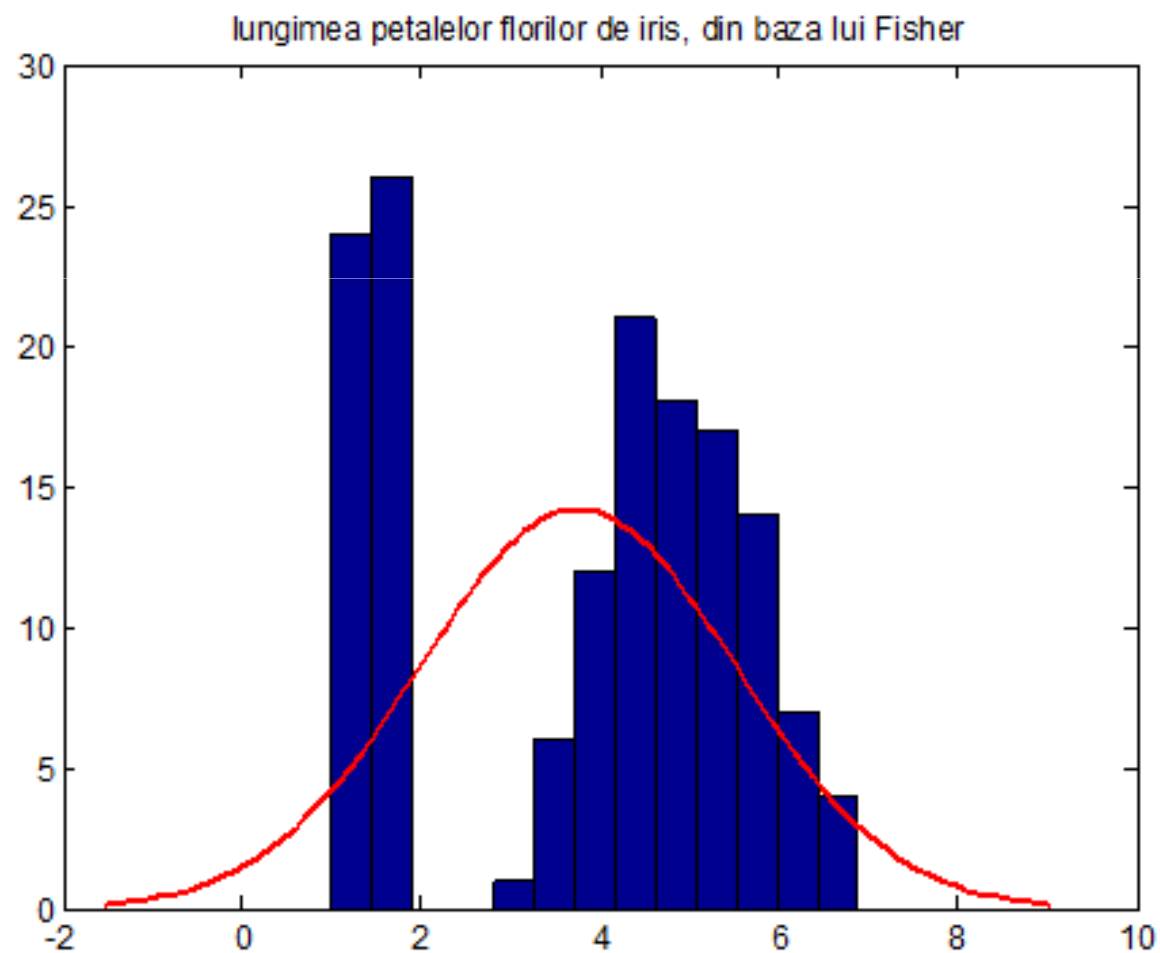
histfit(X) crează o histogramă a valorilor vectorului X, utilizând un număr de intervale egal cu $\lceil \sqrt{\dim X} \rceil + 1$, căreia îi suprapune o densitate de repartiție normală. Reluăm exemplul nr 5:

```
>> X1=[1.40, 1.37, 1.57, 1.46, 1.49, 1.46, 1.39, 1.55, 1.29, 1.58, 1.50,  
1.33]; histfit(X1)
```



