

Modelle regresive

Modelele sunt reprezentări ale obiectelor sau ale situațiilor reale.

Există trei tipuri de modele:

- *modele iconice*, (iconic models) care sunt replici fizice ale obiectelor reale (de exemplu macheta unei construcții sau un camion jucărie).
- *modele analogice*, (analog models) care nu au aceeași înfățișare fizică cu obiectul ce urmează a fi modelat, dar în fond reprezintă același lucru (de exemplu, poziția acului de la vitezometrul unui automobil reprezintă viteza vehiculului)
- *modelele matematice* reprezintă problema printr-un sistem de simboluri și relații matematice.

Herbert A. Simon, deținător al premiului Nobel pentru economie, specialist în teoria deciziilor, spunea că modelele matematice nu trebuie să fie identice cu realitatea ci cât mai apropiate de aceasta și să conducă la rezultate mai bune decât cele obținute prin logica bunului simț

Un *model predictiv* creează, pe baza tuturor informațiilor (datelor) primite, un model statistic, care face o prognoză asupra unui rezultat viitor, cu o anumită acuratețe. Modelul, construit pe baza acelor factori variabili ce influențează rezultatul, cunoscuți sub numele de *predictori*, va fi validat sau revizuit pe baza unor date noi. În construcția modelului, avem de hotărât ce informații putem utiliza, care sunt variabilele asupra cărora putem face predicții, cum vom determina acuratețea modelului și, mai ales, dacă acesta poate fi folosit în lumea reală, deoarece acest tip de modelare este conceput ca să rezolve probleme reale, cu date reale.

Prognoza are drept scop producerea unei *ieșiri* (output) numerice corecte pentru o nouă *intrare* (input).

Să nu uităm însă vorbele lui Nils Bohr: „*Prognoza este foarte dificilă, mai ales dacă este vorba despre viitor*”, pentru a sublinia importanța validării unui model.

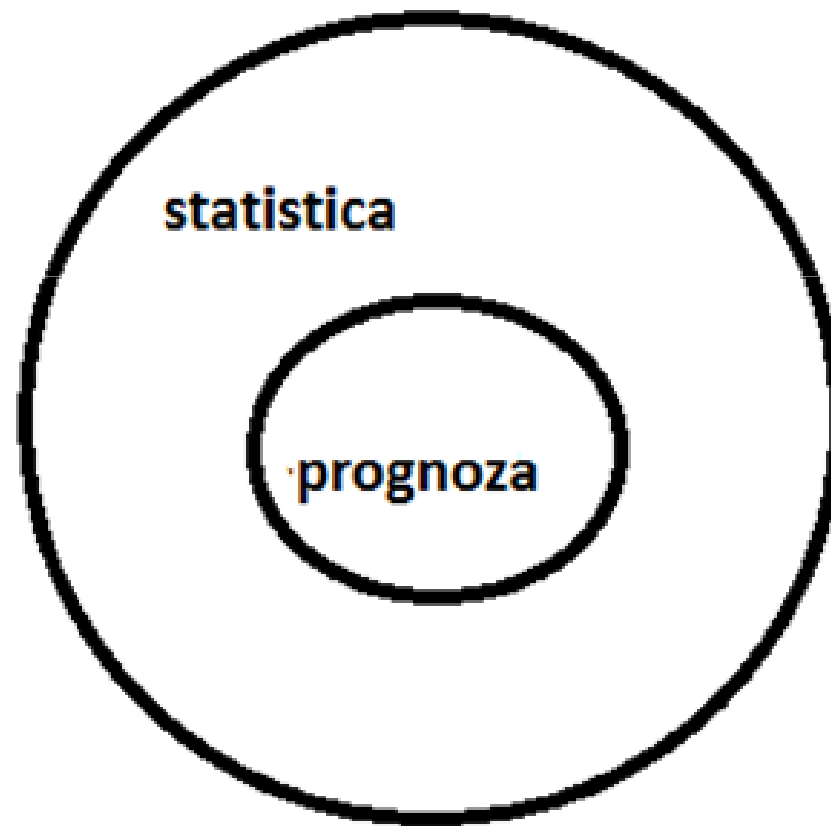
Nu este greu să construim un model care să se potrivească perfect cu datele existente, dar important este să construim acel model care identifică acele caracteristici din datele existente, ce se vor regăsi și în viitor.

Informaticianul și statisticianul au părerii diferite asupra prognozei

Punctul de vedere al informaticianului:



Punctul de vedere al statisticianului:



Informaticianul:

- Ce algoritm să utilizez sau să inventez?
- Care sunt proprietățile algoritmului (complexitate, timp de rulare)?

Statisticianul:

- Care ipoteze asupra datelor sunt convenabile?
- Care este cel mai bun rezultat pe care îl putem obține în aceste condiții?
- Cum să construiesc un predictor/un algoritm care să rezolve problema?

Informaticianul:

- Mulțimea de antrenament \Rightarrow Algoritm \Rightarrow Prognoză
- Proprietățile algoritmului

Statisticianul:

- Model: $(X_1, Y_1), \dots, (X_n, Y_n) \sim P \in \mathcal{P}$
- Notând cu $\hat{m}(X)$ prognoza, trebuie să gășesc \hat{m} optimal, adică pe cel care satiface condiția:

$$\sup_{P \in \mathcal{P}} E(Y - \hat{m}(X))^2 = \inf_{\hat{m}} \sup_{P \in \mathcal{P}} E(Y - \hat{m}(X))^2$$

- Sunt interesat de intervalul de încredere pentru $\hat{m}(X)$

Modele regresive

Modelele regresive sunt modele statistice care stabilesc relația între valorile a două sau mai multe variabile aleatoare, în vederea prognozării valorilor uneia în raport cu valorile celeilalte sau celorlalte.

Această problemă se pune atunci când între variabilele aleatoare considerate există o legătură consistentă, bazată pe natura intimă a fenomenelor care stau la baza lor.

Intuitiv, având o mulțime de valori experimentale din \mathbf{R}^2 , considerate a fi valorile a două variabile aleatoare, încercăm să determinăm curba care este „cea mai apropiată” de punctele respective.

Metoda cea mai des utilizată în acest scop este *metoda celor mai mici pătrate*. Dacă relația care stabilește legătura între variabila dependentă și variabilele independente este una liniară, se vorbește despre *regresia liniară*, altfel fiind vorba de *regresia neliniară* (polinomială, exponențială, logaritmică etc.).

Diagrama de împrăștiere

Prezentăm mecanismele cu care se pot evidenția legăturile între secvențe de date, provenind de la două sau mai multe variabile aleatoare.

În cazul a două serii statistice $\{x_i\}_{1 \leq i \leq n}$ și $\{y_i\}_{1 \leq i \leq n}$ definite pe același lot de obiecte putem considera seria cuplurilor de observații $\{(x_i, y_i)\}_{1 \leq i \leq n}$ definite de cele două variabile statistice pe același obiect i .

Acest cuplu de observații este reprezentat grafic folosind *norul* de puncte definit de reprezentarea bidimensională a punctelor (x_i, y_i) – așa numita *diagramă de împrăștiere*.

1. Exemplu

Într-o fabrică s-a constatat că viteza (m/min) benzii de lucru afectează numărul de defecțiuni descoperite în timpul verificării. Următorul tabel prezintă diferite viteze ale benzii și numărul defectelor găsite.

v	6	6	7	7	8	8	9	9	10	10	11	11	12	12	13	13
n	21	18	19	20	22	19	21	24	23	25	24	26	28	27	27	28

Prezentăm diagrama „norului” de împrăștiere a cuplurilor (viteza, număr defecțiuni) pentru lotul de 16 de observații:

```
>> X=[6 6 7 7 8 8 9 9 10 10 11 11 12 12 13 13];  
>> Y=[21 18 19 20 22 19 21 24 23 25 24 26 25 28 27 28];  
>> plot(X,Y,'ok')
```

diagrama imprastierii: viteza benzii/nr defectiuni

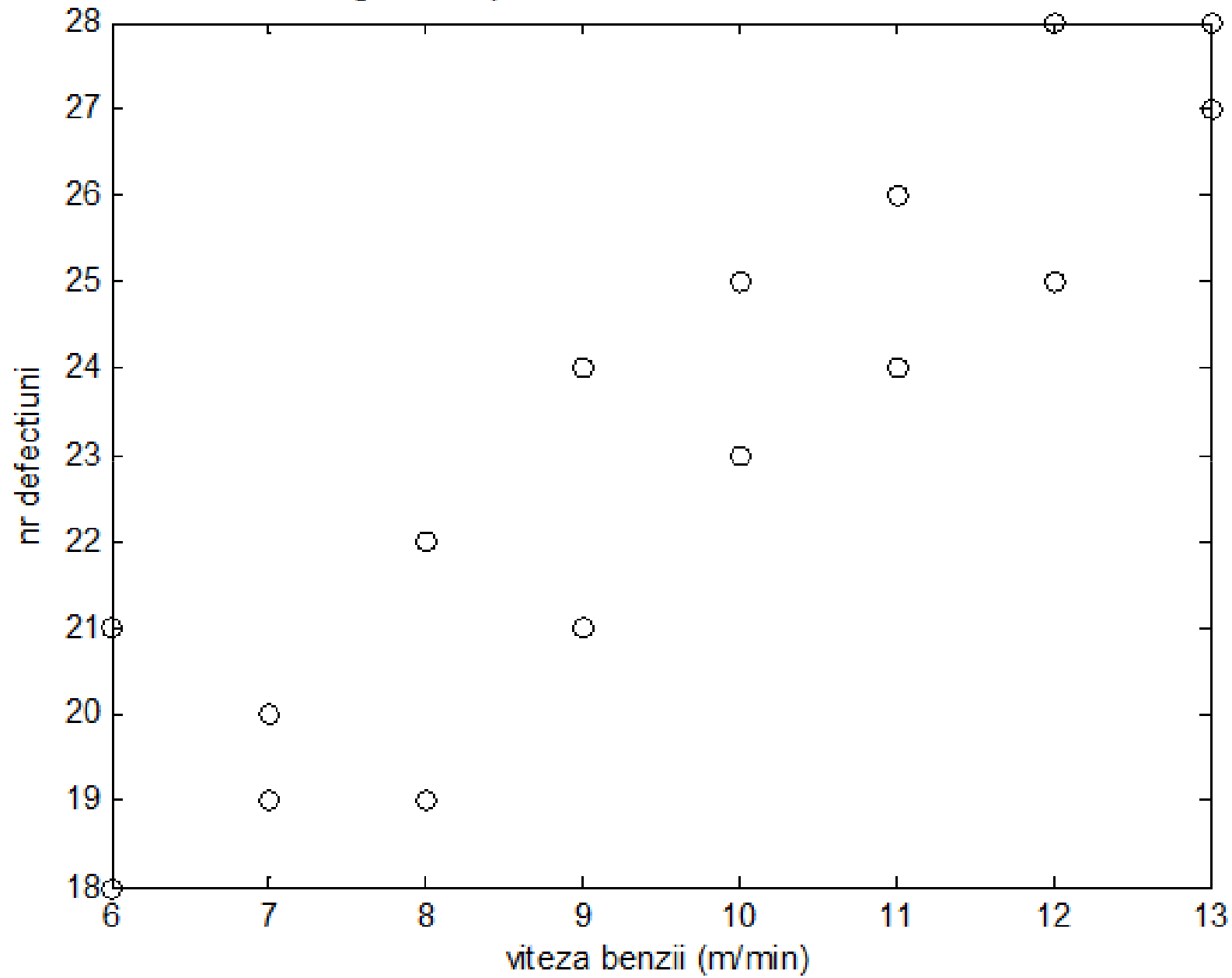


Diagrama împrăștierii dă informații importante privind legătura între cele două serii statistice fiind astfel un instrument util într-o analiză statistică.

Dacă norul are forma unei figuri geometrice alungite, simetrică față de o axă, există o legătură liniară între variabile. Forma de cerc sau pătrat a norului implică independența variabilelor.

Este importantă alegerea unităților de măsură a celor două variabile pentru a nu fi modificată forma norului.

O mare parte a studiilor statistice uzuale se ocupă cu analiza relației între două variabile statistice ce corespund aceluiași grup de subiecți/obiecte.

Cel mai cunoscut exemplu se referă la relația ce există între înălțimea și greutatea unui individ ce corespunde unor anumite standarde geografice, rasiale etc. Pentru a o identifica, se studiază relația dintre cele două caracteristici măsurate pe indivizii dintr-un anumit lot.

Suntem interesați în a descrie relația care ar putea exista între cele două variabile, analizând legătura între cele două serii de observații.

Concret, se analizează dacă tendința ascendentă a uneia implică o tendință ascendentă, descendentă sau nici o tendință a celeilalte. Scopul final, în ipoteza existenței unei legături reale între cele două variabile, este *prognoza* valorilor uneia în raport cu valorile celeilalte pe baza ecuației de regresie.

Coeficientul de corelație

Posibilele asociații între valorile a două variabile statistice continue, prelevate de la același grup de subiecți, sunt date de *coeficientul de corelație*.

Acest coeficient poate fi calculat pentru orice set de date.

Coeficientul de corelație are relevanță statistică, dacă sunt îndeplinite următoarele două condiții:

1. cele două variabile sunt definite de același lot de obiecte, cuplurile de date corespunzând aceluiasi obiect;
2. cel puțin una din variabile să aibă o repartiție aproximativ normală.

Considerăm, două serii statistice $\{x_i\}_{1 \leq i \leq n}$ și $\{y_i\}_{1 \leq i \leq n}$ corespunzătoare variabilelor statistice X și Y , generate de un grup de obiecte.

Prin *coeficientul de corelație* r al celor două variabile, numit și *Pearson's r* , vom înțelege numărul real r , cuprins între -1 și 1 , definit de formula:

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \cdot \sum_{i=1}^n (y_i - \bar{y})^2}},$$

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

Formula poate fi scrisă sub forma:

$$r = \frac{\sum_{i=1}^n x_i \cdot y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i \right) \cdot \left(\sum_{i=1}^n y_i \right)}{\sqrt{\left(\sum_{i=1}^n x_i^2 - \frac{1}{n} \cdot \left(\sum_{i=1}^n x_i \right)^2 \right) \cdot \left(\sum_{i=1}^n y_i^2 - \frac{1}{n} \cdot \left(\sum_{i=1}^n y_i \right)^2 \right)}} .$$

Exemplu

Calculăm coeficientul de corelație a celor două variabile din exemplul nr 1

```
>> X=[6 6 7 7 8 8 9 9 10 10 11 11 12 12 13 13];  
>> Y=[21 18 19 20 22 19 21 24 23 25 24 26 25 28 27 28];  
>> r=(X*Y'-16*mean(X)*mean(Y))/sqrt((X*X'-16*mean(X)^2)*  
(Y*Y'-16*mean(Y)^2))  
r =  
    0.9151
```

Construcția intervalului de încredere 95% pentru r

Intervalul de încredere 95% pentru r este intervalul care-l conține pe r cu o probabilitate de 0.95.

Plecând de la faptul că variabila aleatoare $z = \frac{1}{2} \ln \left(\frac{1+r}{1-r} \right)$ este normal

repartizată, rezultă că intervalul de încredere 95% pentru z are forma (z_1, z_2) ,

unde: $z_1 = z - \frac{1.96}{\sqrt{n-3}}$, $z_2 = z + \frac{1.96}{\sqrt{n-3}}$.

Aplicând transformarea inversă, obținem intervalul de încredere 95% pentru r ,

dat de: $\left(\frac{e^{2z_1} - 1}{e^{2z_1} + 1}, \frac{e^{2z_2} - 1}{e^{2z_2} + 1} \right)$

Exemplu

Construim intervalul de încredere 95% pentru r din exemplul nr 1

```
>> r = 0.9274
>> z = 1/2*log((1+r)/(1-r))
z =
    1.6396
>> z1 = z - 1.96/sqrt(17)
z1 =
    1.1642
>> z2 = z + 1.96/sqrt(17)
z2 =
    2.1149
>> a = (exp(2*z1) - 1) / (exp(2*z1) + 1)
a =
    0.8224
>> b = (exp(2*z2) - 1) / (exp(2*z2) + 1)
b =
    0.9713
```

Așadar, în acest caz intervalul de încredere 95% pentru r este:
(0.8224, 0.9713).

Coeficientul de corelație r (Pearson) ia valori cuprinse între -1 și $+1$.

- O valoare a lui r apropiată de $+1$ indică o corelație pozitivă puternică, (tendința de creștere puternică a unei variabile atunci când cealaltă variabilă crește)
- O valoare a lui r apropiată de -1 indică o corelație negativă puternică, (tendința de descreștere puternică a unei variabile atunci când cealaltă variabilă crește)
- Valoarea 0 a coeficientului de corelație indică independența liniară între cele două variabile.

Coeficientul de corelație r (Pearson) măsoară cât de *puternică* este legătura dintre două variabile.

Este foarte important cât de *semnificativă* este această legătură dintre două variabile.

În statistică un rezultat este *statistic semnificativ*, dacă este improbabil să fie obținut din întâmplare.

De obicei *nivelul de semnificație* (prag de probabilitate) mai mic de 0.05 înseamnă că rezultatul observat are cel mult 5% șanse de a fi obținut printr-o întâmplare, ceea ce ne permite să-l considerăm semnificativ.

Dacă inegalitatea:

$$|r| \cdot \sqrt{n-1} \geq 3$$

este adevărată, putem considera că cele două variabile sunt într-adevăr corelate.

În exemplul nr 1 avem:

```
>>r*sqrt(15)  
ans =  
    3.5442
```

ceea ce înseamnă că legătura între viteza benzii de lucru și numărul de defecțiuni este semnificativă.

Pentru a stabili cât de semnificativă este legătura dintre cele două variabile, se utilizează nivelul de semnificație p asociat calculării coeficientului r :

„dacă $p < 0.05$ atunci legătura este semnificativă”.

Putem calcula în Matlab matricea coeficienților de corelație și matricea nivelului de semnificație p folosind funcția `corrcoef`:

În exemplul nr 1 avem:

```
>> [r,p]=corrcoef(X,Y)
```

```
r =
```

```
    1.0000    0.9274  
    0.9274    1.0000
```

```
p =
```

```
    1.0000    0.0000  
    0.0000    1.0000
```

rezultând $r = 0.9274$ și $p=0$.

Cu cât este mai mic nivelul de semnificație, cu atât este mai *semnificativă* legătura.

Pe de altă parte, cu cât valoarea lui $|r|$ este mai apropiată de 1, cu atât mai *puternică* este legătura.

Factorul esențial este dimensiunea eșantionului.

- În cazul eșantioanelor mici se poate obține din întâmplare o corelație puternică; înainte de a concluziona este absolut necesară consultarea nivelului de semnificație p ;
- În cazul eșantioanelor mari este ușor de obținut un nivel de semnificație mic și astfel este necesar să se determine cât de puternică este legătura.

Prezentarea corelației

Prezentarea corelației dintre două variabile statistice trebuie să urmeze următorul model:

1. Se prezintă mai întâi diagrama de împrăștiere a norului de puncte.
2. Valoarea coeficientului de corelație r trebuie să aibă cel puțin două zecimale și să fie însoțită de nivelul de semnificație p și de intervalul de încredere corespunzător, dacă este posibil.
3. Numărul de observații analizate trebuie menționat.

Analiza matricei corelațiilor

Luând în considerare un set de variabile (atribute), X_1, X_2, \dots, X_k , $k > 2$, matricea de tip $k \times k$ al cărui element (i, j) este coeficientul de corelație al variabilelor (X_i, X_j) se numește *matricea corelațiilor*.

La fel ca și în cazul cuplurilor de variabile, putem analiza corelațiile între fiecare două perechi de variabile, precum și ,norul' împrăștierii lor .

Avantajul prezentării legăturii între atributele obiectelor cu ajutorul matricei corelațiilor și nu cu corelațiile fiecărei perechi de atribute în parte constă în faptul că astfel avem o privire de ansamblu asupra tuturor conexiunilor între atributele analizate.

De asemenea, ,norul' punctelor reprezentate, ca și alte reprezentări grafice sunt mult mai sugestive în cazul reprezentării colective a datelor decât luate separat.

2. Exemplu

Într-o firmă de IT, se presupune că Y , rezultatul vânzărilor, depinde de X_1 , suma cheltuită pentru reclama tradițională la TV și în ziare, și X_2 , suma cheltuită pentru reclama prin Internet, (sumele sunt în zeci de mii de euro).

X_1	5	7	7	5	9	11	12	13	17	17	18	19
X_2	2.9	2.5	3	1.5	4	4.5	5	5.5	6	7.5	8	7
Y	1	2	2.5	1.5	3.5	4	5	6	6.5	8	8	7.5

Prezentăm o asemenea analiză multiplă (multivariată) a acestor date.

Studiem cât de puternică este legătura între variabila răspuns Y și variabilele predictive $X1$, $X2$, scriind matricea corelațiilor și calculând nivelul de semnificație p :

```
>> Y=[5 7 7 5 9 11 12 13 15 17 18 19];  
>> X1=[2 2.5 3 1.5 4 4.5 5 5.5 6 7.5 8 7];  
>> X2=[1 2 2.5 1.5 3.5 4 5 6 6.5 8 8 7.5];  
>> A=[X1' X2' Y'];[r,p]=corrcoef(A)
```

$r =$

1.0000	0.9899	0.9825
0.9899	1.0000	0.9858
0.9825	0.9858	1.0000

$p =$

1.0000	0.0000	0.0000
0.0000	1.0000	0.0000
0.0000	0.0000	1.0000

Diagrama împrăştierii

```
>> plot3(X1,X2,Y,'o')  
>> grid on  
>> axis square
```

diagrama imprastierii: cheltuieli reclama traditionala /reclama on-line/rezultat vanzari

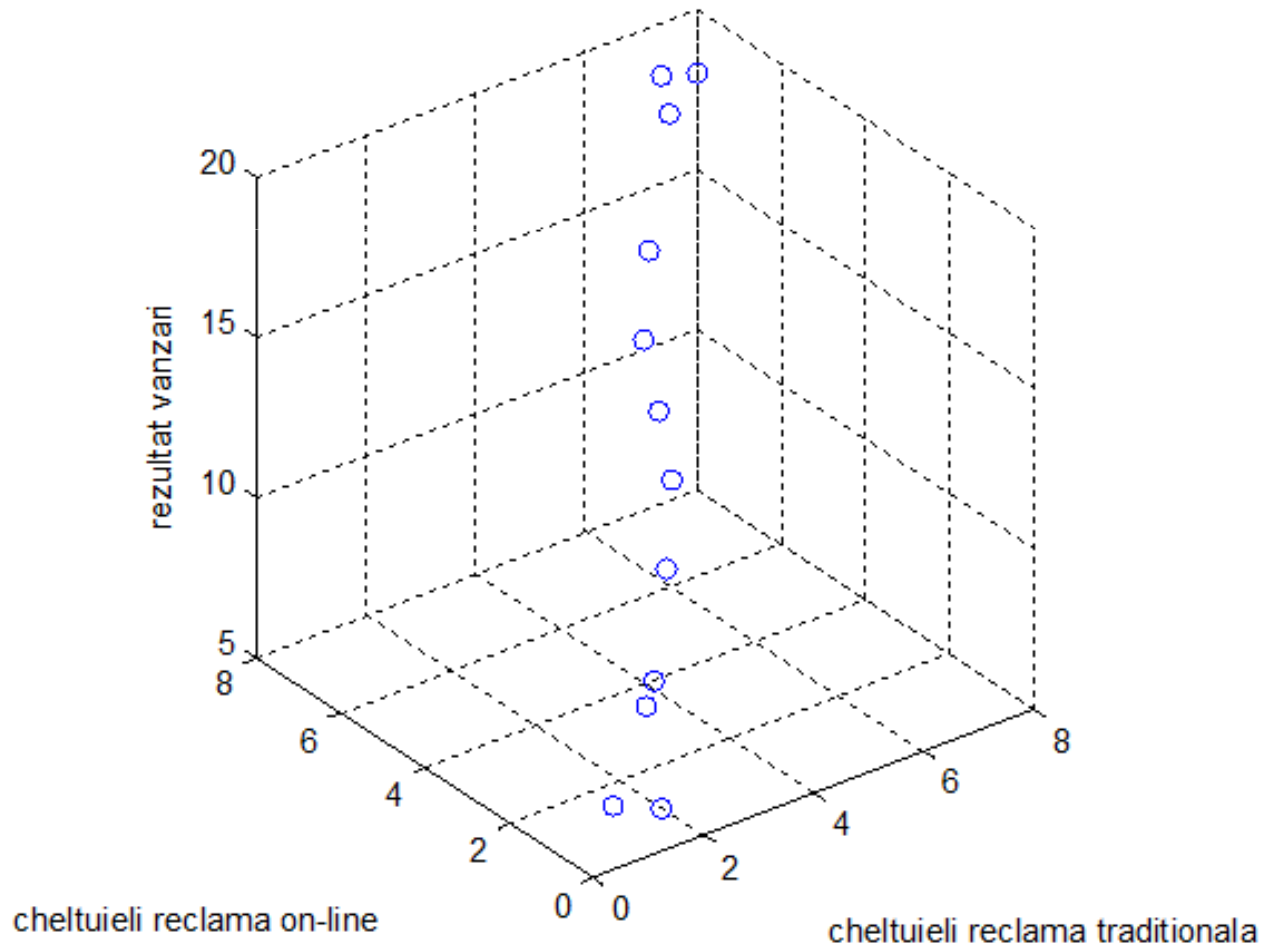
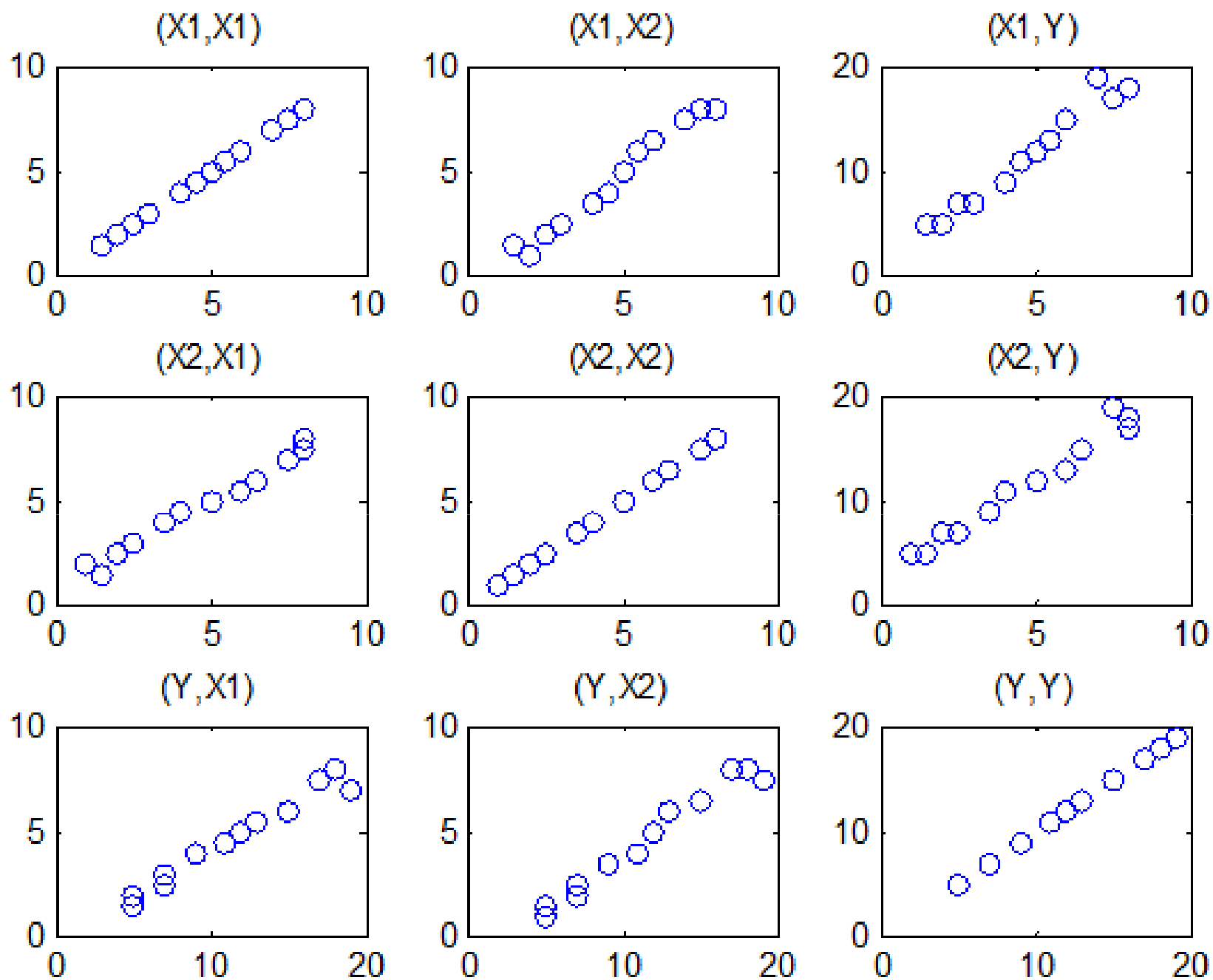


Diagrama corelațiilor multiple

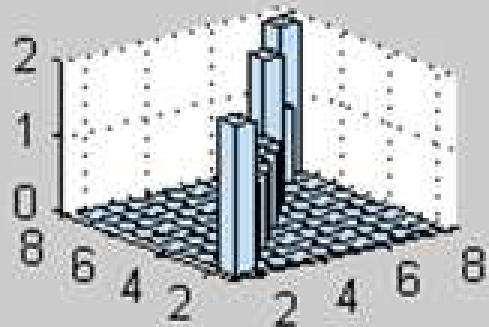
În ceea ce privește ilustrarea grafică a corelațiilor de mai sus, vizualizarea este foarte utilă, rezumând sugestiv toate informațiile numerice prezentate arid în tabel

```
>> subplot(331);plot(X1,X1,'o');  
>> subplot(332);plot(X1,X2,'o');  
>> subplot(333);plot(X1,Y,'o');  
>> subplot(334);plot(X2,X1,'o');  
>> subplot(335);plot(X2,X2,'o');  
>> subplot(336);plot(X2,Y,'o');  
>> subplot(337);plot(Y,X1,'o');  
>> subplot(338);plot(Y,X2,'o');  
>> subplot(339);plot(Y,Y,'o');
```

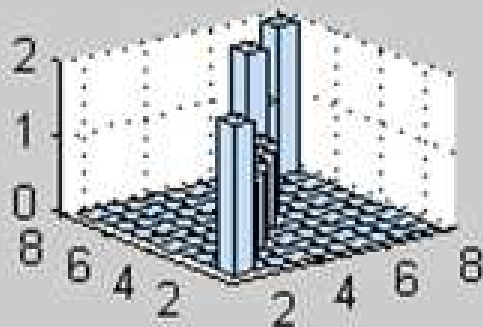


```
>> subplot(331);A1=[X1' X1'];hist3(A1)
>> subplot(332);A2=[X1' X2'];hist3(A2)
>> subplot(333);A3=[X1' Y'];hist3(A3)
>> subplot(334);A4=[X2' X1'];hist3(A4)
>> subplot(335);A5=[X2' X2'];hist3(A5)
>> subplot(336);A6=[X2' Y'];hist3(A6)
>> subplot(337);A7=[Y' X1'];hist3(A7)
>> subplot(338);A8=[Y' X2'];hist3(A8)
>> subplot(339);A9=[Y' Y'];hist3(A9)
```

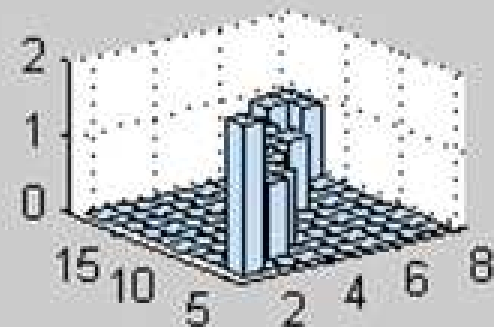
(X1,X1)



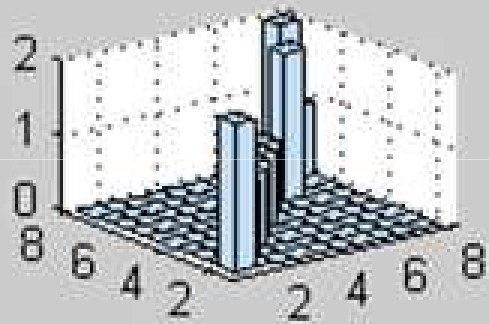
(X1,X2)



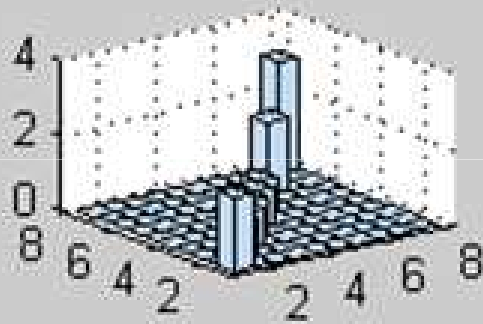
(X1,Y)



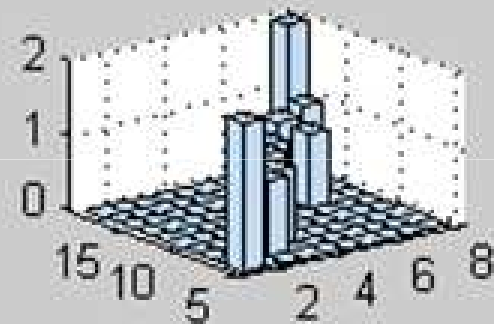
(X2,X1)



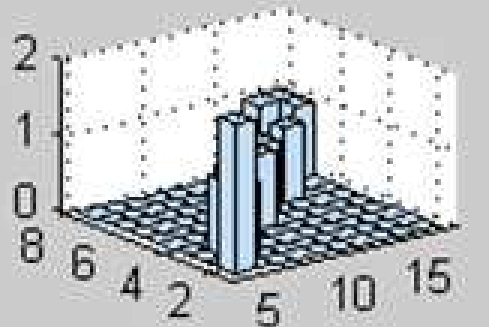
(X2,X2)



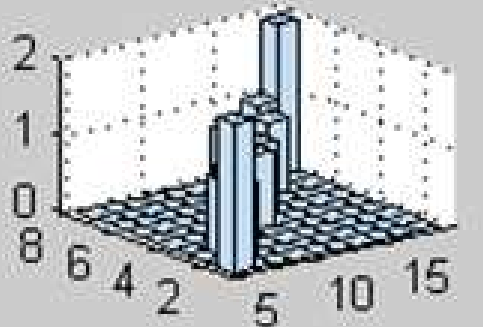
(X2,Y)



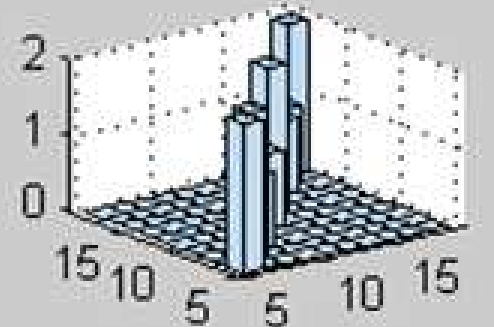
(Y,X1)



(Y,X2)



(Y,Y)



Regresia liniară

Regresia liniară este tehnica statistică ce permite modelarea relației liniare între o variabilă explicativă (notată cu X) și variabila ce urmează a fi explicată (notată cu Y).

Teoremă

Coeficientul de corelație r a două variabile aleatoare X și Y ia valori în intervalul $[-1, 1]$.

Dacă variabilele sunt independente coeficientul de corelație este nul.

Coeficientul de corelație r este egal cu ± 1 dacă și numai dacă variabilele X și Y verifică ecuația:

$$aX + bY = c \Leftrightarrow Y = AX + B, \quad a, b, c, A, B \in \mathbf{R}$$

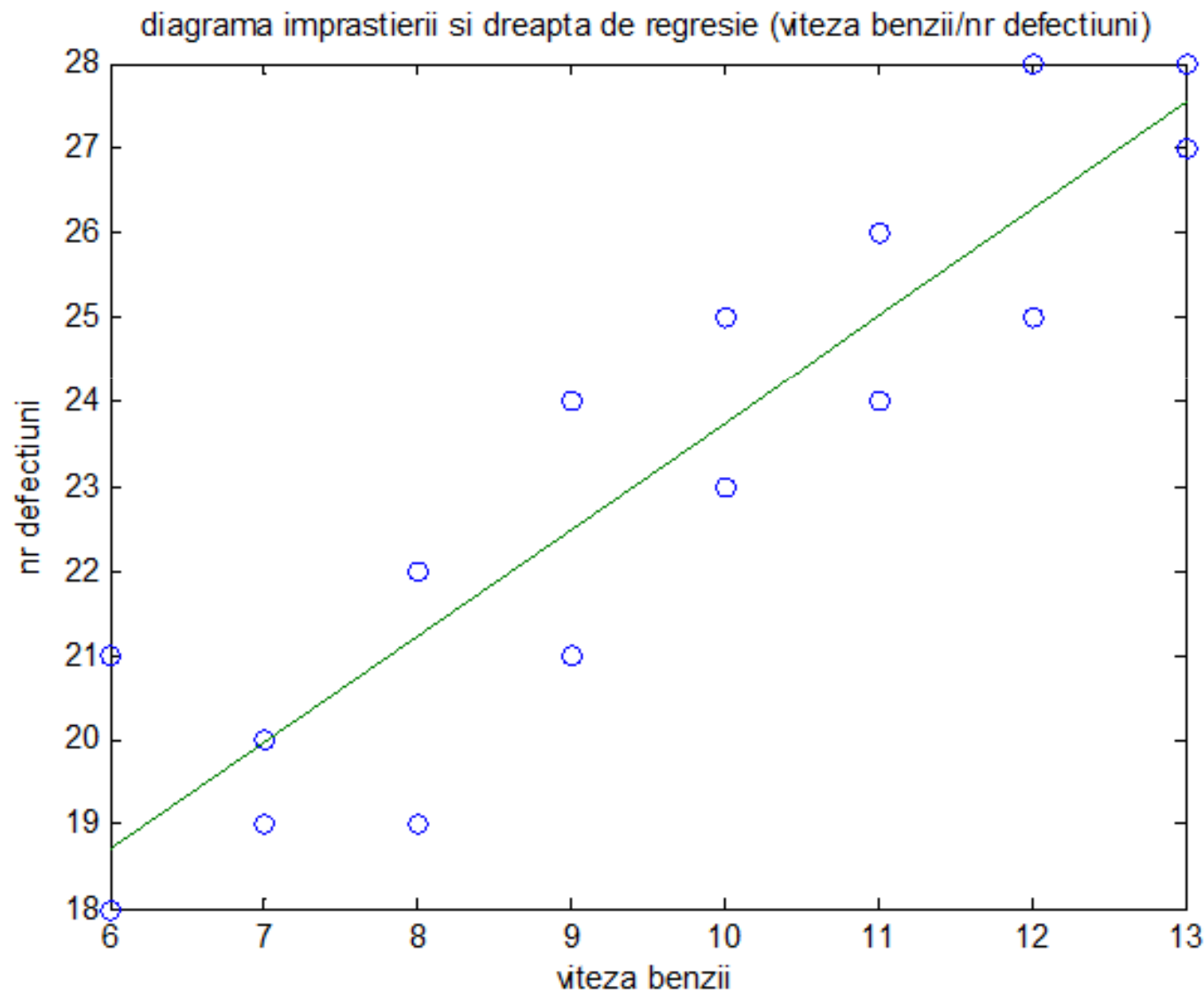
Reciproca nu este totdeauna adevărată, adică există variabile dependente având coeficientul de corelație zero.