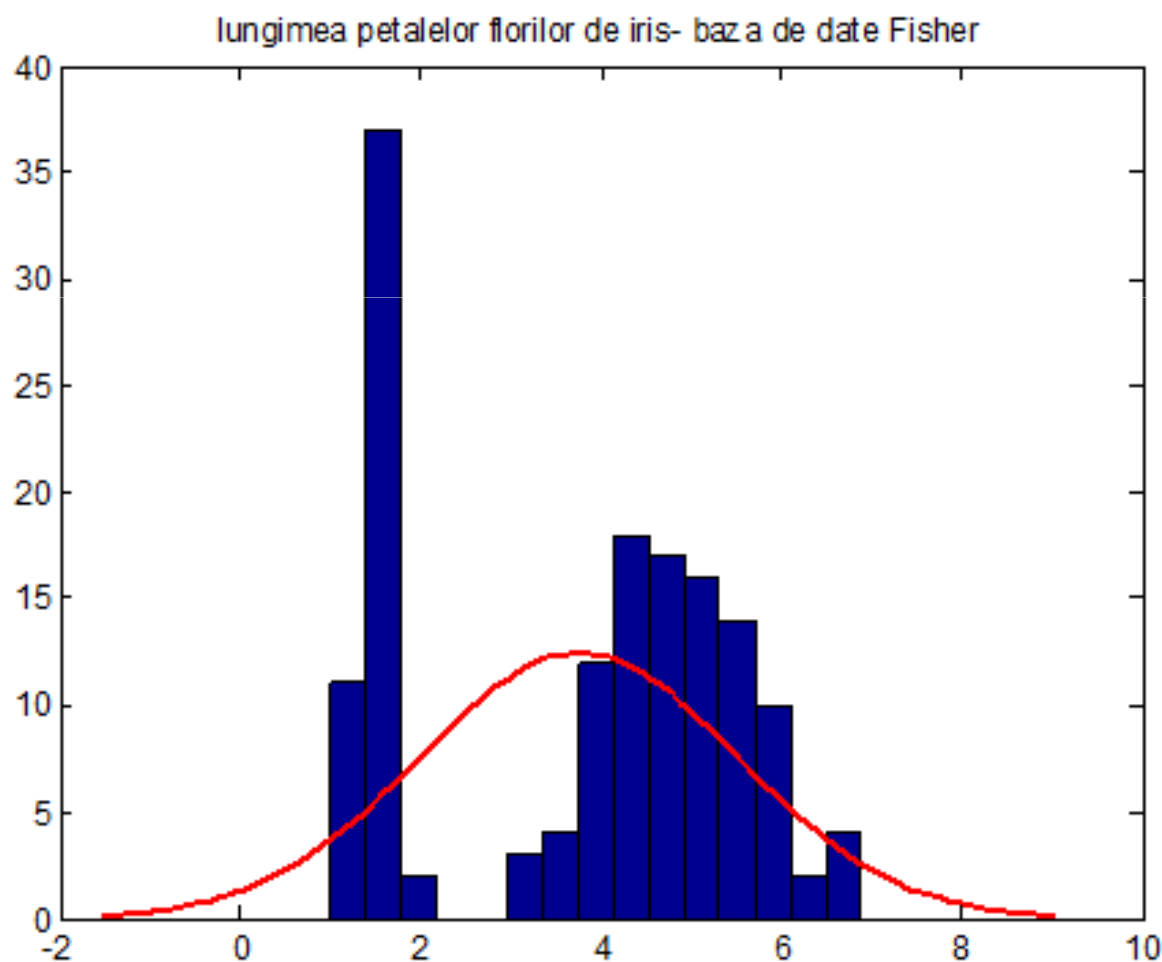
Exemplificăm cu vectorul ce reprezintă lungimea petalelor florilor de iris (Fisher)

```
>> load fisheriris  
>> x=meas(:,3); histfit(x)
```



histfit(X,n) crează o histogramă a valorilor vectorului X, utilizând n intervale, căreia îi suprapune o densitate de repartiție normală. Exemplificăm cu vectorul ce reprezintă lungimea petalelor florilor de iris (Fisher)

```
>> load fisheriris  
>> x=meas(:,3); histfit(x,15)
```



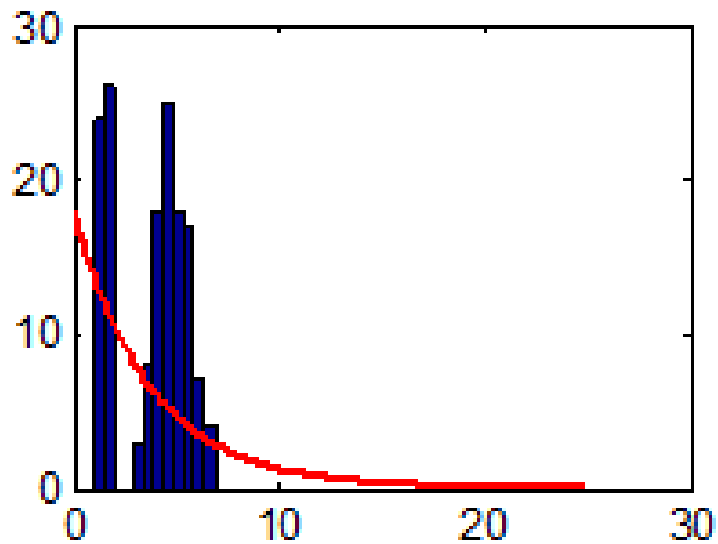
`histfit(X,n,'repartitie')` crează o histogramă a valorilor vectorului `X`, utilizând `n` intervale, căreia îi suprapune o densitate de repartiție impusă:

- 'exponential'
- 'gamma'
- 'lognormal'
- 'normal'
- 'weibull' sau 'wbl'