Ecuția de regresie

În cazul în care există o legătură liniară între două variabile (numerice), se consideră una dintre variabile ca *variabilă independentă (predictor)*, iar cealaltă variabilă ca *variabilă dependentă (răspuns)*.

Legătura liniară dintre cele două variabile este descrisă de o ecuație liniară, așa-numita *ecuația de regresie*, căreia îi corespunde geometric *dreapta de regresie*.

În exemplul nr 1 avem:



Dreapta de regresie

Dreapta de regresie este acea dreaptă (ideală) ce trece prin norul de puncte format de perechile de date ale celor două variabile și care minimizează distanța între date și ea, prin minimizarea sumei pătratelor distanțelor. Ecuația de drepte de regresie este de forma:

$$y = a + bx$$

unde a se numește *interceptor*, iar b se numește *coeficientul de regresie* .

Ecuția dreptei de regresie se stabilește pe baza metodei “*celor mai mici pătrate*”, care în cazul nostru constă în calcularea distanțelor dintre punctele observate (reale) corespunzătoare cuplului de serii statistice și punctele (imaginare) de pe o anumită dreaptă (de regresie), ce trece prin mijlocul norului de puncte generate de cuplurile de date, distanțe cunoscute sub numele de *reziduuri*.

Datele existente satisfac ecuațiile:

$$y_i = a + bx_i + \varepsilon_i \quad 1 \leq i \leq p$$

unde:

- x_i și y_i , $1 \leq i \leq p$ sunt numere reale cunoscute;
- ε_i sunt erorile.

Metoda celor mai mici pătrate constă în aflarea parametrilor a și b pentru care

expresia $S(a, b) = \sqrt{\sum_{i=1}^p (y_i - a - b \cdot x_i)^2}$ este minimă.

Regresia liniară este utilizată, dacă sunt îndeplinite următoarele trei ipoteze de lucru :

- Valorile variabilei dependente Y trebuie să aibă o repartiție normală (gaussiană);
- Variabilitatea variabilei prognozate Y trebuie să fie asemănătoare cu cea a predictorului X (dispersia sau deviația standard asemănătoare);
- Legătura dintre cele două variabile, predictorul și variabila dependentă, trebuie să fie liniară (verificare empirică pe baza norului de puncte, care trebuie să aibă o formă alungită).

Modul standard de a verifica simultan toate cele trei ipoteze de lucru este analiza statistică a reziduurilor. Astfel, se poate demonstra că dacă toate cele trei ipoteze sunt verificate simultan, atunci reziduurile sunt normal repartizate de medie zero.

Dreapta de regresie este dată de ecuația:

$$y = a + b \cdot x, \text{ unde:}$$

$$b = \frac{n(\sum_{i=1}^n x_i y_i - n \cdot \bar{x} \cdot \bar{y})}{\sum_{i=1}^n x_i^2 - n \cdot \bar{x}^2} \quad \text{și} \quad a = \bar{y} - b \cdot \bar{x}$$

cu notațiile $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$, $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$.

3. Exemplu

Următoarea baza de date conține informații despre prețul de vânzare a 15 apartamente dintr-un cartier Ω al unui oraș Z , în trimestrul al III-lea din 2011. Prezintă datele referitoare la preț (în mii de euro) și suprafața locuibilă (în metri pătrați)

supr	30	32	35	48	49	50	55	60	65	70	75	80	82	85	90
preț	30	29	31	36	38	40	48	47	53	49	58	60	59	71	65

```
>> X1=[ 30 32 35 48 49 50 55 60 65 70 75 80 82 85 90];  
>> Y=[30 29 31 36 38 40 48 47 53 49 58 60 59 71 65];  
>> A=[X1' Y'];[r p]=corrcoef(A)  
r =  
    1.0000    0.9729  
    0.9729    1.0000  
p =  
    1.0000    0.0000  
    0.0000    1.0000
```

Se observă că variabilele sunt puternic (pozitiv) corelate, coeficientul de corelație fiind 0.9729 și nivelul de semnificație este $p = 0.0000$.

Calculăm dreapta de regresie, folosind formulele prezentate mai sus, și o desenăm în același sistem de axe cu diagrama împrăștierii:

```
>> b=(X1*Y'-15*mean(X1)*mean(Y))/((norm(X1))^2-15*(mean(X1))^2)
```

```
b =
```

```
0.6572
```

```
>> a=mean(Y)-b*mean(X1)
```

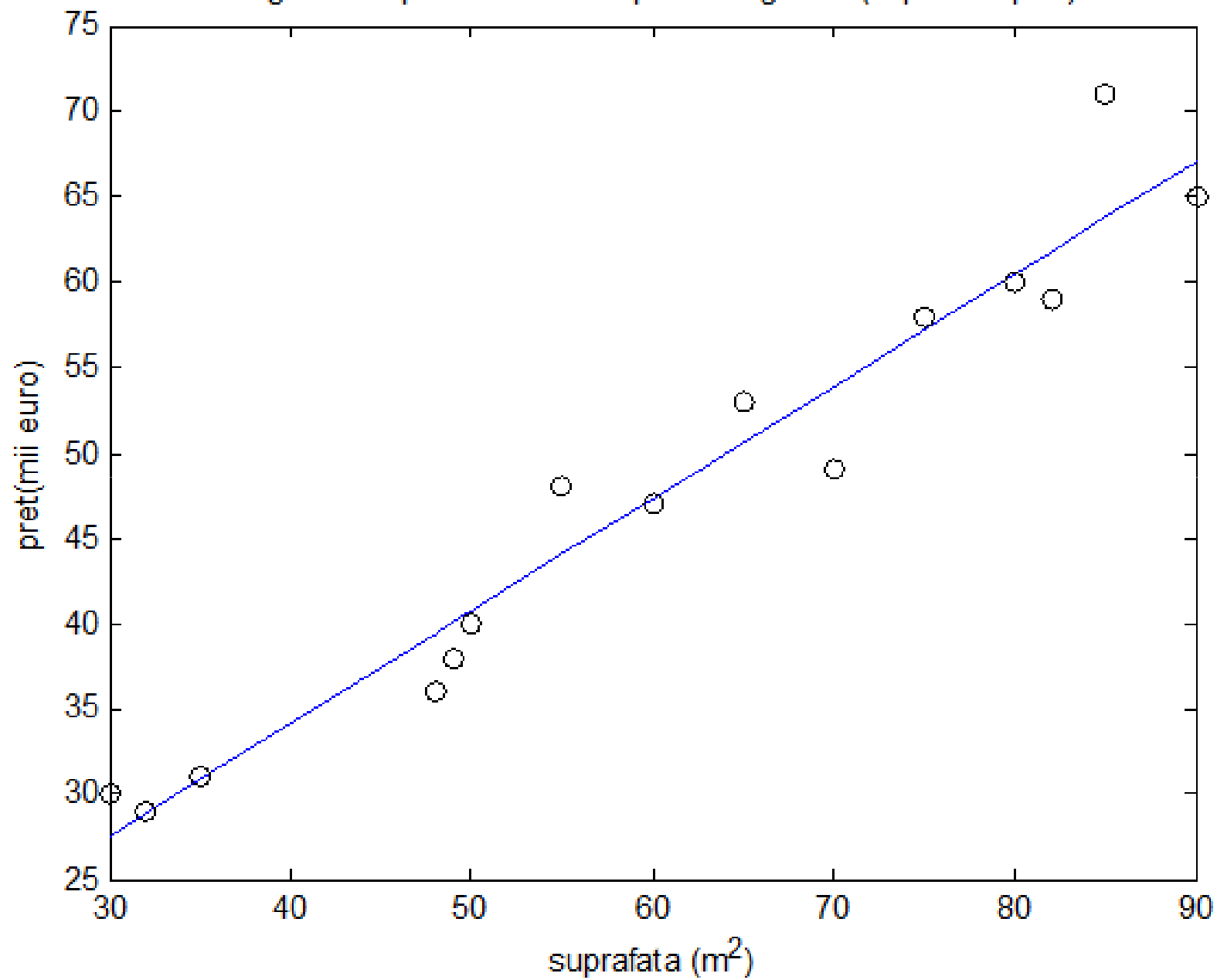
```
a =
```

```
7.9053
```

```
>> plot(X1,a+b*X1);hold on
```

```
>> plot(X1,Y,'ko');hold off
```

diagrama imprastierii si dreapta de regresie (suprafata/pret)



Dreapta de regresie obținută poate fi utilizată în prognoză.

De exemplu, pe baza datelor avute la dispoziție, prețul de vânzare al unui apartament de 72m^2 , reaspectiv de 87 ar fi:

```
>> syms X1
```

```
>> f=a+b*X1;\
```

```
>> p1=subs(f,X1,72)
```

```
p1 =  
55.2235
```

```
>> p2=subs(f,X1,87)
```

```
p2 =  
65.0814
```

Pentru calculul coeficienților a și b putem apela la scrierea matricială a ecuațiilor

$$y_i = a + bx_i + \varepsilon_i \quad 1 \leq i \leq p,$$

și anume:

$$\begin{pmatrix} y_1 \\ \vdots \\ y_p \end{pmatrix} = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_p \end{pmatrix} \cdot \begin{pmatrix} a \\ b \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_p \end{pmatrix}.$$

Presupunând că matricea $\begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_p \end{pmatrix}$ este de rang 2, utilizând metoda

celor mai mici pătrate, putem calcula parametrii a și b :

$$\begin{pmatrix} a \\ b \end{pmatrix} = \left(\begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_p \end{pmatrix}' \cdot \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_p \end{pmatrix} \right)^{-1} \cdot \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_p \end{pmatrix}' \cdot \begin{pmatrix} y_1 \\ \vdots \\ y_p \end{pmatrix}$$

Exemplu

Reluăm exemplul nr.3

```
>> X1=[ 30 32 35 48 49 50 55 60 65 70 75 80 82 85 90];  
>> i=1; B=[1]; for i=2:15 B=[1 [B]]; end  
>> X=[B' X1'];  
>> a=inv(X'*X)*X'*Y'  
a =  
    7.9053  
    0.6572
```

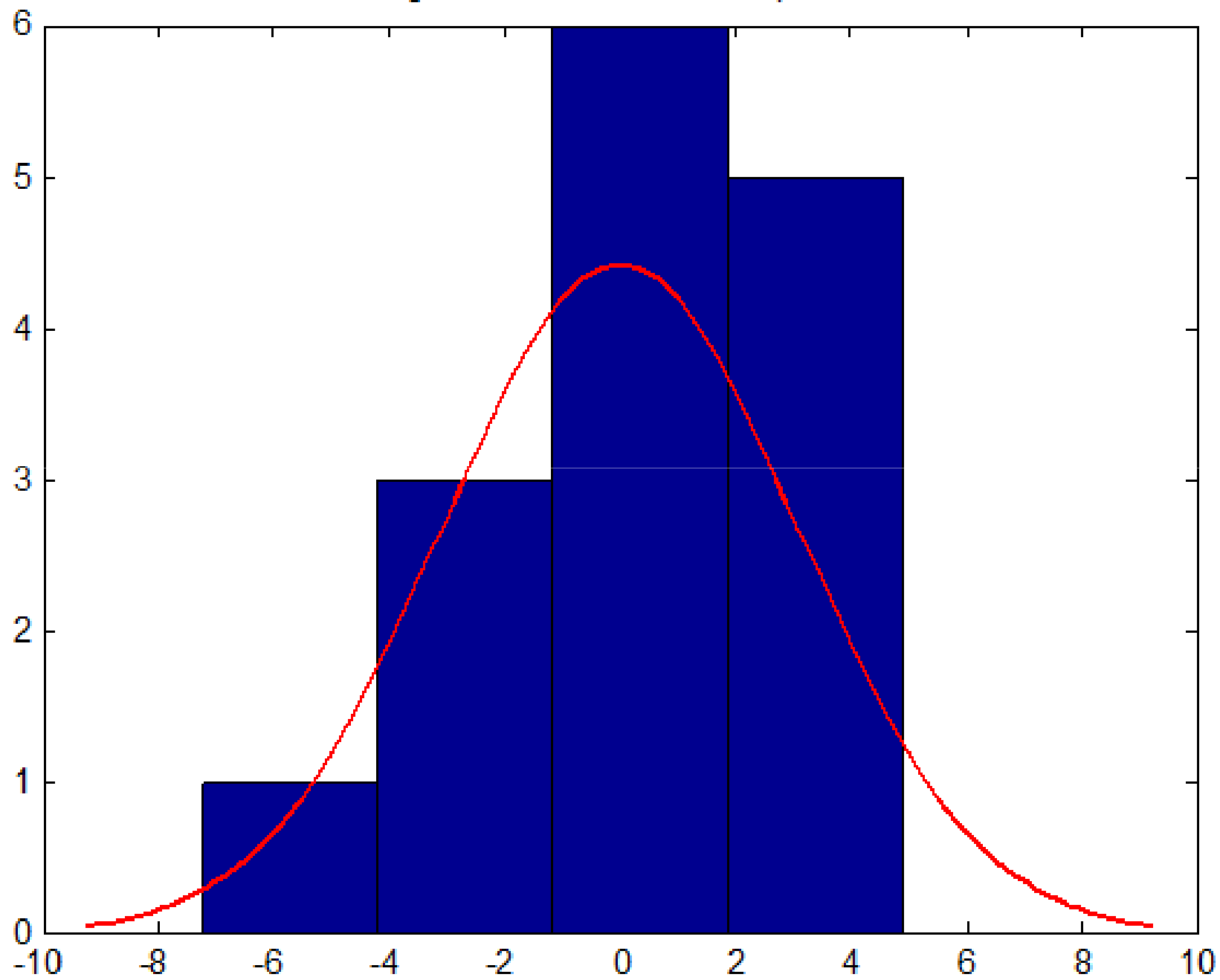
În Matlab, problema se poate rezolva mai simplu:

```
>> a=X\Y'  
a =  
    7.9053  
    0.6572
```

Vom calcula reziduurile, media acestora și vom desena histograma lor, suprapunând graficul densității de repartiție normală

```
>> Rez=a(1)+a(2)*X1-Y;  
>> mean(Rez)  
ans =  
-2.2974e-01  
>> histfit(Rez)
```

histograma reziduurilor-exemplul nr 3



Regresia neliniară

În cazul când legătura dintre cele două variabile statistice nu este liniară și totuși există, avem de a face cu o *regresie neliniară*.

Variabila dependentă este o combinație neliniară a variabilelor independente.

Ne vom limita la modelarea relației dintre o variabilă explicativă și variabila ce urmează a fi explicată printr-o ecuație a unei curbe de regresie:

- ecuație polinomială,
- ecuație mixt exponențială.

Regresia polinomială

Regresia polinomială va furniza o ecuație de forma:

$$Y = b_0 + b_1X + b_2X^2 + \dots + b_nX^n,$$

unde, ca de obicei, X este variabila predictoare, iar Y este variabila prognozată, adică variabila răspuns. Pentru a determina coeficienții de regresie și interceptorul vom aplica *metoda celor mai mici pătrate*:

Presupunem că datele existente satisfac ecuațiile:

$$Y_i = b_0 + b_1 \cdot X_i + b_2 \cdot X_i^2 + \dots + b_n \cdot X_i^n + \varepsilon_i, \quad 1 \leq i \leq p,$$

unde:

- $X_i, Y_i, \quad 1 \leq i \leq p$ sunt numere reale cunoscute (reținem că $p > n$),
- ε_i sunt variabile aleatoare necunoscute.

Sub formă matriceală, modelul devine:

$$Y = X \cdot b + \varepsilon,$$

unde:

- $Y = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_p \end{pmatrix}$ este un vector aleator p -dimensional,

- $X = \begin{pmatrix} 1 & X_1 & \dots & X_1^n \\ 1 & X_2 & \dots & X_2^n \\ \dots & \dots & \dots & \dots \\ 1 & X_p & \dots & X_p^n \end{pmatrix},$

- b este vectorul $(n+1)$ -dimensional, ale cărui componente sunt parametrii necunoscuți ai modelului,

- ε este vectorul p -dimensional al erorilor.