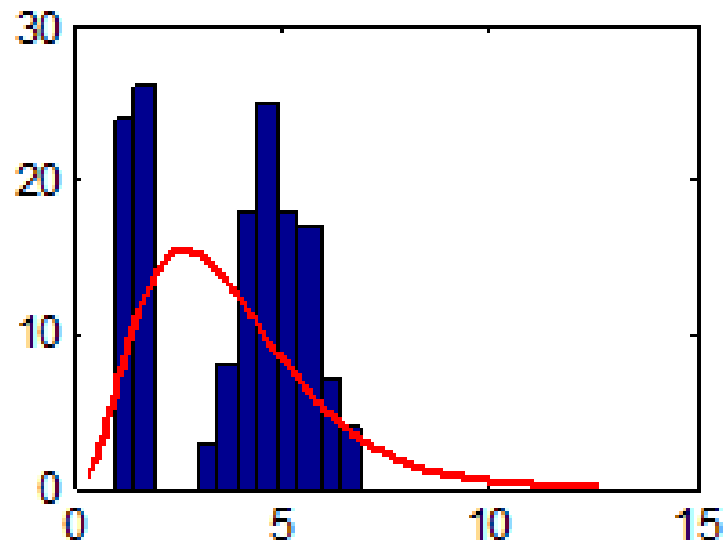
Prentăm histograma lungimii petalelor irișilor din baza de date Fisher, suprapunând câte o densitate de repartiție:

```
>> load fisheriris  
>> x=meas(:,3);  
>> subplot(221);histfit(x,12,'exponential')  
>> subplot(222);histfit(x,12,'gamma')  
>> subplot(223);histfit(x,12,'lognormal')  
>> subplot(224);histfit(x,12,'wbl')
```

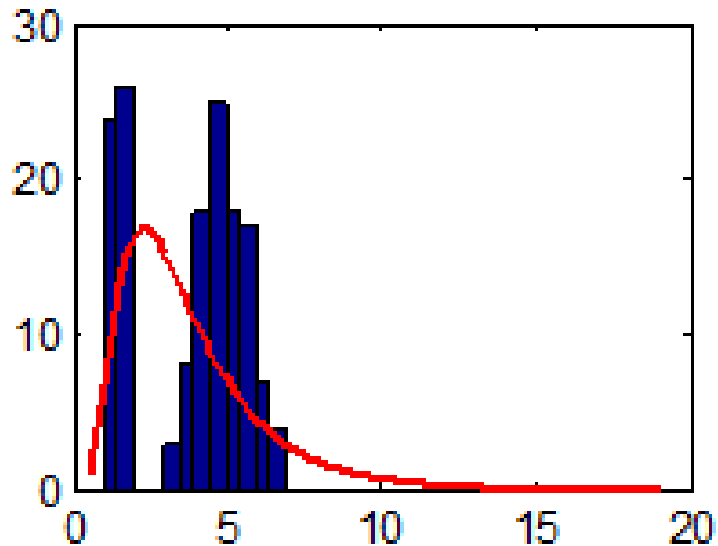
pdf exponential



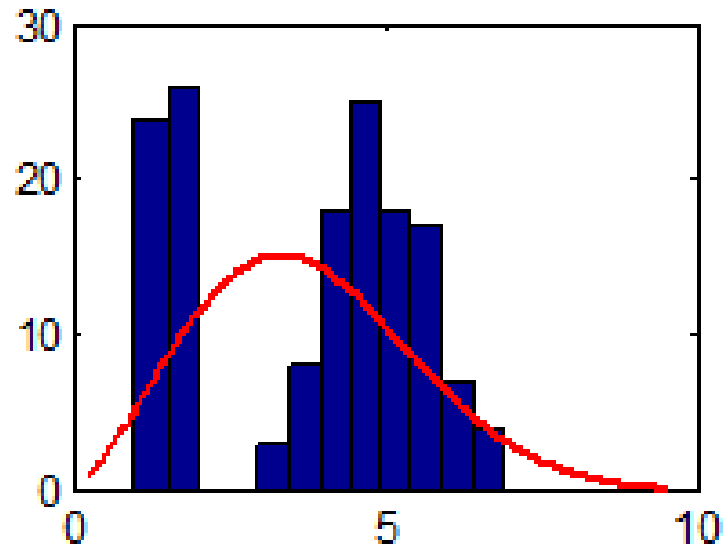
pdf gamma



pdf lognormal



pdf Weibull

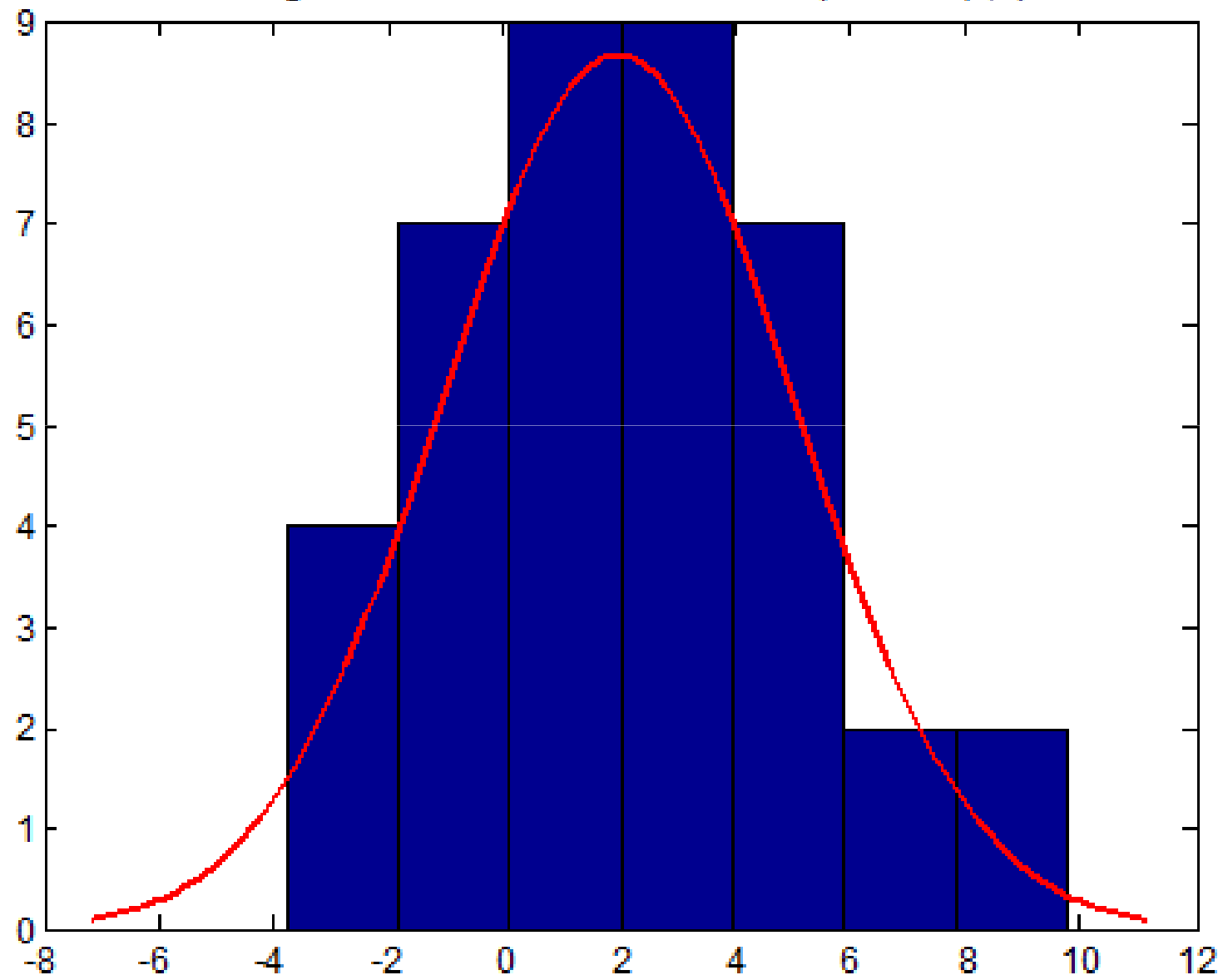


18. Exemplu

Generăm aleator 40 de numere, cu repartiție normală $N(2,3)$, desenăm histograma lor și suprapunem graficul densității de repartiție normale:

```
>> R = normrnd(2,3,1,40);  
>> histfit(R)
```

histograma a 40 numere aleatoare cu repartitie $N(2,3)$



Teste statistice

Plecând de la realitatea înconjurătoare, formulăm diferite ipoteze pe baza informațiilor primite, ipoteze trebuiesc verificate (testate). După testare, care este un proces complex, ipoteza respectivă devine ipoteză *adevărată*.

Din punct de vedere statistic ipotezele fac referire la fenomene și relații ce privesc populațiile statistice. Este vorba de ipoteze privind repartiții de probabilitate, valoarea unor parametri statistici sau legăturile dintre aceștia, plecând de la observațiile făcute pe un eșantion de populație.

Există două tipuri de teste statistice:

- testele parametrice se referă la ipoteze statistice ce privesc parametrii statistici (medie, dispersie);
- testele non-parametrice se folosesc pentru studiul datelor care nu sunt guvernate de repartiția normală, aplicate mai ales la date care sunt guvernate de legi de probabilitate asimetrice.

Formularea și testarea ipotezelor este o parte esențială a inferenței statistice dinspre eșantion spre populația generală.

În scopul testării se formulează o ipoteză de studiu pe care trebuie să o demonstrăm. Poate fi vorba despre un raționament care este considerat corect sau care poate fi folosit ca bază de comparație.

De exemplu o ipoteză poate fi:

un tratament nou este mai bun decât tratamentul folosit curent pentru aceeași afecțiune.

O testare statistică se bazează pe confruntarea a două ipoteze:

- ipoteza nulă H_0 care afirmă că nu există nicio diferență semnificativă între datele comparate, diferența obținută datorându-se doar întâmplării.
 - Când comparăm media de sondaj cu media cunoscută a populației, H_0 afirmă că nu există diferență semnificativă între cele două valori, deci eșantionul este reprezentativ pentru populația originală.
 - În cazul a două eșantioane, în contextul comparației unor anumiți parametri, H_0 afirmă că cele două eșantioane aparțin aceleiași populații.

- ipoteza alternativă H_1 care afirmă că efectul comparației este nenul, existând diferențe semnificative, nedatorate hazardului privind datele considerate.

De exemplu, în studiul unui nou medicament se spune:

H_0 = noul medicament nu este mai eficient, în medie, decât medicamentul curent.

Trebuie acordată o mare importanță ipotezei nule pentru că este în strânsă relație cu ipoteza alternativă, care trebuie confirmată.

Neacceptarea ipotezei nule, confirmă aplicarea ipotezei alternative H_1 .care este în fapt scopul studiului:

H_1 : noul medicament este mai eficient, în medie, decât medicamentul curent.