Metoda celor mai mici pătrate constă în aflarea vectorului b , pentru care expresia:

$$S(b) = \|Y - Xb\|^2 = \sqrt{\sum_{i=1}^p (Y_i - b_0 - b_1 \cdot X_i - b_2 \cdot X_i^2 - \dots - b_n \cdot X_i^n)^2},$$

este minimă, presupunând că matricea X are rangul n . Astfel $X' \cdot X$ este inversabilă și avem:

$$b = (X' \cdot X)^{-1} \cdot X' \cdot Y.$$

4. Exemplu

Considerăm situația numărului de pacienți internați într-o secție a unui spital pe parcursul a 11 săptămâni, situație reflectată în următorul tabel:

săptămâna	1	2	3	4	5	6	7	8	9	10	11
nr pacienți	50	92	116	125	135	150	145	130	120	95	85

Determinăm curba de regresie polinomială, considerând pentru început că datele existente satisfac relațiile:

$$Y_i = b_0 + b_1 \cdot t_i + b_2 \cdot t_i^2 + \varepsilon_i, \quad 1 \leq i \leq 6.$$

Utilizând metoda celor mai mici pătrate, vom rezolva ecuația matriceală:

$$\begin{pmatrix} Y_1 \\ Y_2 \\ Y_3 \\ Y_4 \\ Y_5 \\ Y_6 \\ Y_7 \\ Y_8 \\ Y_9 \\ Y_{10} \\ Y_{11} \end{pmatrix} = \begin{pmatrix} 1 & t_1 & t_1^2 \\ 1 & t_2 & t_2^2 \\ 1 & t_3 & t_3^2 \\ 1 & t_4 & t_4^2 \\ 1 & t_5 & t_5^2 \\ 1 & t_6 & t_6^2 \\ 1 & t_7 & t_7^2 \\ 1 & t_8 & t_8^2 \\ 1 & t_9 & t_9^2 \\ 1 & t_{10} & t_{10}^2 \\ 1 & t_{11} & t_{11}^2 \end{pmatrix} \cdot \begin{pmatrix} b_0 \\ b_1 \\ b_2 \end{pmatrix}; \text{ notam } X = \begin{pmatrix} 1 & t_1 & t_1^2 \\ 1 & t_2 & t_2^2 \\ 1 & t_3 & t_3^2 \\ 1 & t_4 & t_4^2 \\ 1 & t_5 & t_5^2 \\ 1 & t_6 & t_6^2 \\ 1 & t_7 & t_7^2 \\ 1 & t_8 & t_8^2 \\ 1 & t_9 & t_9^2 \\ 1 & t_{10} & t_{10}^2 \\ 1 & t_{11} & t_{11}^2 \end{pmatrix}$$

și astfel obținem coeficienții polinomului de gradul al doilea:

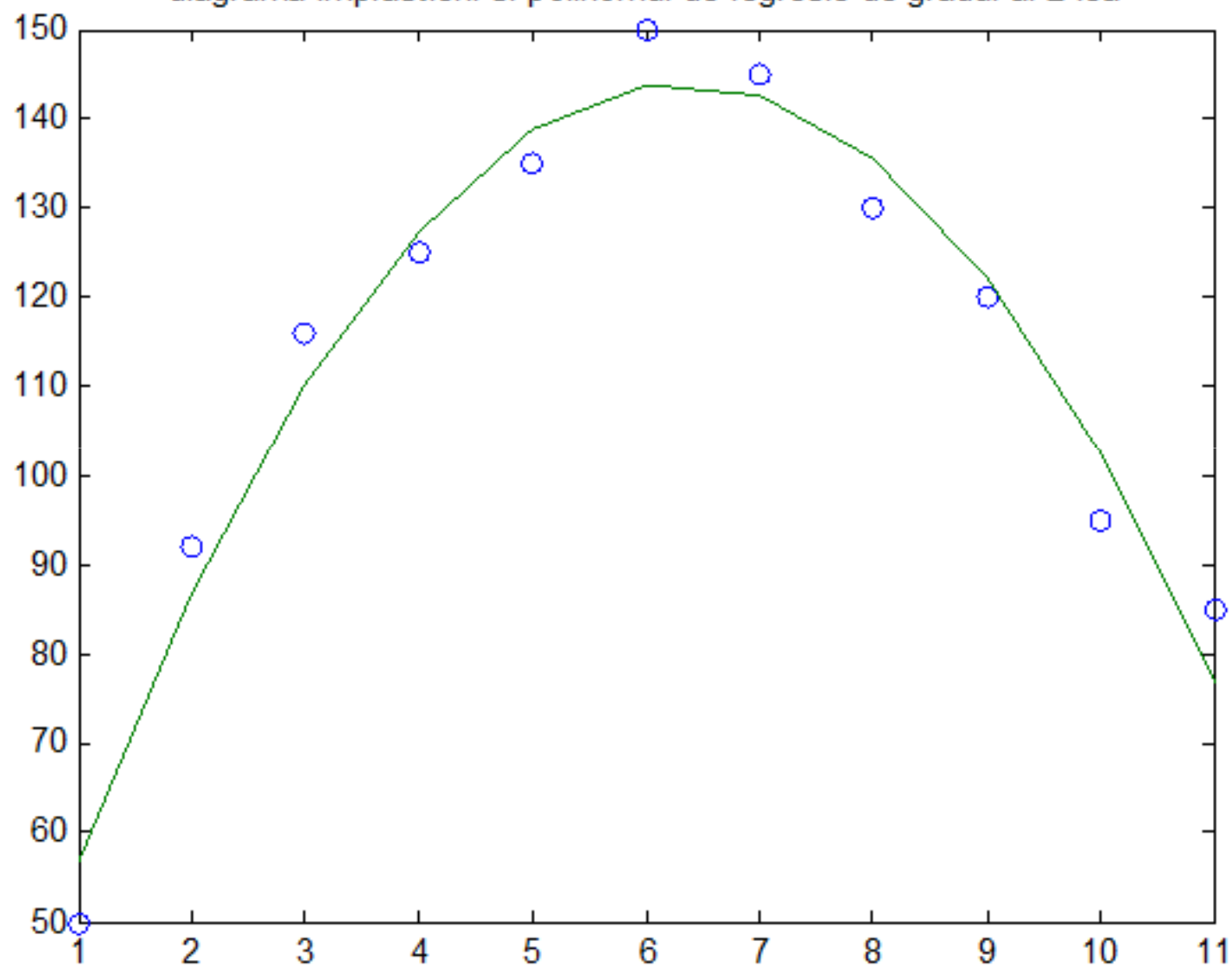
$$b = (X' \cdot X)^{-1} \cdot X' \cdot Y ,$$

```
>> t=1:11; P=[50 92 116 125 135 150 145 130 120 95 85];  
>> i=1;B=1;for i=2:11 B=[1 B];end  
>> for i=1:11 t2(i)=t(i)^2;end  
>> X=[B' t' t2'];b2=X\P';  
>> X=[B' t' t2'];b2=X\P'  
b2 =  
    21.0848  
    38.9000  
   -3.0758
```

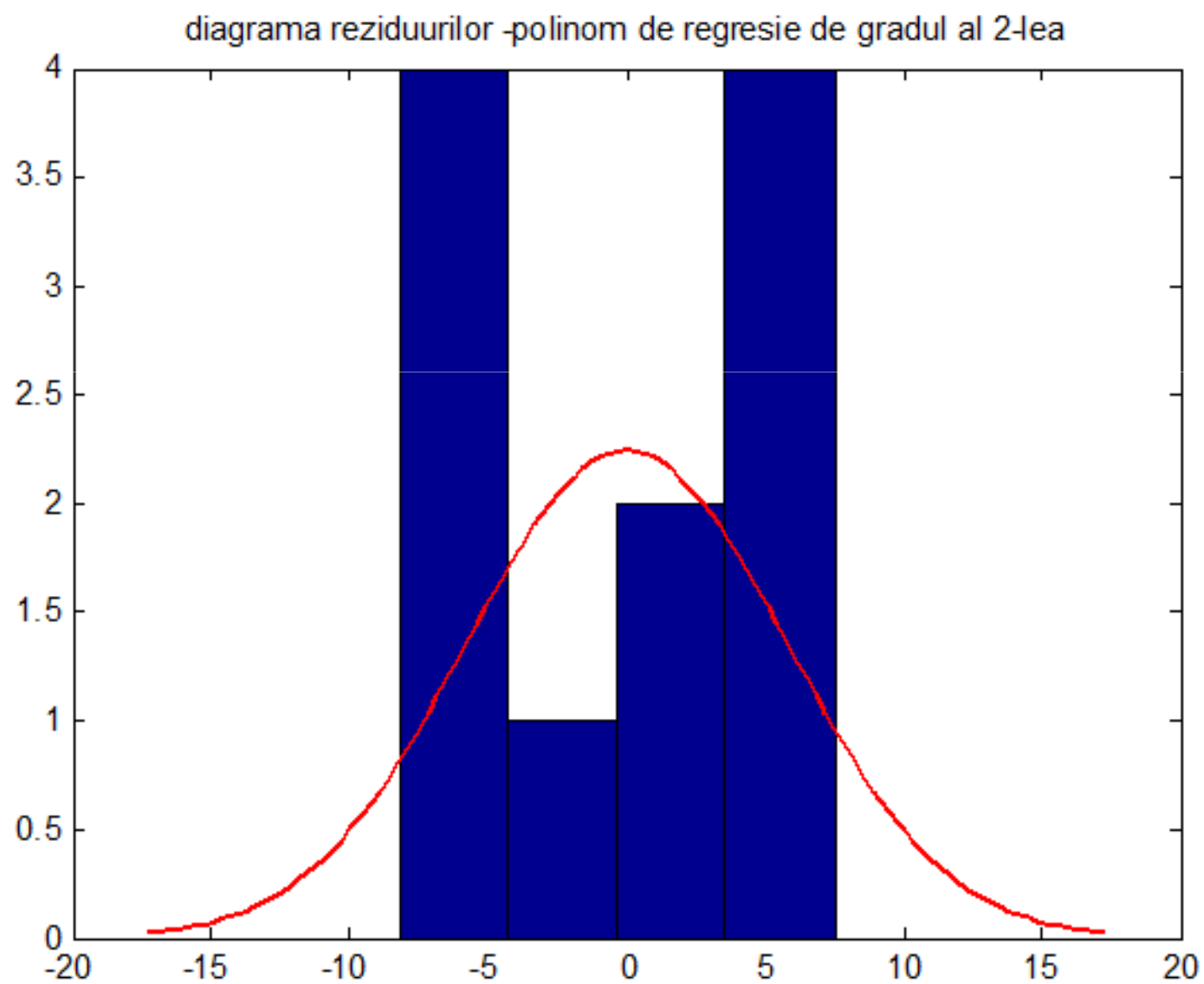
Desenăm diagrama împrăștierei și polinomul de regresie obținut:

```
>> P2=b2(1)+b2(2)*t+b2(3)*t.^2;plot(t,P,'O',t,P2)
```

diagrama imprastierii si polinomul de regresie de gradul al 2-lea



```
>> Rez=P2-P;  
>> mean(Rez)  
ans =  
 3.6173e-014  
>> histfit(Rez)
```

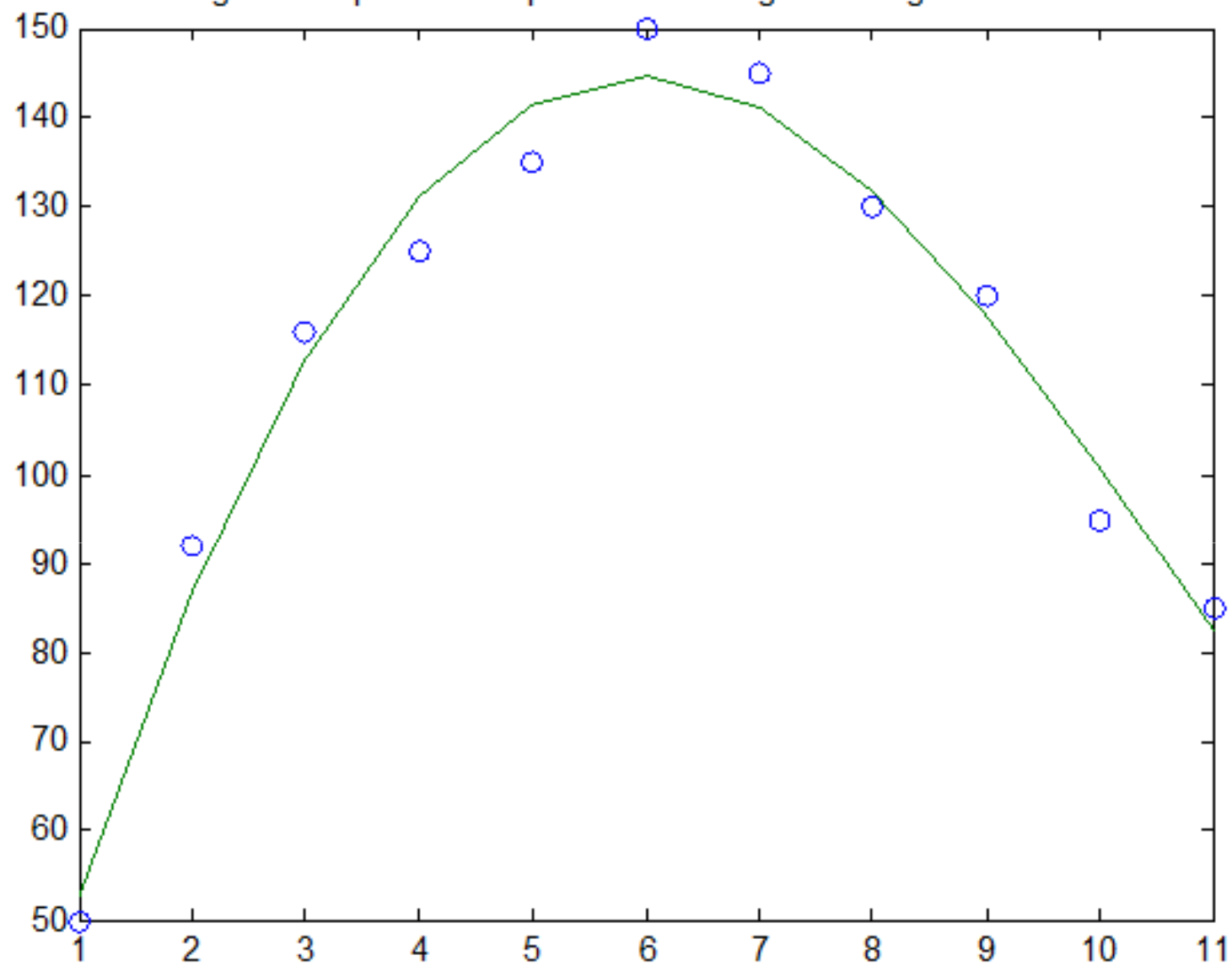


Încercăm să determinăm un polinom de regresie de grad 4, ceea ce înseamnă să presupunem că:

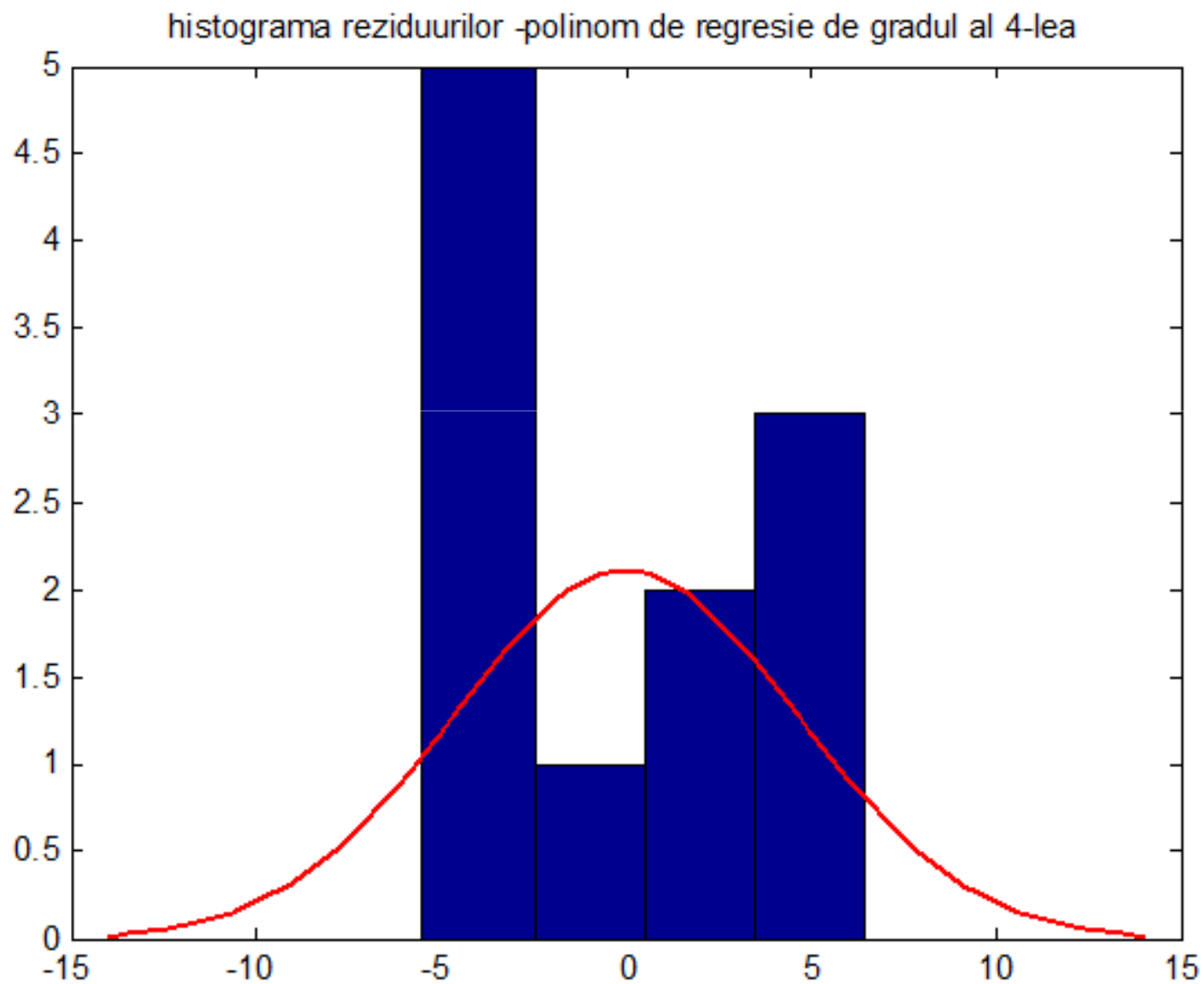
$$Y_i = b_0 + b_1 \cdot t_i + b_2 \cdot t_i^2 + b_3 \cdot t_i^3 + b_4 \cdot t_i^4 + \varepsilon_i, \quad 1 \leq i \leq 6,$$

```
>> t=1:11; P=[50 92 116 125 135 150 145 130 120 95 85];  
>> i=1;B=1;for i=2:11 B=[1 B];end  
>> for i=1:11 t2(i)=t(i)^2;end  
>> for i=1:11 t3(i)=t(i)^3;end  
>> for i=1:11 t4(i)=t(i)^4;end  
>> X=[B' t' t2' t3' t4'];b4=X\P'  
b4 =  
    11.5152  
    44.8918  
    -3.4257  
    -0.1249  
     0.0108  
>> P4=b4(1)+b4(2)*t+b4(3)*t.^2+ b4(4)*t.^3+b4(5)*t.^4;plot(t,P,'O',t,P4)
```

diagrama imprastierii si polinomul de regresie de gradul al 4-lea




```
>> Rez=P4-P;mean(Rez)
ans =
    2.0670e-014
>> histfit(Rez)
```



Polinoamele de regresie pot fi utilizate pentru a prognoza numărul de pacienți din următoarele săptămâni.

Pentru a decide care polinom merită utilizat, vom calcula mediile pătratelor reziduurilor (mediile pătratelor erorii) pentru a stabili care curbă de regresie dă o aproximare mai bună.

```
>> MSE2=mean((P-P2).^2)
MSE2 =
    30.2788
```

```
>> MSE4= mean(( P-P4).^2)
MSE4 =
    19.7839
```

Vom prognoza totuși numărul pacienților ce urmează a fi internați în săptămâna a 12-a folosind ambele polinoame:

$$>> Yp2=b2(1)+b2(2)*12+b2(3)*12^2$$

$$Yp2 =$$

44.9758

$$>> Yp4=b4(1)+b4(2)*12+b4(3)*12^2+b4(4)*12^3+b4(5)*12^4$$

$$Yp4 =$$

64.6364

În general, pentru a obține un rezultat bun este necesar un număr mare de date, să zicem n .

Cu ajutorul a $\frac{3}{4}n$ din aceste date, aleator alese scriem ecuația polinomului de regresie de gradul al 2-lea respectiv al 4-lea. Cu ajutorul acestor ecuații obținem prognoza pentru celelalte $\frac{n}{4}$ date care nu au fost utilizate și calculăm media pătratelor reziduurilor obținute pentru cele două situații considerate.

Reluăm procedeul încă de 4 ori, alegând de fiecare dată aleator alte $\frac{3}{4}n$ date. Vom alege aceea curbă polinomială de regresie pentru care media pătratelor reziduurilor obținute este cea mai mică.