Concluzia finală, după aplicarea unui test de semnificație se exprimă întotdeauna în termenii ipotezei nule:

Respingem H_0 , caz ce sugerează că ipoteza alternativă poate fi adevărată

Nu respingem H_0 , caz ce sugerează că nu există destule argumente pentru respingerea ipotezei nule.

Ipoteza alternativă este o afirmație pe care aplicarea testului de semnificație statistică dorește să o confirme, prin neinfirmare.

Formularea acestei ipoteze poate lua trei forme:

$H_1 \neq H_0$, de exemplu: efectul noului medicament este diferit de cel curent

$H_1 > H_0$, de exemplu: efectul noului medicament este mai bun decât cel curent

$H_1 < H_0$, de exemplu: efectul noului medicament este mai prost decât cel curent

Intrebări:

- Care este riscul de a respinge ipoteza nulă, ea fiind totuși adevărată?
- Care este riscul de a o accepta , ipoteza nulă dovedindu-se în final falsă?

Pot fi făcute două tipuri de erori, noțiunile statistice următoare precizând acest fapt:

- Eroarea (riscul) de primă speță, notată α , este probabilitatea ca să respingem ipoteza nulă, atunci când ea este adevărată; valoarea lui α este considerată uzual 5% și este aleasă înainte de a începe analiza. Decizia este greșită, fiind un rezultat *fals pozitiv*.
- Eroarea (riscul) de speța a doua, notată β , se refera la probabilitatea de a accepta ipoteza nulă, atunci când ea este falsă; valoarea lui β depinde de diferența obținută în analiza cât și de volumul eșantionului. Decizia este greșită, fiind un rezultat *fals negativ*.

Puterea testului, notată cu π și egală cu $1 - \beta$ reprezintă probabilitatea de a respinge ipoteza nulă, atunci când este falsă. Un interval de încredere mare indică o putere scăzută.

Procedura de testare constă în maximizarea puterii testului, ceea ce înseamnă minimizarea erorii de speța a doua, atunci când eroarea de prima speță este limitată a priori arbitrar.

- se stabilește că ipoteza nulă este adevărată, adică nu există nicio diferență semnificativă în comparația făcută;
- interpretarea rezultatului obținut după procesarea datelor: Care este probabilitatea de a fi obținut rezultatul respectiv, dacă ipoteza nulă ar fi adevărată? Această probabilitate este cunoscută sub numele de *nivelul de semnificație p*.

Presupunem că s-au colectat date din două eșantioane și că mediile acestora sunt diferite. A observa că cele două medii sunt diferite nu este suficient pentru a concluziona că populațiile au medii diferite.

E posibil ca populațiile să aibă aceeași medie, iar diferența observată să fie o coincidență în urma eșantionării aleatorii.

Nu există nici un mod de a fi siguri că diferența observată reflectă o diferență reală sau este o consecință a eșantionării aleatorii.

Valoarea p răspunde la această întrebare: dacă populațiile au avut într-adevăr aceeași medie, care este probabilitatea de a observa o asemenea diferență (sau una mai mare) între mediile eșantioanelor într-un experiment de dimensiunea acestuia?

Dacă valoarea p este mică, concluzia este că diferența are șanse mici să fie provocată de eșantionarea aleatorie și se poate concluziona că populațiile au medii diferite.

Pentru detectarea semnificației statistice trebuie urmate etapele de inferență:

1. Stabilirea unei valori prag, înainte de a efectua experimentul; în mod tradițional valoarea - prag (numită α) este stabilită la 0,05 sau 0,01.
2. Definirea ipotezei nule.
3. Se aplică testul statistic aferent pentru a calcula valoarea p .
4. Se compară valoarea p cu valoarea – prag.
 - Dacă valoarea p este mai mică decât pragul stabilit, se declară că *se respinge ipoteza nulă* și că diferența *e semnificativă din punct de vedere statistic*.
 - Dacă valoarea p este mai mare decât pragul, se declară că *nu se respinge ipoteza nulă* și că diferența *nu e semnificativă din punct de vedere statistic*. Nu se poate concluzia că ipoteza nulă este adevărată. Se poate spune doar că nu avem probe suficiente pentru a respinge ipoteza nulă.

Testul de semnificație calculează valoare p , probabilitatea de a obține o valoare corectă a parametrului la nivel de populație, prin extrapolare (inferență) de la valoarea estimatorului (calculat pe baza eșantionului).

Testarea ipotezelor implică posibile erori:

Decizia statistică	H_0 real adevărată	H_0 real falsă
H_0 respinsă	Eroare tip I (α)	Corect
H_0 acceptată	Corect	Eroare tip II (β)

Eroarea de tip I este numită nivel de semnificație, valoarea ei fiind stabilită de cercetător.

Populațiile sunt identice, deci nu există de fapt nici o diferență. Din întâmplare s-au obținut valori mai mari în cazul unui grup și mai mici în cazul celuilalt.

Când apare un rezultat semnificativ din punct de vedere statistic – în cazul în care populațiile sunt identice ($H_0 = H_1$) înseamnă că a intervenit o eroare de tip I.

Dacă se definește $p < 0.05$ ca fiind semnificația statistică, va apare eroare de tip I în 5% din experimentele în care nu există nici o diferență.

Nivelul de semnificație p este o cuantificare a influenței hazardului asupra diferenței obținute prin comparație, eșantionul fiind presupus reprezentativ pentru populația originară.

$p = 0.05$ nivel statistic semnificativ la limită

$p = 0.01$ nivel statistic semnificativ

$p = 0.005$ nivel înalt semnificativ

Reamintim că:

Regiunea critică reprezintă mulțimea valorilor x ale variabilei statistice corespunzătoare X pentru care se respinge ipoteza nulă.

Regiunea de încredere a unui parametru este mulțimea valorilor θ_0 pentru care se acceptă ipoteza nulă.

Problema practică:

Cum putem calcula probabilitatea P de a obține efectul respectiv, plecând de la datele eșantionului, presupunând ipoteza nulă adevărată?

Se calculează testul statistic, care este o funcție de variabilele de sondaj, plecând de la o anumită cantitate calculată pe baza sondajului și comparată cu o valoare ipotetică, propusă, presupunând că ipoteza nulă este adevărată.

$$\text{test statistic} = \frac{\text{valoare de sondaj} - \text{valoare ipotetică}}{\frac{\text{deviatia st.}}{\sqrt{n}}}$$

unde prin valoare de sondaj înțelegem estimatorul parametrului căutat.

Algoritm de lucru

1. Se formulează setul de ipoteze H_0, H_1 .
2. Se calculează dintr-un eșantion o statistică (test statistic)
3. Dacă în lumea reală are loc un eveniment, se calculează în ipoteza H_0 , probabilitatea p de apariție a valorii calculate.
4. Dacă p este mică apare o contradicție, pentru rezolvarea căreia se va respinge H_0 în favoarea lui H_1 (p este mică deoarece în calculul acesteia s-a acceptat ipoteza H_0).
5. Dacă p este mare nu se respinge H_0 , neexistând nici un motiv pentru a lua decizia contrară.

Rămâne o singură întrebare: începând de unde o probabilitate este Considerată drept “mică”? Pentru a nu introduce subiectivismul în această decizie, se fixează, anterior deciziei în test, un prag sub care o probabilitate este considerată “mică”.

Această valoare se numește prag de semnificație și se notează uzual cu α .

Regula de decizie în test poate fi formulată atunci:

- dacă $P \leq \alpha$, atunci se respinge ipoteza nulă, H_0 , în favoarea ipotezei alternative, H_1 ;
- dacă $P > \alpha$, atunci nu se respinge ipoteza nulă H_0 .

Reamintim:Repartiția t-Student - W.S Gosset alias *student* 1908

Repartiția Student (*t-distribution*) este o familie de curbe ce depind de un singur parametru n , care reprezintă gradele de libertate. Dacă $n \rightarrow \infty$, atunci repartiția t se apropie de repartiția normală standard

Dacă x este un eșantion de dimensiune n dintr-o repartiție normală de medie ,

μ , atunci testul statistic $t = \frac{\bar{X} - \mu}{\frac{\bar{\sigma}}{\sqrt{n}}}$ are o repartiție Student cu $n-1$ grade de

libertate.

Funcția de repartiție Student este:

$$f(t) = \frac{\Gamma\left(\frac{n+1}{2}\right)}{\sqrt{n\pi} \cdot \Gamma\left(\frac{n}{2}\right)} \cdot \left(1 + \frac{t^2}{n}\right)^{-\frac{n+1}{2}},$$

unde $\Gamma(n) = \int_0^{\infty} t^{n-1} \cdot e^{-t} dt$, (funcția gamma)