Polinoame de regresie in Matlab

Prezentăm un studiu de caz în Matlab, privind aproximarea unui set de date $\{(x_i, y_i), 1 \leq i \leq n\}$ printr-o funcție polinomială. Funcțiile folosite sunt:

- `polyfit` - funcție ce returnează coeficienții polinomului în ordinea $b_n, b_{n-1}, \dots, b_1, b_0$.
- `polyval` - funcție ce returnează valorile polinomului în punctele $\{x_i, 1 \leq i \leq n\}$.

Reluăm exemplul nr 4 calculând pentru început coeficienții polinoamelor de regresie de gradul al 2-lea, respectiv al 4-lea, și apoi reziduurile în aceste cazuri:

```
>> t=1:11; P=[50 92 116 125 135 150 145 130 120 95 85];
```

```
>> b2=polyfit(t,P,2)
```

```
b2 =
```

```
 -3.0758  38.9000  21.0848
```

```
>> b4=polyfit(t,P,4)
```

```
b4 =
```

```
 0.0108 -0.1249 -3.4257  44.8918  11.5152
```

```
>> pol2 = polyval(b2,t,P);
```

```
>> pol4 = polyval(b4,t,P);
```

```
 8.1818  2.5874
```

```
>> MSE2=mean((P-pol2).^2);MSE4=mean((P-pol4).^2);[MSE2 MSE4]
```

```
ans =
```

```
 30.2788  19.7839
```

Regresie mixt exponențială

Considerând un set de date $\{(x_i, y_i), 1 \leq i \leq n\}$, vom construi funcția de regresie *mixt exponențială*:

$$y = b_0 + b_1 \cdot e^{-x} + b_2 \cdot x \cdot e^{-x}.$$

Presupunem că datele existente satisfac relațiile:

$$Y_i = b_0 + b_1 \cdot e^{-X_i} + b_2 \cdot X_i e^{-X_i} + \varepsilon_i, \quad 1 \leq i \leq p$$

unde:

- $X_i, Y_i, \quad 1 \leq i \leq p$ sunt numere reale cunoscute,
- ε_i sunt variabile aleatoare necunoscute.

Sub formă matriceală modelul devine:

$$Y = X \cdot b + \varepsilon ,$$

unde:

- $Y = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_p \end{pmatrix}$ este un vector aleator p -dimensional,

- $X = \begin{pmatrix} 1 & e^{-X_1} & X_1 e^{-X_1} \\ 1 & e^{-X_2} & X_2 e^{-X_2} \\ \vdots & \vdots & \vdots \\ 1 & e^{-X_p} & X_p e^{-X_p} \end{pmatrix}$,

- b este vectorul 3-dimensional ale cărui componente sunt parametri necunoscuți ai modelului, anume b_0, b_1, b_2 ,

- ε este vectorul p -dimensional al erorilor.

Metoda celor mai mici pătrate constă în aflarea vectorului b pentru care expresia:

$$S(b) = \|Y - Xb\|^2 = \sqrt{\sum_{i=1}^p (Y_i - b_0 - b_1 \cdot e^{-X_i} - b_2 \cdot X_i e^{-X_i})^2},$$

este minimă, presupunând că matricea X are rangul n .

Astfel $X' \cdot X$ este inversabilă și avem:

$$b = (X' \cdot X)^{-1} \cdot X' \cdot Y.$$

Exemplu

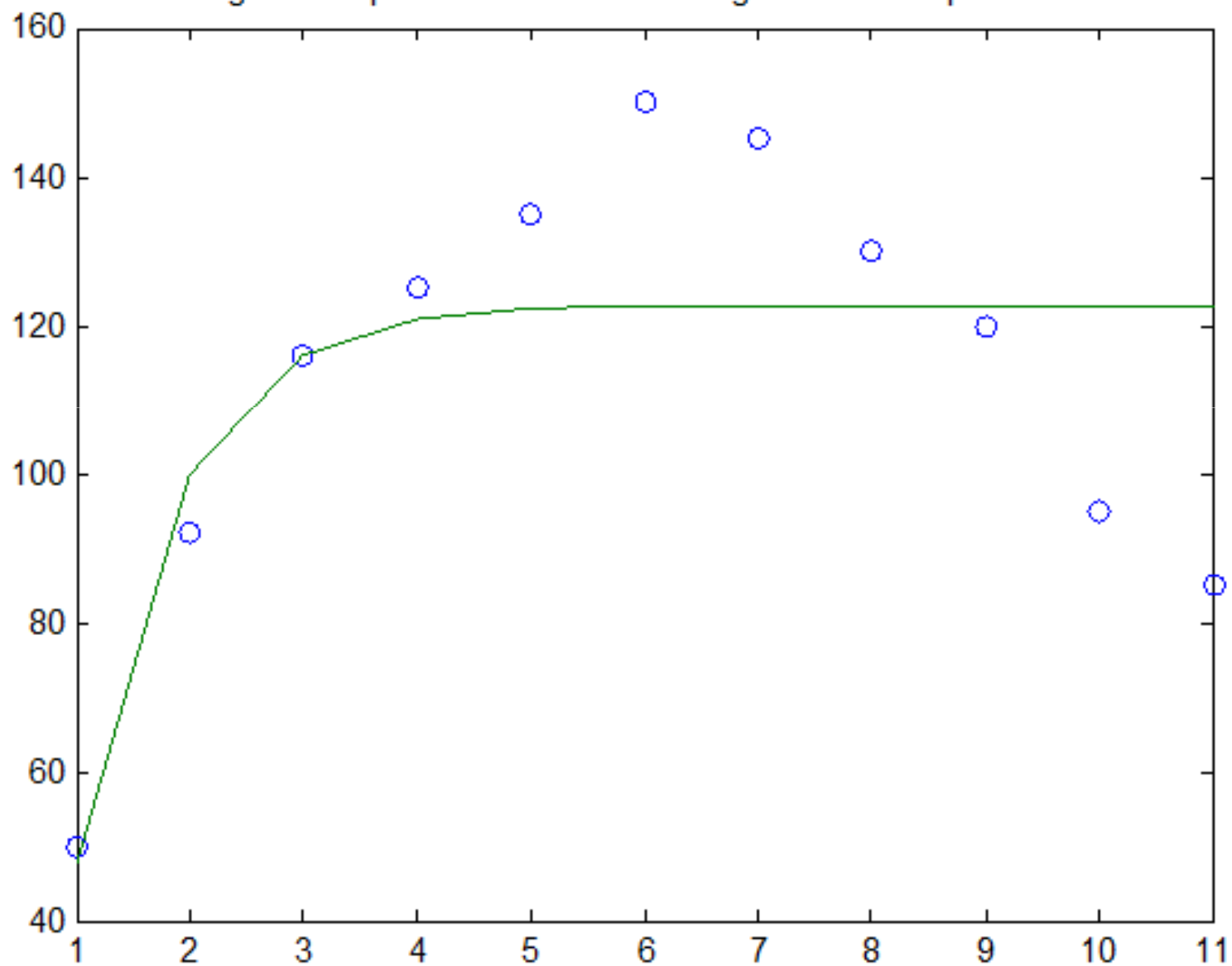
Reluăm exemplul anterior construind cu o funcție de regresie mixt exponențială de forma:

$$f(X) = b_0 + b_1 \cdot e^{-X} + b_2 \cdot X e^{-X} .$$

Pentru a determina parametrii b_0, b_1, b_2 vom utiliza metoda celor mai mici pătrate:

```
>> t=1:11; P=[50 92 116 125 135 150 145 130 120 95 85];  
>> i=1;B=1;for i=2:11 B=[1 B];end  
>> for i=1:11 t1(i)=exp(-t(i));end  
>> for i=1:11 t2(i)=t(i)*exp(-t(i));end  
>> X=[B' t1' t2'];b2=X\P'  
b2 =  
    122.7532  
   -236.5601  
     33.0078  
>> C=b2(1)+b2(2)*exp(-t)+b2(3)*t.*exp(-t);plot(t,P,'O',t,C)
```

diagrama imprastierii si functia de regresie mixt-exponentiala

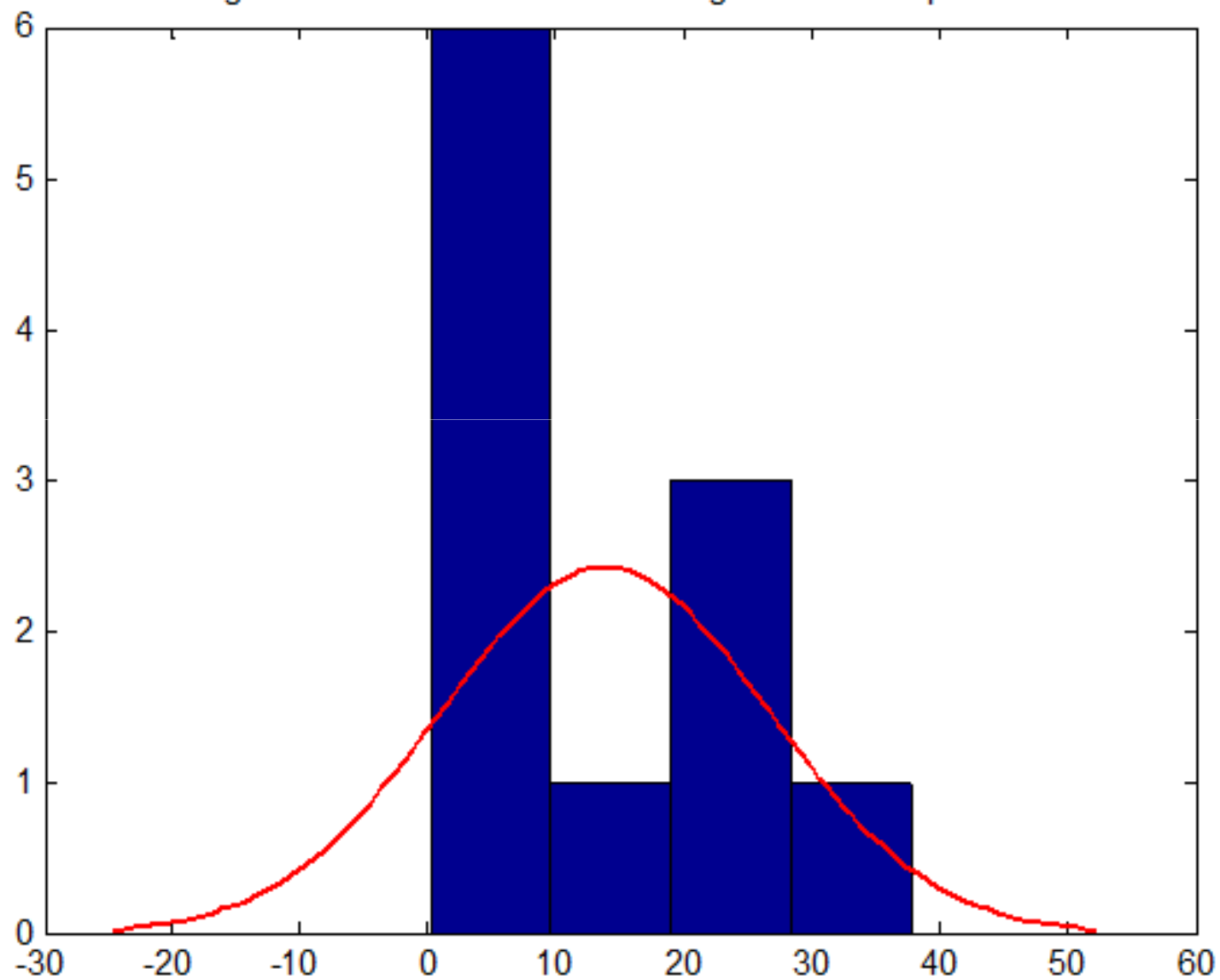


```
>> rez=abs(P-C); MSE=mean((P-C).^2)
MSE =
    340.1303
```

Comparând media pătratelor reziduurilor cu cele obținute în cazul folosirii regresiei polinomiale, se observă că rezultatele obținute sunt mai puțin bune, lucru evident și din figură (diagrama împrăștierii vs. curba de regresie mixt exponențială).

```
>>histfit(rez)
```

histograma reziduurilor- functia de regresie mixt exponentiala



Problema utilizării regresiei neliniare la un set de date este o procedură euristică, fiind necesare mai multe experimente și folosirea cunoștințelor specifice domeniului pentru obținerea unor rezultate bune.

Regresia liniară multiplă

Metoda regresiei liniare se poate extinde de la cupluri de două variabile la mai multe variabile prin metoda *regresiei liniare multiple*, caz în care avem o variabilă dependentă și mai multe variabile predictive.

Să considerăm un set de n variabile aleatoare X_1, X_2, \dots, X_n , considerate independente în sensul că, din punct de vedere probabilistic, reprezintă variabilele ce guvernează n factori predictivi care acționează independent unul față de altul.

Variabilă dependentă (prognozată), notată cu Y , este variabilă aleatoare ale cărei valori vor fi estimate pe baza analizei regresive liniare multiple, plecând de la valorile factorilor predictivi.

5. Exemplu

Prezentăm baza de date completă ce conține informații despre prețul de vânzare a 15 apartamente dintr-un cartier Ω al unui oraș Z , în trimestrul al III-lea din 2011 (parțial, această bază a fost prezentată în exemplul nr 3).

supr	30	32	35	48	49	50	55	60	65	70	75	80	82	85	90
preț	30	29	31	36	38	40	48	47	53	49	58	60	59	71	65
nr cam	1	1	1	2	2	2	2	3	3	3	3	3	4	4	4
nr bai	1	1	1	1	1	1	1	1	2	2	2	2	1	2	2

Variabilele predictive sunt caracteristicile apartamentelor (suprafață, număr camere, număr băi) și variabila dependentă este reprezentată de prețul apartamentului.

Variabilele predictive sunt caracteristicile apartamentelor (suprafață, număr camere, număr băi) și variabila dependentă este reprezentată de prețul apartamentului.

```
>> Y=[30 29 31 36 38 40 48 47 53 49 58 60 59 71 65];  
>> X1=[ 30 32 35 48 49 50 55 60 65 70 75 80 82 85 90];  
>> X2=[1 1 1 2 2 2 2 3 3 3 3 3 4 4 4];  
>> X3=[ 1 1 1 1 1 1 1 1 2 2 2 2 1 2 2 ];  
>> A=[X1' X2' X3' Y'];[r,p]=corrcoef(A)
```

r =

1.0000	0.9687	0.7321	0.9729
0.9687	1.0000	0.6378	0.9408
0.7321	0.6378	1.0000	0.7436
0.9729	0.9408	0.7436	1.0000

p =

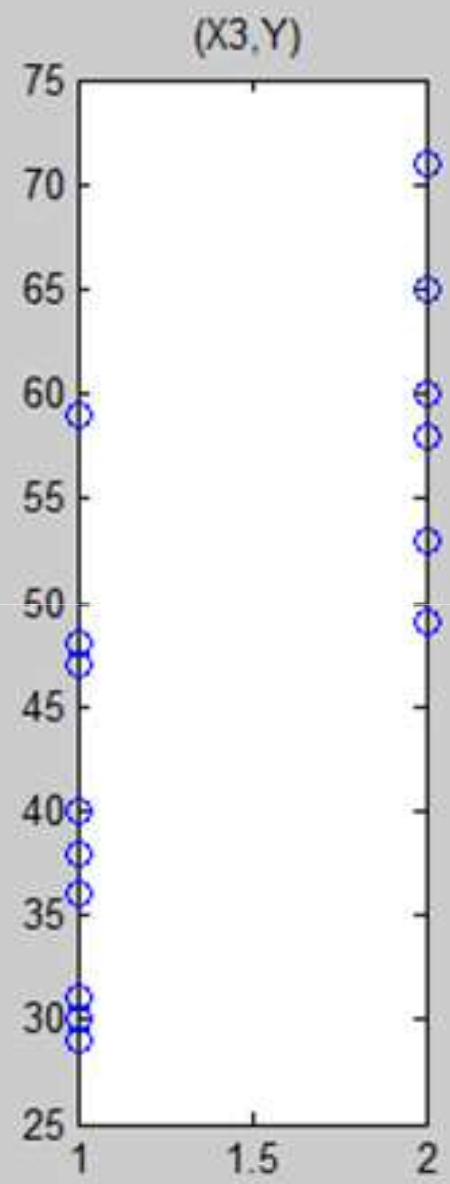
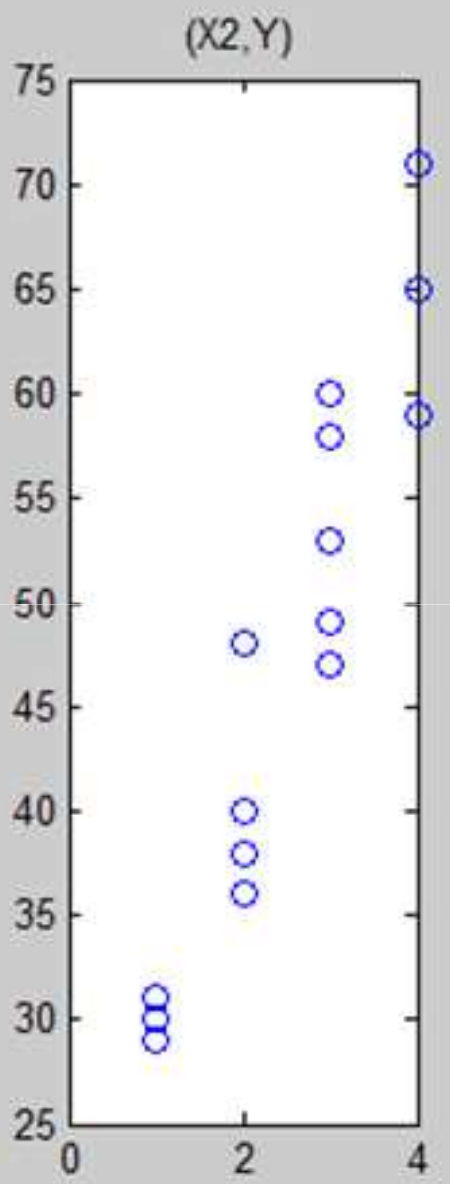
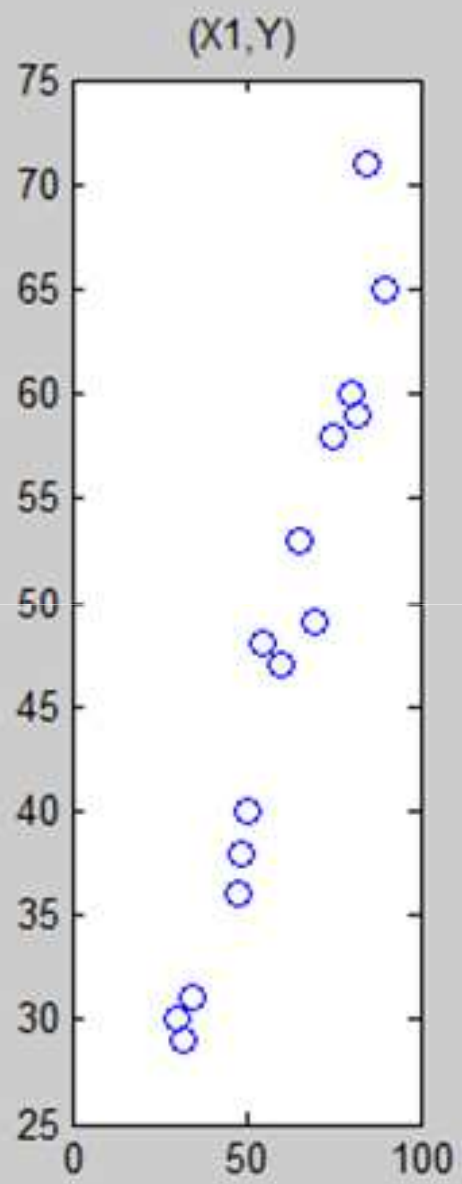
1.0000	0.0000	0.0019	0.0000
0.0000	1.0000	0.0105	0.0000
0.0019	0.0105	1.0000	0.0015
0.0000	0.0000	0.0015	1.0000

Legătura între variabila răspuns și variabilele predictive este puternică (în special X1 și Y, respectiv X2 și Y, și semnificativă.

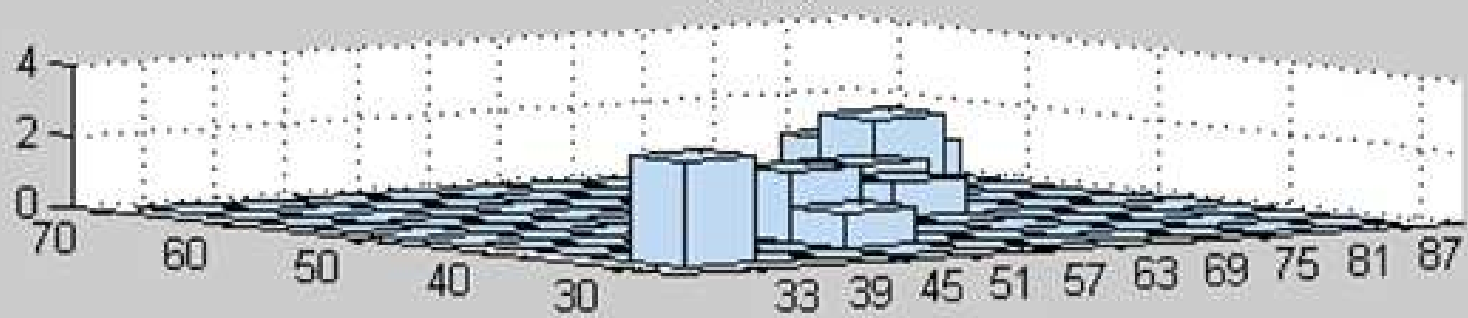
Nefiind interesați de legăturile existente între variabilele predictivice, vom desena diagramele împrăștierii și histogramele corespunzătoare doar în cazurile (variabilă predictivă, predictor):

```
>> subplot(131);plot(X1,Y,'o');  
>> subplot(132);plot(X2,Y,'o');  
>> subplot(133);plot(X3,Y,'o');
```

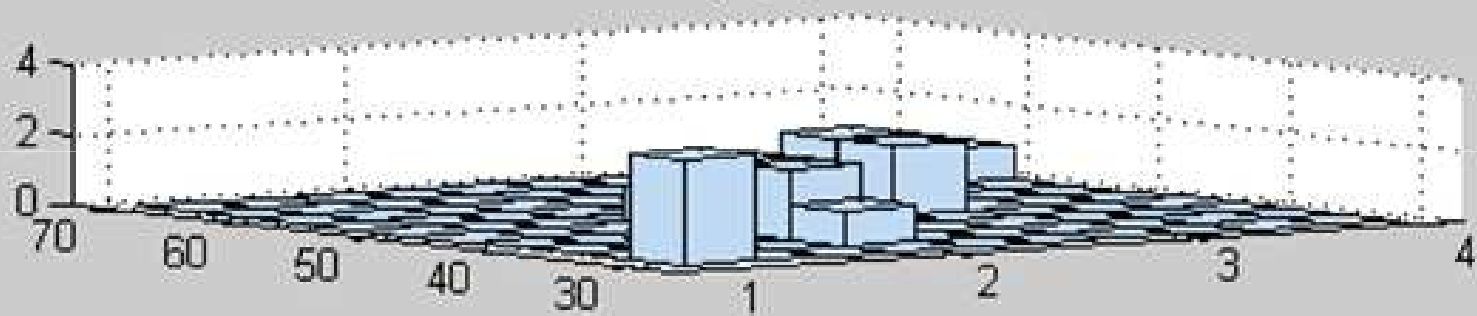
```
>> subplot(311);A1=[X1',Y']; hist3(A1);  
>> subplot(312);A2=[X2',Y']; hist3(A2);  
>> subplot(313);A3=[X3',Y']; hist3(A3);
```



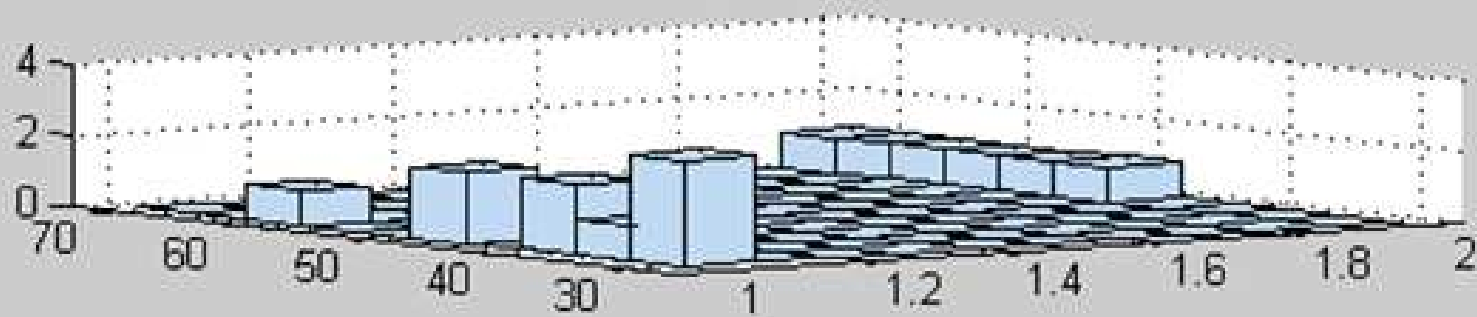
(X1,Y)



(X2,Y)



(X3,Y)



Practic, deducând ecuația de regresie pe baza datelor culese de la acest lot de apartamente, se poate prognoza prețul de vânzare al unui nou apartament, introducând în ecuație valorile particulare ale acestuia pentru fiecare variabilă predictivă în parte.

Așadar, pentru regresia multi-liniară, dispunem de n variabile predictive sau explicative, a căror efecte sunt independente, în principiu, unele de altele, și de o variabilă prognozată, numită *răspuns*, pe care vrem să o deducem, cunoscând valorile celor n variabile predictive.

Legătura între variabila răspuns Y și variabilele predictive X_i trebuie să fie una liniară.

Odată stabilită validitatea condițiilor de aplicare a metodei, se trece la obținerea ecuației de regresie liniară multiplă, care este de forma:

$$Y = b_0 + b_1 \cdot X_1 + b_2 \cdot X_2 + \dots + b_n \cdot X_n ,$$

unde b_1, b_2, \dots, b_n se numesc *coeficienții de regresie*, iar b_0 *interceptor*.

Metoda regresiei multiple standard reprezintă o generalizare a regresiei simple, obținându-se hiperplanul care trece prin „norul” multidimensional al tuturor variabilelor.

Metoda regresiei multiple standard reprezintă o generalizare a regresiei simple, obținându-se hiperplanul care trece prin „norul” multidimensional al tuturor variabilelor.

Pentru a determina coeficienții de regresie și interceptorul, vom aplica *metoda celor mai mici pătrate*. Presupunem că datele existente satisfac ecuațiile:

$$Y_i = b_0 X_0^i + b_1 \cdot X_1^i + b_2 \cdot X_2^i + \dots + b_n \cdot X_n^i + \varepsilon_i, \quad 1 \leq i \leq p,$$

unde:

- $X_j^i, 1 \leq j \leq n, 1 \leq i \leq p$ și $Y_i, 1 \leq i \leq p$ sunt numere reale cunoscute (reținem că $p > n$).
- $X_0^i = 1, 1 \leq i \leq p$.
- ε_i sunt variabile aleatoare necunoscute.

Sub formă matriceală modelul devine:

$$Y = X \cdot b + \varepsilon ,$$

unde:

- $Y = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_p \end{pmatrix}$ este un vector aleator p -dimensional,

- $X = \begin{pmatrix} 1 & X_1^1 & \dots & X_n^1 \\ 1 & X_1^2 & \dots & X_n^2 \\ \dots & \dots & \dots & \dots \\ 1 & X_1^p & \dots & X_n^p \end{pmatrix}$ este matricea obținută prin

concatenarea celor p variabile X_i ,

- b este vectorul $(n+1)$ -dimensional ale cărui componente sunt parametri necunoscuți ai modelului,
- ε este vectorul p -dimensional al erorilor.

Deoarece X este o matrice de rang n , $X' \cdot X$ este inversabilă și astfel:

$$b = (X' \cdot X)^{-1} \cdot X' \cdot Y .$$

Exemplu

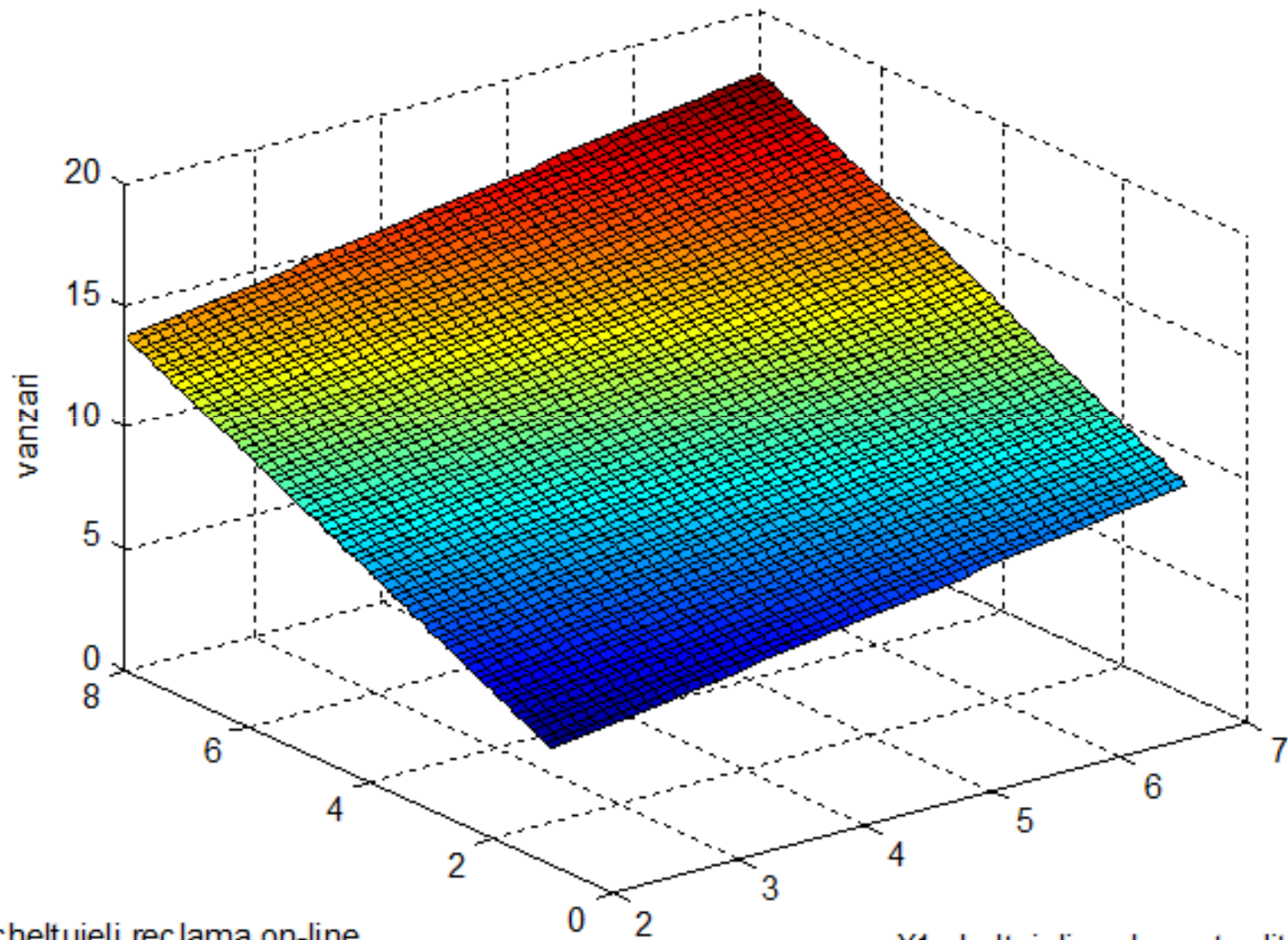
Reluăm exemplul nr 2 cu firma de IT, unde se presupune că Y, rezultatul vânzărilor, depinde de X1, suma cheltuită pentru reclama tradițională la TV și în ziare, și X2, suma cheltuită pentru reclama prin Internet. Vom scrie ecuația planului de regresie și îl vom desena.

```
>> Y=[5 7 7 5 9 11 12 13 15 17 18 19];
>> X1=[2 2.5 3 1.5 4 4.5 5 5.5 6 7.5 8 7];
>> X2=[1 2 2.5 1.5 3.5 4 5 6 6.5 8 8 7.5];
>> i=1; B=[1]; for i=2:12 B=[B [1]];end
>> X=[ B' X1' X2' ];
>> b=X\Y'
b =
    1.9921
    0.7672
    1.2748
>> [X1,X2]=meshgrid(2:.1:7,1:.1:8);
>> surf(X1,X2, b(1)+b(2)*X1+b(3)*X2)
```

Reamintim că în MATLAB nu este nevoie să utilizăm formula

$b = (X' \cdot X)^{-1} \cdot X' \cdot Y$, obținând același rezultat cu $b=X\backslash Y'$:

planul de regresie - exemplul nr 2



X2 cheltuieli reclama on-line

X1 cheltuieli reclama traditionala

Dacă se utilizează 85000 EUR pentru publicitate clasică și 60.000 EUR pentru publicitate prin Internet, suntem interesați la cât se prognozează cifra vânzărilor.

Vom apela la Symbolic Math:

```
>> syms X1 X2  
>> f=b(1)+b(2)*X1+b(3)*X2;  
>> f1=subs(f,[X1,X2],[8.5,6.0])  
f1 =  
    16.1616
```

Exemplu

Reluăm exemplul nr 5, determinând coeficienții hiperplanului de regresie:

```
>> Y=[30 29 31 36 38 40 48 47 53 49 58 60 59 71 65];  
>> X1=[ 30 32 35 48 49 50 55 60 65 70 75 80 82 85 90];  
>> X2=[1 1 1 2 2 2 2 3 3 3 3 3 4 4 4];  
>> X3=[ 1 1 1 1 1 1 1 1 2 2 2 2 1 2 2 ];  
>> i=1; B=[1]; for i=2:15 B=[B [1]];end  
>> X=[ B' X1' X2' X3']  
>> b=inv(X'*X)*X'*Y'
```

```
b =  
 7.8582  
 0.5769  
 0.8070  
 2.0385
```

Ecuția hiperplanului de regresie este:

$$Y = 7.8582 + 0.5769 \cdot X_1 + 0.8070 \cdot X_2 + 2.0385 \cdot X_3$$

Folosind ecuația hiperplanului de regresie, putem estima valoarea de vânzare a unui apartament din cartierul Ω al orașului Z : să presupunem un apartament cu 3 camere, două băi, cu suprafața totală de 74 m^2 :

```
>> syms X1 X2 X3
>> a=subs(b(1)+b(2)*X1+b(3)*X2+b(4)*X3, [X1, X2, X3], [74,3,2])
a =
  57.0453
```

În cazul regresiei liniare multiple apare problema ierarhizării predictorilor în scopul păstrării numai a aceluia care, într-adevăr, au o influență semnificativă asupra variabilei efect, renunțând astfel la cei ne semnificativi. Astfel se iau cei care sunt puternic și mai ales semnificativ corelați cu variabila răspuns.

Reducerea numărului atributelor unei baze de date este o problemă serioasă, care implică tehnici variate.

6. Exemplu

În exemplul următor vom reduce numărul atributelor, vom determina hiperplanele de regresie în cazul numărului total se atribute și în cazul numărului redus de atribute și vom calcula sumă pătratelor reziduurilor în cele două cazuri

Să considerăm analiza regresivă multiplă utilizată în exprimarea (predicția) *indexului* de rezistență a mușchiului respirator PEmax (exprimat în cm H₂O) în funcție de variabilele predictoare reprezentate de înălțime (H -cm), greutate (G-kg), vârstă (ani), sex, procentul masei corporale (%), BMP), volumul respirator forțat per secundă (FEV₁), volumul rezidual (RV), capacitatea funcțională reziduală (FRC) și capacitatea totală a plămânului (TLC), pentru un lot de 25 bolnavii cu fibroză cistică . Variabila dependentă a acestui model este reprezentată de indexul de rezistență a mușchiului respirator (PEmax). Tabelul următor sintetizează toate aceste caracteristici.

Vårstã	Sex	H	G	BMP	FEV ₁	RV	FRC	TLC	PE _{max}
7	0	109	13.1	68	32	258	183	137	95
7	1	112	12.9	65	19	449	245	134	85
8	0	124	14.1	64	22	441	268	147	100
8	1	125	16.2	67	41	234	146	124	85
8	0	127	21.5	93	52	202	131	104	95
9	0	130	17.5	68	44	308	155	118	80
11	1	139	30.7	89	28	305	179	119	65
12	1	150	28.4	69	18	369	198	103	110
12	0	146	25.1	67	24	312	194	128	70
13	1	155	31.5	68	23	413	225	136	95
13	0	156	39.9	89	39	206	142	95	110
14	1	153	42.1	90	26	253	191	121	90
14	0	160	45.6	93	45	174	139	108	100
15	1	158	51.2	93	45	158	124	90	80
16	1	160	35.9	66	31	302	133	101	134
17	1	153	34.8	70	39	204	118	120	134
17	0	174	44.7	70	49	187	104	103	165
17	1	176	60.1	92	29	188	129	130	120
17	0	171	42.6	69	38	172	130	103	130
19	1	156	37.2	72	21	216	119	81	85
19	0	174	54.6	86	37	184	118	101	85
20	0	178	64.0	86	34	225	148	135	160
23	0	180	73.8	97	57	171	108	98	165
23	0	175	51.5	71	33	224	131	113	95
23	0	179	71.5	95	52	225	127	101	195

Suntem interesați doar de ultima coloană din matricea corelațiilor și matricea valorilor nivelului de semnificație, în această coloană fiind coeficienții de corelație ai predictorilor cu variabila răspuns.

r =

Column 10

0.6135
-0.2886
0.5992
0.6329
0.2295
0.3061
-0.3156
-0.4172
-0.1816
1.0000

p=

Column 10

0.0011
0.1618
0.0015
0.0007
0.2698
0.1367
0.1244
0.0380
0.3849
1.0000

Calculăm hiperplanul de regresie și suma pătratelor reziduurilor utilizând toți predictorii:

```
>> i=1; B=[1]; for i=2:25 B=[1 B];end
>> B2=[B' V' s' H' G' BMP' FEV' RV' FRC' TLC']; b2=B2\PE'
b2 =
 254.9404
 -3.9563
-11.8768
 -0.6649
  3.6155
 -1.7278
  0.4066
  0.2586
 -0.5103
  0.1220
>>r2=PE-b2(1)-b2(2)*V-b2(3)*s-b2(4)*H-b2(5)*G-b2(6)*BMP-b2(7)*FEV-
b2(8)*RV-b2(9)*FRC-b2(10)*TLC;
>> MSE2=mean(r2.^2)
MSE2 =
 419.2110
```

Calculăm hiperplanul de regresie și suma pătratelor reziduurilor utilizând predictorii puternic și semnificativ corelați cu variabila răspuns:
V, H,G, FRC'

```
>> B1=[B' V' H' G' FRC']; b1=B1\PE'
```

```
b1 =
```

```
65.7249
```

```
1.5651
```

```
-0.0609
```

```
0.8295
```

```
-0.0119
```

```
>> r1=PE-b1(1)-b1(2)*V-b1(3)*H-b1(4)*G-b1(5)*FRC;
```

```
>> MSE1=mean(r1.^2)
```

```
MSE1 =
```

```
633.5354
```

Utilizăm ecuația de regresie multiplă pentru a obține valorile variabilei dependente pentru orice valori individuale ale variabilelor explicative. În acest mod, pentru un anumit obiect cu atributele predictive cunoscute, se deduce valoarea atributului necunoscut, considerat ca atribut răspuns (*outcome*).