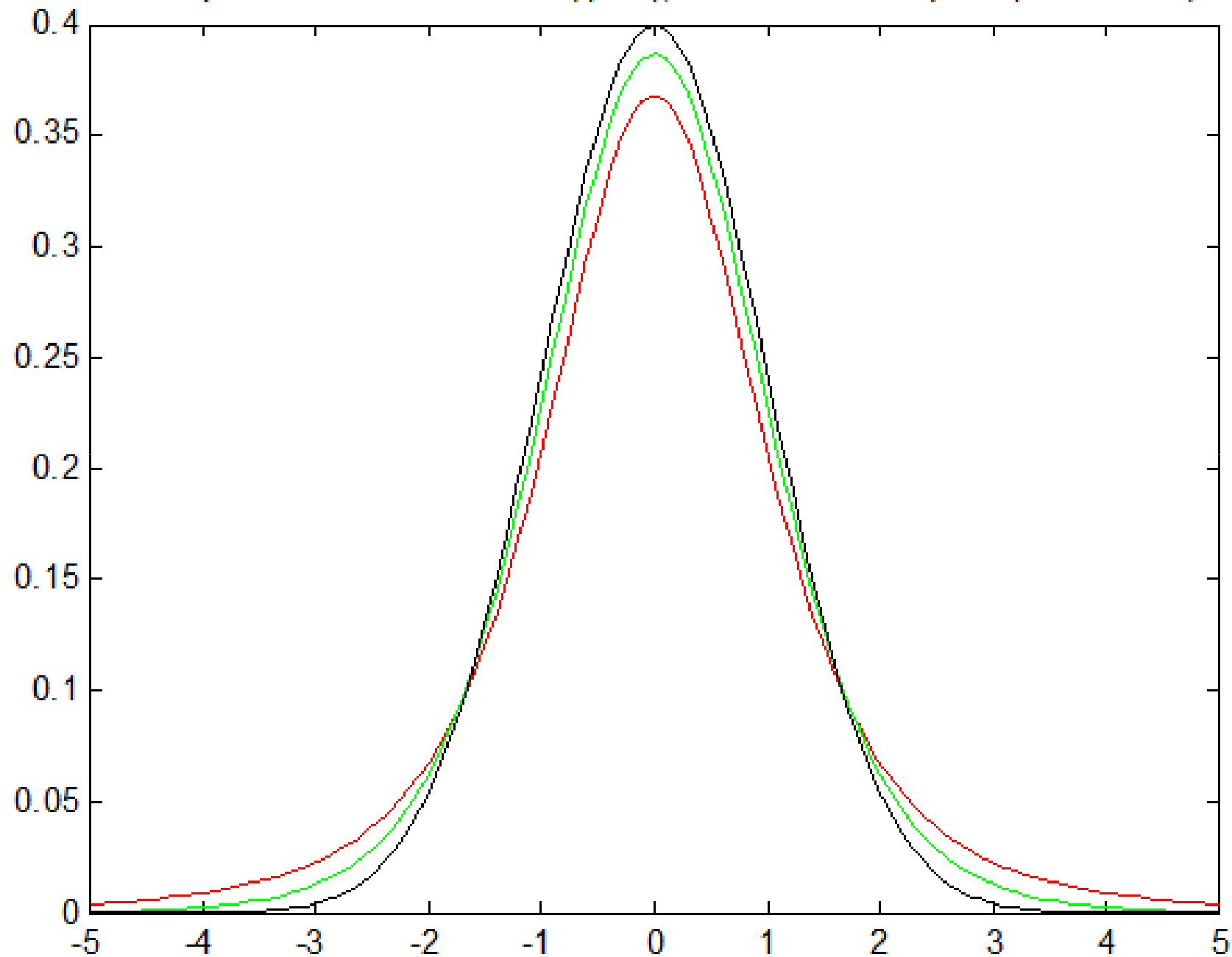
19.Exemplu

Vom desena în același cadran funcția de repartiție Student, pentru $n = 3$, respectiv pentru $n = 8$ și funcția de repartiție normală.

$y = \text{tpdf}(x,n)$ calculează valoarea funcției de repartiție Student, în fiecare componentă a vectorului x , folosind numărul de grade de libertate date de n .

```
>> x = -5:0.1:5;  
>> y1 = tpdf(x,3);  
>> y2 = tpdf(x,8);  
>> z = normpdf(x,0,1);  
>> plot(x,y1,'r',x,y2,'g',x,z,'k')
```

functii de repartitie: Student cu $n=3$ (rosu), Student cu $n=8$ (verde) i normala (negru)



20. Exemplu

În cazul testării dacă media populației originare are o anumită valoare ipotetică, propusă într-un studiu statistic, testul statistic este:

$$t = \frac{\bar{X} - \mu}{\frac{\bar{\sigma}}{\sqrt{n}}}, \text{ unde:}$$

- \bar{X} este media de sondaj,
- μ este valoarea ipotetică a mediei populației;
- $\bar{\sigma}$ este deviația standard.

Concret, se alege se alege ca medie ipotetică propusă o anumită valoare μ_0 , presupunându-se adevărată ipoteza nulă $H_0 : \mu = \mu_0$, ceea ce implică faptul că testul statistic urmează o lege Student cu un anumit număr de grade de libertate. Plecând de la valoarea concretă a lui t , se obține valoarea corespunzătoare a lui P , din tabelul repartiției Student.

Concret, dacă $\bar{x} = 6753.6$, $\mu_0 = 7725$, $\bar{\sigma} = 1142.1$, $n = 11$, obținem $t = 2.821$.

Prezentăm o parte dintr-un tabel corespunzător unei repartiții Student cu grade de libertate între 1 și 14.

<i>One Sided</i>	75%	80%	85%	90%	95%	97.5%	99%	99.5%	99.75%	99.9%	99.95%
<i>Two Sided</i>	50%	60%	70%	80%	90%	95%	98%	99%	99.5%	99.8%	99.9%
1	1.000	1.376	1.963	3.078	6.314	12.71	31.82	63.66	127.3	318.3	636.6
2	0.816	1.061	1.386	1.886	2.920	4.303	6.965	9.925	14.09	22.33	31.60
3	0.765	0.978	1.250	1.638	2.353	3.182	4.541	5.841	7.453	10.21	12.92
4	0.741	0.941	1.190	1.533	2.132	2.776	3.747	4.604	5.598	7.173	8.610
5	0.727	0.920	1.156	1.476	2.015	2.571	3.365	4.032	4.773	5.893	6.869
6	0.718	0.906	1.134	1.440	1.943	2.447	3.143	3.707	4.317	5.208	5.959
7	0.711	0.896	1.119	1.415	1.895	2.365	2.998	3.499	4.029	4.785	5.408
8	0.706	0.889	1.108	1.397	1.860	2.306	2.896	3.355	3.833	4.501	5.041
9	0.703	0.883	1.100	1.383	1.833	2.262	2.821	3.250	3.690	4.297	4.781
10	0.700	0.879	1.093	1.372	1.812	2.228	2.764	3.169	3.581	4.144	4.587
11	0.697	0.876	1.088	1.363	1.796	2.201	2.718	3.106	3.497	4.025	4.437
12	0.695	0.873	1.083	1.356	1.782	2.179	2.681	3.055	3.428	3.930	4.318
13	0.694	0.870	1.079	1.350	1.771	2.160	2.650	3.012	3.372	3.852	4.221
14	0.692	0.868	1.076	1.345	1.761	2.145	2.624	2.977	3.326	3.787	4.140

În cazul nostru având 10 grade de libertate, valoarea lui t corespunde la un nivel de semnificație (two tailed) între 98% și 99%.