În cazul exemplului nr 6, pentru un anumit pacient căruia i se cunosc valorile celor nouă parametri medicali predictivi, i se poate prognoza, cu o acuratețe suficientă, valoarea PEmax, prin introducerea în ecuația de regresie a valorilor sale individuale.

Spunem că, astfel, se obține o *valoare prognostic (index prognostic)*, pe baza datelor cunoscute.

În analiza regresivă se va evita un model cu multe variabile predictive provenind dintr-un eșantion mic.

Ca regulă generală, numărul de predictorii nu trebuie să depășească valoarea $n/10$, unde n este volumul eșantionului. În plus dacă există variabile explicative puternic corelate între ele, acestea nu vor fi incluse toate în model, ci doar un reprezentant al lor.

Pentru mai multă siguranță, se verifică capacitatea modelului pe alt eșantion, dacă acest lucru este posibil.

Regresia multiliniara in Matlab

$\mathbf{b} = \text{regress}(\mathbf{y}, \mathbf{X})$ returnează un vector coloană de p componente ce reprezintă coeficienții hiperplanului de regresie definit de variabila răspuns y și de cei p predictorii.

\mathbf{X} este matricea cu n linii și p coloane a celor p predictorii, caracterizați de n attribute, în timp ce \mathbf{y} este un vector coloană cu n componente, care reprezintă cele n răspunsuri observate.

Dacă în \mathbf{X} sau \mathbf{y} există date lipsă, notate NaN, `regress` le ignoră.

În calcule, \mathbf{X} trebuie să conțină o coloană ale cărei elemente sunt 1.

7.Exemplu

Considerăm baza de date din Matlab, numită **carsmall**, ce conține 92 de automobile ce au ca predictorii greutatea și puterea motorului și ca răspuns kilometrajul

```
>> load carsmall
>> x1 = Weight;x2 = Horsepower;y = MPG;
>> A=[x1 x2 y];
>> X = [ones(size(x1)) x1 x2];
>> b = regress(y,X)
b =
  47.7694
 -0.0066
 -0.0420
```

Regresia logistică

În tehnicile regresive prezentate, variabila răspuns care urma a fi dedusă din variabilele explicative, este o variabilă numerică.

În anumite cazuri această variabilă este categorială, având valorile: DA sau NU, caz în care este vorba de un tip special de clasificare.

De exemplu,

- în medicină, această variabilă poate reprezenta diagnosticarea unei boli (DA sau NU) a unui individ, pe baza unor analize clinice și biochimice ale acestuia. Pentru fiecare combinație a diferitelor valori ale unor variabile care ar putea influența riscul de îmbolnăvire, trebuie estimată probabilitatea apariției maladiei.
- altă situație se referă la încadrarea în categoria „fraudă” a unor tranzacții comerciale, pe baza anumitor caracteristici ale lor.

Principiul de bază rămâne același ca și la regresia multiplă, cu deosebirea că în acest caz estimăm o *transformare* a variabilei dependente.

Astfel, dacă variabila dependentă are ca valori binare afirmațiile DA și NU, codate cu valorile 1 și 0, atunci media acesteia reprezintă proporția indivizilor din populație cu caracteristica respectivă.

De exemplu presupunem că avem pentru 100 de indivizi, 73 răspunsuri DA și 27 răspunsuri NU, adică 73 de 1 și 27 de 0. Media acestor valori va fi:

$$(73 * 1 + 27 * 0)/100 = 73/100$$

adică 73% de DA.

Modelul regresiv logistic estimează probabilitatea ca un obiect oarecare din populație să aibă o anumită caracteristică, de exemplu o anumită boală, tentativă de fraudă etc.

O combinație oarecare a predictorilor din ecuația de regresie poate lua valori în afara intervalului $[0, 1]$ și din acest motiv se ia în considerare o transformare a intervalului obținut, în intervalul $[0, 1]$, situație în care valorile variabilei răspuns pot fi considerate probabilități.

Metoda regresiei logistice constă în folosirea transformării *logit* scrisă ca *logit(p)*.

Concret, aici p reprezintă probabilitatea ca un obiect oarecare să aibă caracteristica cerută, deci, $(1 - p)$ va reprezenta probabilitatea ca acesta să nu o aibă.

În exemplele prezentate, p este probabilitatea ca un subiect din populație să aibă cancer, sau să producă o fraudă, iar $(1 - p)$ să nu aibă cancer sau nu producă fraudă.

Raportul $\frac{p}{1 - p}$ se numește *șansă*, iar transformarea:

$$\text{logit}(p) = \ln \frac{p}{1 - p}$$

este numită *logaritmul (log-ul) șansei*.

Pe baza predictorilor X_1, X_2, \dots, X_k , $\text{logit}(p)$ este dat de ecuația de regresie logistică:

$$\text{logit}(p) = \ln \frac{p}{1-p} = b_0 + b_1 X_1 + \dots + b_k X_k.$$

Prin transformarea inversă, se găsește valoarea probabilității p ca un anumit obiect să aibă sau nu o anumită *etichetă* de clasificare:

$$p = \frac{1}{1 + e^{-(b_0 + b_1 X_1 + \dots + b_k X_k)}}.$$

Pentru a calcula coeficienții b_0, b_1, \dots, b_k utilizăm metoda celor mai mici pătrate.

Presupunem că datele existente satisfac relațiile:

$$\text{logit}(p_i) = \ln \frac{p_i}{1-p_i} = b_0 + b_1 X_1^i + \dots + b_k X_k^i + \varepsilon_i, 1 \leq i \leq m,$$

unde:

- $X_j^i, 1 \leq j \leq k, 1 \leq i \leq m$ sunt numere reale cunoscute;

reținem că $m > k$,

- $p_i, 1 \leq i \leq m$ reprezintă probabilitatea ca obiectul indice i să aibă o anumită etichetă de clasificare,
- ε_i sunt variabile aleatoare necunoscute.

Sub formă matriceală, modelul devine: $\text{logit}(p) = X \cdot b + \varepsilon$, unde

$$\text{- logit}(p) = \begin{pmatrix} \ln \frac{p_1}{1-p_1} \\ \ln \frac{p_2}{1-p_2} \\ \vdots \\ \ln \frac{p_m}{1-p_m} \end{pmatrix} \text{ este un vector aleator } m\text{-dimensional,}$$

$$\text{- } X = \begin{pmatrix} 1 & X_1^1 & \dots & X_k^1 \\ 1 & X_1^2 & \dots & X_k^2 \\ \dots & \dots & \dots & \dots \\ 1 & X_1^m & \dots & X_k^m \end{pmatrix} \text{ este matricea obținută prin}$$

concatenarea celor m variabile X_i ,

- b este vectorul $(k+1)$ -dimensional ale cărui componente sunt parametri necunoscuți ai modelului,

- ε este vectorul m -dimensional al erorilor.

Metoda celor mai mici pătrate constă în aflarea vectorului b pentru care expresia:

$$S(b) = \|\text{logit}(p) - Xb\|^2 = \sqrt{\sum_{i=1}^m \left(\ln \frac{p_i}{1-p_i} - b_0 - b_1 X_1^i + \dots - b_k X_k^i \right)^2},$$

este minimă, presupunând că matricea X are rangul k .

Astfel $X' \cdot X$ este inversabilă și avem:

$$b = (X' \cdot X)^{-1} \cdot X' \cdot \text{logit}(p).$$

8.Exemplu

Pentru stabilirea legăturilor existente între apariția hipertensiunii și următorii factori de risc: *fumat*, *obezitate*, *vârstă* (sub sau deasupra valorii alese de 40 ani –codat 0 sau 1), vom aplica o analiză regresivă logistică. În vederea utilizării sale, să considerăm următoarea situație statistică, descrisă în tabelul de mai jos

fumat	obezitate	vârstă	Nr. subiecți	Nr subiecți hipertensivi
0	0	0	60	5(8%)
1	0	0	17	2(11%)
0	1	0	8	1(13%)
0	0	1	187	35(19%)
1	1	0	2	0(0%)
1	0	1	60	5(8%)
0	1	1	51	15(29%)
1	1	1	23	8(35%)
			Total=433	total=79 (18%)

Plecând de la acest tabel general, vom construi un tabel de lucru, care rezumă toate posibilitățile existente în felul următor:

- în primele două linii ale sale se iau în considerație cele două situații posibile de clasificare ale indivizilor (hipertensiune = 1, normal = 0), pentru situația în care variabilele explicative fumat, obezitate și vârstă au valorile egale cu 0, indivizi care nu fumează, nu sunt obezi și au vârsta sub 40 de ani.

Mai departe se consideră toate combinațiile posibile ale predictorilor.

De reținut: trebuie luate în considerație toate cazurile posibile privind combinațiile variabilelor explicative.

fumat	obezitate	vârstă	Nr subiecți	hipertensiune
0	0	0	55	0
0	0	0	5	1
1	0	0	15	0
1	0	0	2	1
0	1	0	7	0
0	1	0	1	1
0	0	1	152	0
0	0	1	35	1
1	1	0	2	0
1	1	0	0	1
1	0	1	55	0
1	0	1	5	1
0	1	1	36	0
0	1	1	15	1
1	1	1	15	0
1	1	1	8	1

```

>> A=[ 0 0 0; 1 0 0; 0 1 0; 0 0 1; 1 1 0; 1 0 1; 0 1 1; 1 1 1 ];
>> i=1; B=[1]; for i=2:8 B=[1 B];end
>> X=[B' A];
>> logitp =[log(5/55) log(2/15) log(1/7) log(35/152) 0 log(5/55) log(15/36)
log(8/15)];
>> b=inv(X'*X)*X'*logitp'
b =
-2.3991
-0.4116
1.2073
0.2471

```

Ecuatia de regresie logistica corespunzatoare este data de:

$$\text{logit}(p) = -2.3991 - 0.4116 \cdot \text{fumat} + 1.2073 \cdot \text{obezitate} + 0.2471 \cdot \text{varsta} .$$

Se observă utilitatea acestei ecuații atunci când trebuie rezolvată problema clasificării unui individ oarecare ca fiind predispus sau nu la hipertensiune, pe baza atributelor (caracteristicilor): fumat, obezitate și vârstă.

Astfel, se consideră valorile concrete ale aceluși individ privind cei trei factori de risc enumerați mai sus, care se vor introduce în ecuația de regresie.

Va rezulta valoarea $\text{logit}(p)$ corespunzătoare, care va estima riscul ca individul respectiv să fie sau nu predispus la hipertensiune.

De exemplu:

- dacă respectivul este nefumător, nu este obez și are sub 40 de ani, rezultă:

$$\text{logit}(p) = -2.3991,$$

deci, probabilitatea să aibă hipertensiune este egală cu 8.32%.

- dacă pacientul este nefumător, obez și are peste 40 de ani, rezultă

$$\text{logit}(p) = -0.9447,$$

și astfel probabilitatea ca să aibă hipertensiune este egală cu 28%.

Ecuatia de regresie logistica se mai poate folosi și în scopul de a compara probabilitățile de predicție a hipertensiunii pentru diferite grupuri, de exemplu pentru cei cu vârsta sub 40 ani față de cei peste 40 ani.

Astfel, codând ca mai sus, cu 0 subiecții sub 40 ani și cu 1 pe cei peste 40 ani, și utilizând ecuația de regresie de mai sus, obținem cele două variante ale sale:

$$\text{logit}(p_{\text{var sta} > 40}) = - 2.3991 - 0.4116 \cdot \text{fumat} + 1.2073 \cdot \text{obezitate} + 0.2471$$

și:

$$\text{logit}(p_{\text{var sta} < 40}) = - 2.3991 - 0.4116 \cdot \text{fumat} + 1.2073 \cdot \text{obezitate} \quad |$$

Se obține astfel că:

$$\text{logit}(p_{\text{var sta} > 40}) - \text{logit}(p_{\text{var sta} < 40}) = 0.2471.$$

Rezultă că raportul șanselor hipertensiunii asociat cu nivelul de *vârstă* considerat mai sus (de 40 ani) este $e^{0.2471} = 1.2803$, valoare care poate fi interpretată astfel: *riscul de a face hipertensiune peste vârsta de 40 de ani este de 1.28 de ori mai mare decât sub această vârstă.*

În exemplul nr 8 am considerat cazul în care variabilele explicative sunt, de asemenea, categoriale.

Nu există niciun impediment în a considera modelul regresiv logistic și în cazul valorilor numerice ale acestora. Singura variabilă care trebuie să fie categorială este doar variabila dependentă, adică cea care dă eticheta de clasă (categorie).

9. Exemplu

Fibroza ficatului este un indicator important al bolilor hepatice.

Tehnic fiecare dintre cele 6 stadii al fibrozei corespunde unui anumit scor MetavirF, începând cu MetavirF = 0 nu există fibroză și încheind cu MetavirF = 4 ciroza.

În ultimii ani se caută o manieră de evaluare neinvazivă a fibrozei, evitând biopsia prin combinarea testelor biochimice cu metode de imagistică medicală.

Fibroscanul este o tehnică de ultimă generație, ce permite cuantificarea fibrozei hepatice pe baza analizei deplasării unei unde elastice de soc care se se propaga în țesutul hepatic.

In contractul nr 2076/2007 s-a realizat o baza de date formată din 722 de pacienți de la Clinica medicală nr 3, UMF Cluj-Napoca, pacienți cu diferite stadii ale fibrozei ficatului. Fiecare pacient este caracterizat de 26 de parametrii clinici și biochimici, de valoarea rigidității ficatului (stiffness_data de fibroscan) și de stadiul fibrozei, obținut prin biopsie. Această bază de date va constitui *mulțimea de antrenament*.

Prezentăm 5 pacienți, în diferite stadii ale fibrozei, cu parametrii clinici și biochimici și valorile rigidității ficatului.

	P ₁	P ₂	P ₃	P ₄	P ₅
Metavir F	0	1	2	3	4
Stiffness	4.9	5.3	6	10.6	27
Sex	2	2	1	1	1
Age	45	56	31	37	46
BMI (body mass index)	23.1	24.3	34	25.8	30.1
Glycemia	108	110	96	99	84
Triglycerides	349	54	154	134	93
Cholesterol	220	133	197	162	152
Aspartate aminotransferase	27	62	46	48	105
Alanin aminotransferase	42	61	117	60	167
Gamma glutamyl transpeptidase	117	19	41	46	187
Total bilirubin	0.4	1.2	0.9	0.5	1
Alkaline phosphatase	156	283	248	184	246
Prothrombin index	92.2	102	107	79.1	42.7
Tqs (tocopheryl quinones)	15.3	15.2	15	18.4	25.4
INR(prothrombin time ratio)	1.06	0.98	0.95	1.19	1.85
Prolonged activated partial thromboplastin time	27.9	30.9	28.4	29	30.9
Haematids (erythrocytes)	4.64	4.91	4.81	4.1	5.11
Hemoglobin	13.5	14.5	14.9	14.1	15.2
Hematocrit	40.9	40.9	42.6	40.5	45.2
Medium eritrocity volume	88.2	83.2	88.6	89.9	88.5
Average eritrocitary hemoglobin	29.2	29.5	31	33.4	29.7
Av. concentr. of hemoglobin in a red blood cell	33.1	35.6	35	37.2	33.5
Leukocytes	6.3	83.2	88	5.98	88.6
Thrombocytes	236	29.5	31	130	29.7
Sideraemia	70	35.5	35	132	33.5

Am aplicat un model de regresie logistică, în care am găsit probabilitatea ca un pacient să fie

- în stradiul $\text{MetavirF} = 0$ vs $\text{MetavirF} \in \{1,2,3,4\}$
- în stradiul $\text{MetavirF} \in \{0,1\}$ vs $\text{MetavirF} \in \{2,3,4\}$
- în stradiul $\text{MetavirF} \in \{0,1,2\}$ vs $\text{MetavirF} \in \{3,4\}$
- în stradiul $\text{MetavirF} \in \{0,1,2,3\}$ vs $\text{MetavirF} = 4$

studiu care se poate face și cu ajutorul curbelor ROC, despre care vom vorbi ulterior.

Pentru diagnosticul automat rezultatul nu este deosebit, fiind interesați de clasificarea pacienților într-un anumit stadiu, lucru ce s-a realizat cu alte metode de clasificare.

Ceea ce ne interesează este că putem face o reducere a atributelor, lucru important atât pentru rularea programului, cât și pentru serioase economii în domeniul medical.

Acele atribute care nu au o legătură puternică și semnificativă cu stadiul fibrozei vor fi eliminate.

Vom stabili probabilitatea P_1 ca un pacient să fie în în stradiul

MetavirF = 0 vs MetavirF $\in \{1,2,3,4\}$.

$$\text{logit}(P_1) = -29.27 - 1.21 * \text{stiffness} - 24.1 * \text{trigliceride} - 0.24 * \text{hematocrit} + 0.47 * \text{vem} - 1.65 * \text{hem} + 1.17 * \text{chem} + 362.28 * \text{tg-uln}$$

Nivelul de semnificație al parametrilor predictivi ce apar in această ecuație:

	p-level
stiffness	0.0005
trigliceride	0.0174
hematocrit	0.0424
vem	0.0014
hem	0.0001
chem	0.0001
tg_uln	0.0175

Vom stabili probabilitatea P_2 , ca un pacient să fie în stadiul MetavirF $\in \{0,1\}$ vs MetavirF $\in \{2,3,4\}$

$$\text{logit}(P_2) = -7.70 - 0.42 * \text{stiffness} - 0.32 * \text{hem} + 0.34 * \text{chem} + 0.004 * \text{trombocite}$$

Nivelul de semnificație al parametrilor ce apar în ecuație:

	p-level
stiffness	0.0000
hem	0.0022
chem	0.0001
trombocite	0.0445

Stabilim probabilitatea P_3 , ca un pacient să fie în stadiul

MetavirF $\in \{0,1,2\}$ vs MetavirF $\in \{3,4\}$

$$\text{logit}(P_3) = 4.48 - 0.35 * \text{stiffness} + 0.60 * \text{trigliceride} - 2.87 * \text{alat} + 0.010 * \text{sideremia} + 120.99 * \text{alat_uln} - 132.84 * \text{cst_uln}$$

Nivelul de semnificație al parametrilor predictivi ce apar in ecuație

	p-level
stiffness	0.0000
trigliceride	0.0165
alat	0.0206
sideremie	0.0155
alat_uln	0.0203
cst_uln	0.0162

Stabilim probabilitatea P_4 ca un pacient să fie în stradiul

MetavirF $\in \{0,1,2,3\}$ vs MetavirF = 4

$$\text{logit}(P_4) = 3.03 - 0.23 * \text{stiffness} + 0.047 * \text{ind.d} + 0.42 * \text{hemoglobin} - 0.39 * \text{hem} + \\ + 0.009 * \text{trombocite} - 0.472 * \text{sex}$$

Nivelul de semnificație al parametrilor predictivi ce apar in ecuație:

	p-level
stiffness	0.0000
ind.d	0.0095
hemoglobin	0.0469
hem	0.0102
trombocite	0.0052
sex	0.0323

În concluzie un model de regresie logistică permite prognozarea unui anumit răspuns în funcție de o serie de factori predictivi, realizând în acest mod o clasificare a unor obiecte în două clase distincte în raport cu acest răspuns.