Dintr-un tabel complet am calculat $P = 0.018$, deci diferența obținută este semnificativă statistic.

Alegerea numărului de grade de libertate se bazează pe faptul că ipoteza H_0 este adevărată.

Odată testul statistic calculat, alegerea lui P se poate face în două moduri: *one-tailed* sau *two-tailed*, în funcție de coada sau cozile repartițiilor de probabilitate ce apar pentru valorile extreme.

Testul chi pătrat de independență - test neparametric

Reamintim că dacă X_1, \dots, X_n sunt variabile gaussiene standard, atunci variabila $X^2 = X_1^2 + \dots + X_n^2$ urmează legea χ^2 cu n grade de libertate. Densitatea sa de repartiție este

$$f_n(x) = \frac{2^{1-\frac{n}{2}} \cdot x^{n-1} \cdot e^{-\frac{x^2}{2}}}{\Gamma\left(\frac{n}{2}\right)}$$

Tabelul de contingență reprezintă o clasificare a datelor în funcție de 2 criterii în cadrul carora datele sunt în continuare divizate în 2 sau mai multe categorii discrete și mutual exclusive.

De exemplu:

Considerăm o populație statistică formată din copii, împărțiți pe sexe: fete și băieți și vom studia care este preferința acestora față de trei tipuri de jucării A, B și C, pe baza analizei unui eșantion de copii.

Tabelul de contingență corespunzător este:

		Tipuri de jucării		
		A	B	C
Sex	Fete	50	20	15
	Băieți	0	30	54

Analiza acestui tabel se bazează pe testarea ipotezelor.

Ipoteza nulă se referă la presupunerea că nu există nicio relație semnificativă la nivelul populației de copii între cele două clasificări : sexul copiilor și tipurile de jucării, adică variabilele sunt considerate independente.

In acest caz vom compara frecvențele observate cu cele așteptate, dacă ipoteza nulă ar fi adevărată.

Notăm cu O_{ij} frecvențele observate și cu E_{ij} frecvențele așteptate dacă ipoteza nulă ar fi adevărată, unde i indică numărul liniei, iar j indică numărul coloanei din tablou.

Sub auspiciile ipotezei nule statistica
$$\chi^2 = \sum_{i=1}^p \sum_{j=1}^q \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

urmează o repartiție χ^2 cu $(p - 1) \cdot (q - 1)$ grade de libertate, unde p este numărul de linii și q este numărul de coloane din tabelul de contingență.

Prezentăm o parte din tabelul atașat testului χ^2 de independență, pâna la 16 grade de libertate.

n	Probability less than the critical value				
	0.90	0.95	0.975	0.99	0.999
1	2.706	3.841	5.024	6.635	10.828
2	4.605	5.991	7.378	9.210	13.816
3	6.251	7.815	9.348	11.345	16.266
4	7.779	9.488	11.143	13.277	18.467
5	9.236	11.070	12.833	15.086	20.515
6	10.645	12.592	14.449	16.812	22.458
7	12.017	14.067	16.013	18.475	24.322
8	13.362	15.507	17.535	20.090	26.125
9	14.684	16.919	19.023	21.666	27.877
10	15.987	18.307	20.483	23.209	29.588
11	17.275	19.675	21.920	24.725	31.264
12	18.549	21.026	23.337	26.217	32.910
13	19.812	22.362	24.736	27.688	34.528
14	21.064	23.685	26.119	29.141	36.123
15	22.307	24.996	27.488	30.578	37.697
16	23.542	26.296	28.845	32.000	39.252

21. Exemplu

Tabelul de contingență următor ilustrează legătura între statutul marital și consumul de cafea:

		Consum cafeină (mg/zi)				total
		0	1-150	151-300	>300	
Status marital	căsătorit	652	1537	598	242	3029 77.5%
	divorțat sau văduv	58	46	38	21	163 4.2%
	necăsătorit	218	327	106	67	718 18.3%
total		928	1910	742	330	3910

Din cei ce nu consumă cafea

- 928*77.5/100 ar trebui să fie căsătoriți,
- 928*4.2/100 ar trebui să fie divorțați sau văduvi
- 928*18.3/100 ar trebui să fie necăsătoriți

s.a.m.d

Obținem astfel tabelul cu frecvențele așteptate :

		Consum cofeină (mg/zi)			
		0	1-150	151-300	>300
Status marital	căsătorit	719	1480	575	256
	divorțat sau văduv	39	80	31	14
	necăsătorit	170	350	136	60
total		928	1910	742	330

Aplicând formula de mai sus obținem: $\chi^2 = 51.61$, ceea ce corespunde pentru repartiția χ^2 cu $3 \cdot 2 = 6$ grade de libertate, unui nivel de semnificație $P=0.00$, ceea ce înseamnă că există o asociere înalt semnificativă între statutul marital și consumul de cafea (variabile dependente).