Pe de altă parte, prin utilizarea unui model de regresie logistică putem elimina o serie de atribute cu nivel de semnificație scăzut, în problema noastră, deci regresia logistică este și o tehnică de *Features selection*.

Analiza supraviețuirii (Survival Analysis)

Una dintre ramurile importante ale Statisticii, cu aplicații deosebit de interesante mai ales în domeniul medical și al sistemelor mecanice, este *analiza supraviețuirii (Survival Analysis)*.

În medicină este cunoscută chiar sub acest nume, în timp ce în științele ingineresti se numește *teoria siguranței în funcționare (Reliability Theory)*, iar în economie este cunoscută ca *analiza duratei (Duration Analysis)*.

Indiferent de context, elementul de bază al acestei teorii este termenul de „moarte”, „eșec”, „cădere”, „absență”, „ieșire din funcțiune” etc., care este privit ca un *eveniment* în analiza supraviețuirii.

Pentru a înțelege mai bine despre ce este vorba vom considera cazul medical, de la care de fapt s-a și pornit în dezvoltarea acestei ramuri statistice.

Astfel, problema analizei timpilor de supraviețuire, așa cum o indică și numele, se referă, în principiu, la supraviețuirea unui pacient în urma unei operații serioase, tratamentului unei anumite boli cu sfârșit letal, e.g. cancer, SIDA etc. Practic, se înregistrează timpul din momentul începutului procesului medical (operație, tratament etc.) până în momentul decesului, timpul respectiv fiind numit *timp de supraviețuire*, iar analiza sa făcând obiectul analizei supraviețuirii.

Analiza supraviețuirii utilizează așa-numitele *date cenzurate*, adică unele observații sunt incomplete.

Spre exemplu, să ne imaginăm că un grup de pacienți cu cancer sunt urmăriți în cadrul unui experiment o anumită perioadă de timp (*follow-up*). După trecerea acestei perioade, pacienții care au supraviețuit nu mai sunt supravegheați și, atunci când se analizează timpul de supraviețuire, nu se mai știe cu exactitate dacă mai sunt încă în viață. Pe de altă parte, unii pacienți pot părăsi grupul în timpul perioadei de supraveghere, fără a se mai cunoaște situația lor ulterioară.

Datele privind aceste tipuri de pacienți sunt *date cenzurate*.

Să menționăm faptul că, în multe cercetări clinice, prin timp de supraviețuire se înțelege și timpul până la apariția unui simptom, până la revenirea bolii, până la remisiune etc.

În acest context, prin *probabilitate de supraviețuire* se înțelege proporția indivizilor unui grup, supus unui anumit experiment medical comun, care ar putea supraviețui o anumită perioadă de timp (e.g. operație de transplant de inimă, chimioterapie etc.), în anumite circumstanțe date.

Tehnica clasică de calculare a timpului de supraviețuire este următoarea:

Se notează variabila aleatoare care reprezintă timpul de supraviețuire cu X .

Probabilitatea de supraviețuire se calculează prin împărțirea timpului în mici intervale $(0, t_1), \dots, (t_{k-1}, t_k), \dots$, și estimarea apoi a probabilității:

$$P\{X \leq t_n\} = P\{X \leq t_1\}P\{X \leq t_2 | X = t_1\} \dots P\{X \leq t_n | X = t_{n-1}\}.$$

De exemplu probabilitatea ca un pacient să supraviețuiască două zile în urma unui transplant de ficat este produsul dintre probabilitatea ca pacientul să supraviețuiască o zi și probabilitatea ca pacientul să supraviețuiască a doua zi, condiționată de faptul că a supraviețuit prima zi.

O problemă comună atât domeniului cercetării medicale - timpul de supraviețuire, cât și cercetării inginerești - timpul de funcționare sigură a unui mecanism este determinarea efectului unor variabile continue (independente) asupra timpului de supraviețuire, în particular identificarea existenței corelației între predictorii și timpul de supraviețuire.

Metoda clasică de analiza supraviețuirii este bazată pe tehnici ca, de exemplu, *curba de supraviețuire Kaplan-Meier (Kaplan-Meier survival curve)*, *analiza tabelelor de viață (life table analysis)*, *testul logrank (logrank test)*, *raportul hazardului (hazard ratio)*.

Curba de supraviețuire *Kaplan-Meier*

Am stabilit anterior modul de calcul al timpului de supraviețuire

Notând cu p_{100} probabilitatea de a supraviețui a 100-a zi condiționată de faptul că a supraviețuit primele 99 de zile, atunci probabilitatea de a supraviețui 100 de zile după transplant este:

$$p_1 \cdot p_2 \cdot \dots \cdot p_{99} \cdot p_{100}$$

Probabilitatea de a supraviețui a 100-a zi este estimată a fi raportul dintre cei ce au supraviețuit a 100-a și cei ce supraviețuiseră primele 99 de zile.

Probabilitatea p este 1 în zilele când nimeni nu moare, ceea ce simplifică mult calculele, fiind necesar calcularea probabilităților în zilele în care moare cel puțin o persoană.

10.Exemplu

Vom prezenta un experiment pentru prognoza apariției răului de mare:

21 de subiecții au fost plasați într-un simulator, timp de 2 ore.

Timpul de supraviețuire al unui subiect este calculat din minutul 0 până când i se făcea rău.

Unii subiecți au renunțat la experiment înainte de a li se face rău, în timp ce alții au rezistat doua ore fără a avea rău de mare, ambele cazuri ducând la observații cenzurate.

Tabelul următor, care prezintă situația acestui experiment, observațiile cenzurate sunt notate cu * este un așa numit *tabel de viață* (life table):

Am notat: T=timp de supraviețuire (min); p = proporția de supraviețuire

nr	1	2	3	4	5	6	7	8	9	21
T	30	50	50*	51	66*	82	92	120*	120*	120*
p	0.95	0.9		0.85		0.79	0.74				

$$\text{unde } p_1 = \frac{20}{21} = 0.9524, \quad p_2 = \frac{19}{21} = 0.9048, \quad p_4 = \frac{17}{20} = 0.8500,$$

$$p_6 = \frac{15}{19} = 0.7895, \quad p_7 = \frac{14}{19} = 0.7368$$

Proporțiile de supraviețuire se calculează doar pentru cele 5 date necenzurate Vom avea doar 5 estimări ale probabilității de supraviețuire.

La minutul 30 probabilitatea de supraviețuire (fracția de subiecți care au supraviețuit) este $P_{30} = \frac{20}{21} = 0.9524$

La începutul intervalului (30',50') procentul de subiecți rămași în simulator este 0.9524. Pentru a calcula probabilitatea de supraviețuire la minutul 50, conform formulei prezentate vom înmulți P_{30} cu procentul de pacienți rămași în simulator la sfârșitul intervalului (30',50'), adică $\frac{19}{21}$.

$$P_{50} = 0.9524 \cdot \frac{19}{21} = 0.8615.$$

Se ține seama de subiecții cenzurați, atât numărătorul cât și numitorul sunt micșorați, în ziua când apare un astfel de subiect.

Probabilitatea de supraviețuire după 51' este $P_{51} = 0.8615 \cdot \frac{17}{20} = 0.7323$

Probabilitatea de supraviețuire după 82' este $P_{82} = 0.7323 \cdot \frac{15}{19} = 0.5781$

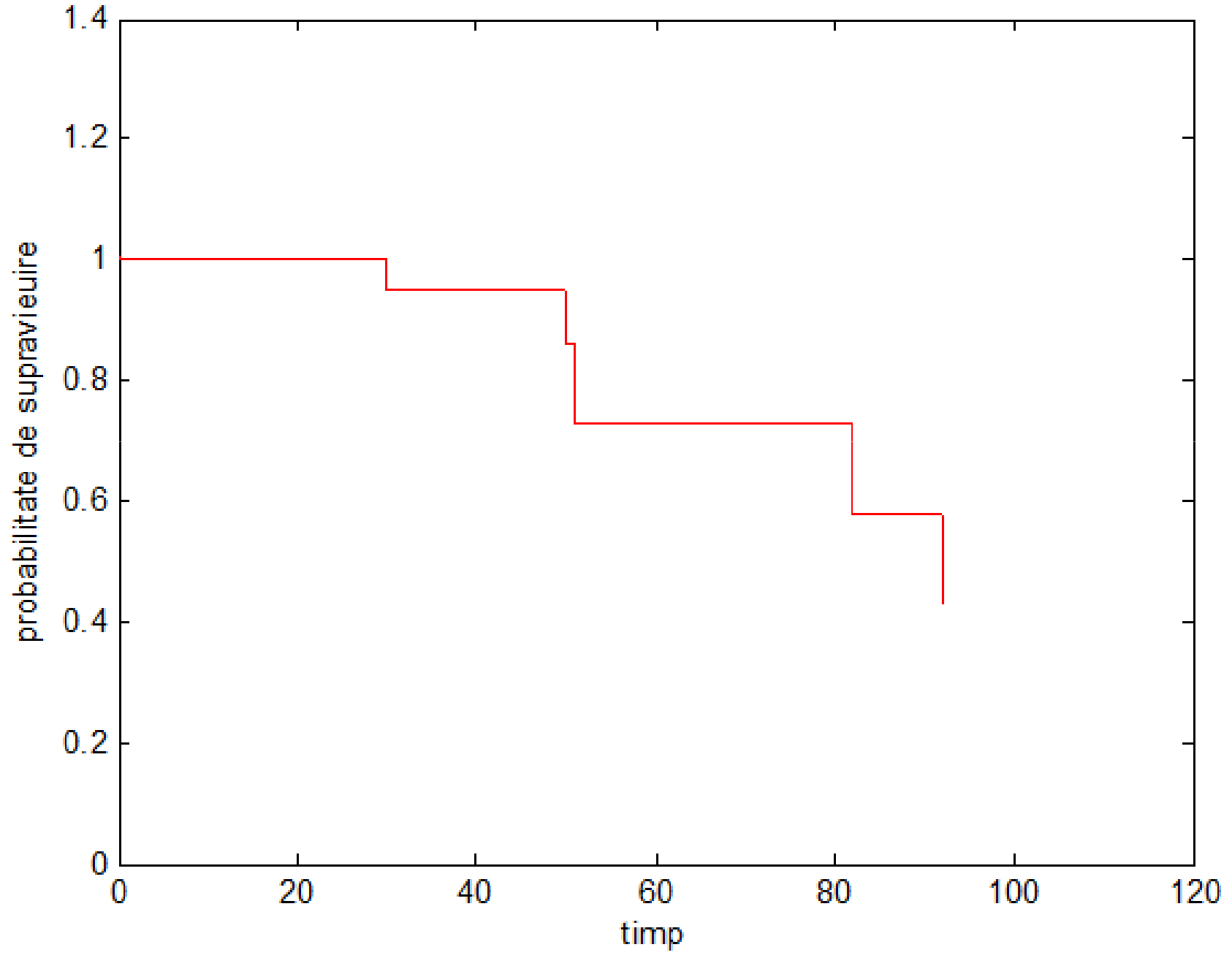
Probabilitatea de supraviețuire după 92' este $P_{92} = 0.5781 \cdot \frac{14}{19} = 0.4260$

Reținem că în primele 30 min, probabilitatea este 1, și nu putem estima supraviețuirea dincolo de ultima observație necenzurată din minutul 92.

Vom desena curba de supraviețuire corespunzătoare:

```
>>x=[0 0.01 30 50 51 82 92];y=[ 1.01 1 0.95 0.86 0.73 0.58 0.43];  
>>x1=[ 92 120];y1=[0 0];  
>> [xb, yb]=stairs(x, y); [x1b y1b]=stairs(x1,y1);  
>>plot(xb,yb,'r',x1b,y1b,'k')
```

curba de supravietuire



Estimator

Un estimator de parametru necunoscut θ a unui model este o funcție ce face să corespundă unui set de observații x_1, x_2, \dots, x_n provenite din model, valoarea $\hat{\theta}$, numită estimare.

Estimatorul este o variabilă aleatoare deoarece valoarea lui depinde de eșantionul obținut în mod aleatoriu.

11. Exemple

— Dacă dorim să calculăm înălțimea medie a copiilor de 10 ani, se poate alege un eșantion de școli situate în medii diferite, și media înălțimii copiilor de 10 ani din acest eșantion, numită și medie empirică, este un estimator al înălțimii medii a copiilor de 10 ani.

Dacă dorim să determinăm procentajul electorilor hotărâți să voteze pentru candidatul A vom face un sondaj pe un eșantion reprezentativ. Procentajul obținut pe acest eșantion este un estimator al procentajului evotanților candidatului A.

De cele mai multe ori un estimator este o medie, o proporție sau o dispersie.

ecdf - funcția cumulativă de distribuție empirică (Matlab)

ecdf - funcția cumulativă de distribuție empirică

Funcția cumulativă de distribuție empirică a unui set de valori numerice este definită pentru orice valoare reală x ca proporția de observații mai mici sau egale cu x .

Este o funcție în scară. Intuitiv, lățimea unei trepte depinde de distanța dintre date consecutive și înălțimea depinde de numărul de valori egale cu x .

Funcția este monotonă. Limita sa la $-\infty$ este 0, iar limita la $+\infty$ este 1.

`[f,x] = ecdf(y)` calculează estimatorul Kaplan-Meier a funcției cumulative de distribuție empirice (cdf = cumulative distribution function), unde

- y este un vector ce are drept componente valorile datelor
- f este vectorul ce are drept componente valorile funcției cumulative de distribuție empirice în X .

12. Exemplu

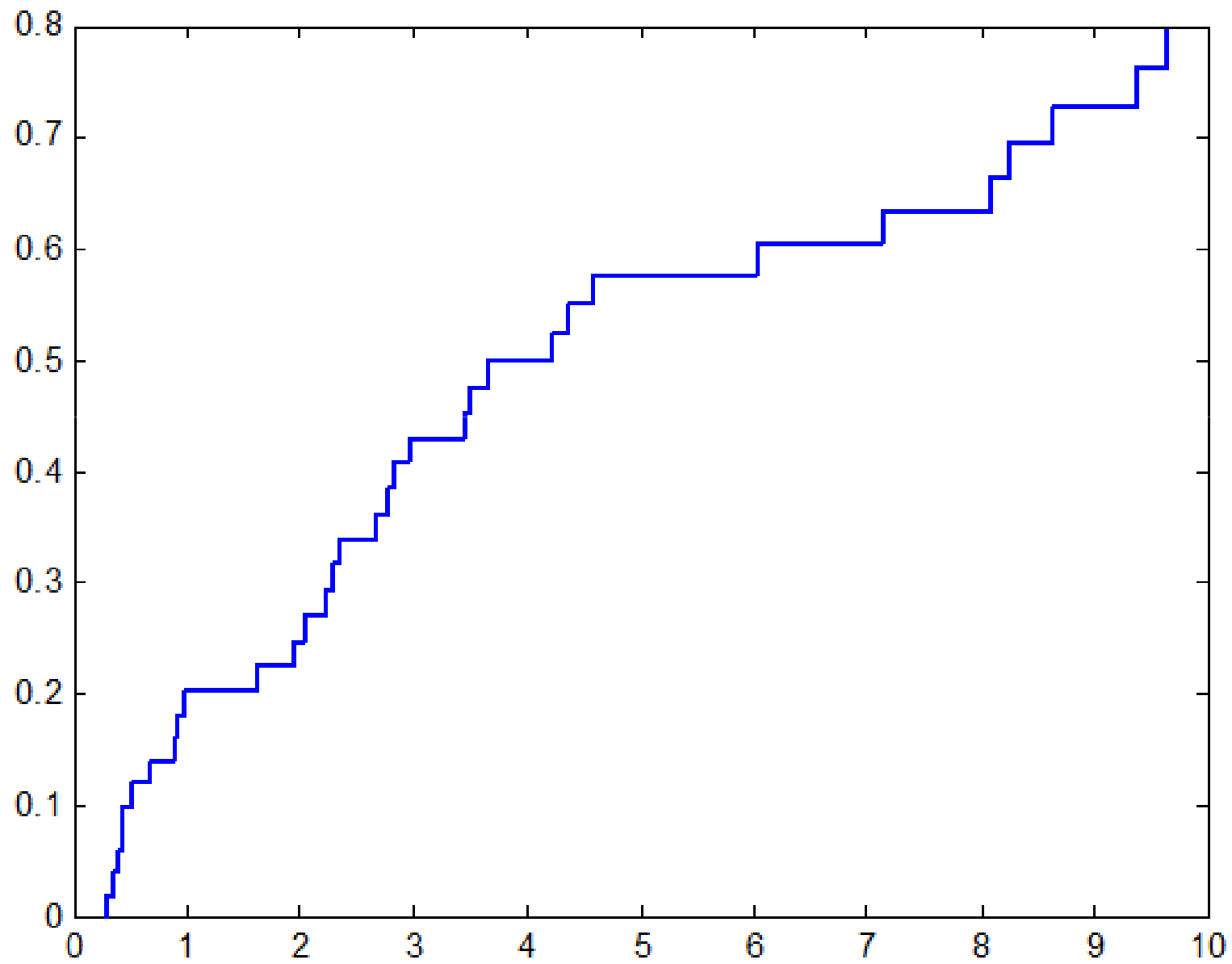
Vom genera aleator timpi de supraviețuire și timpi cenzurați, date ce vor avea o repartție exponențială.

```
>> y = exprnd(10,50,1); % generează un vector coloană cu 50 de  
componente, ce au o repartție exponențială de medie  $\mu = 10$ , componente ce  
reprezintă timpii de supraviețuire  
>>d = exprnd(20,50,1); % timpii de renunțare (drop-out) care au o repartție  
exponențială de medie  $\mu = 20$ .  
>>t = min(y,d);  
>>censored = (y>d); % se obține un vector boolean de aceeași dimensiune cu  
y; avem 1 pentru observațiile cenzurate.
```

`% Vom calcula și desena funcția cumulativă de distribuție empirică;
vom folosi parametrul 'censoring', a cărui valoare este vectorul boolean
de aceeași dimensiune cu X.`

Pentru datele cenzurate se atribuie valoarea 1.

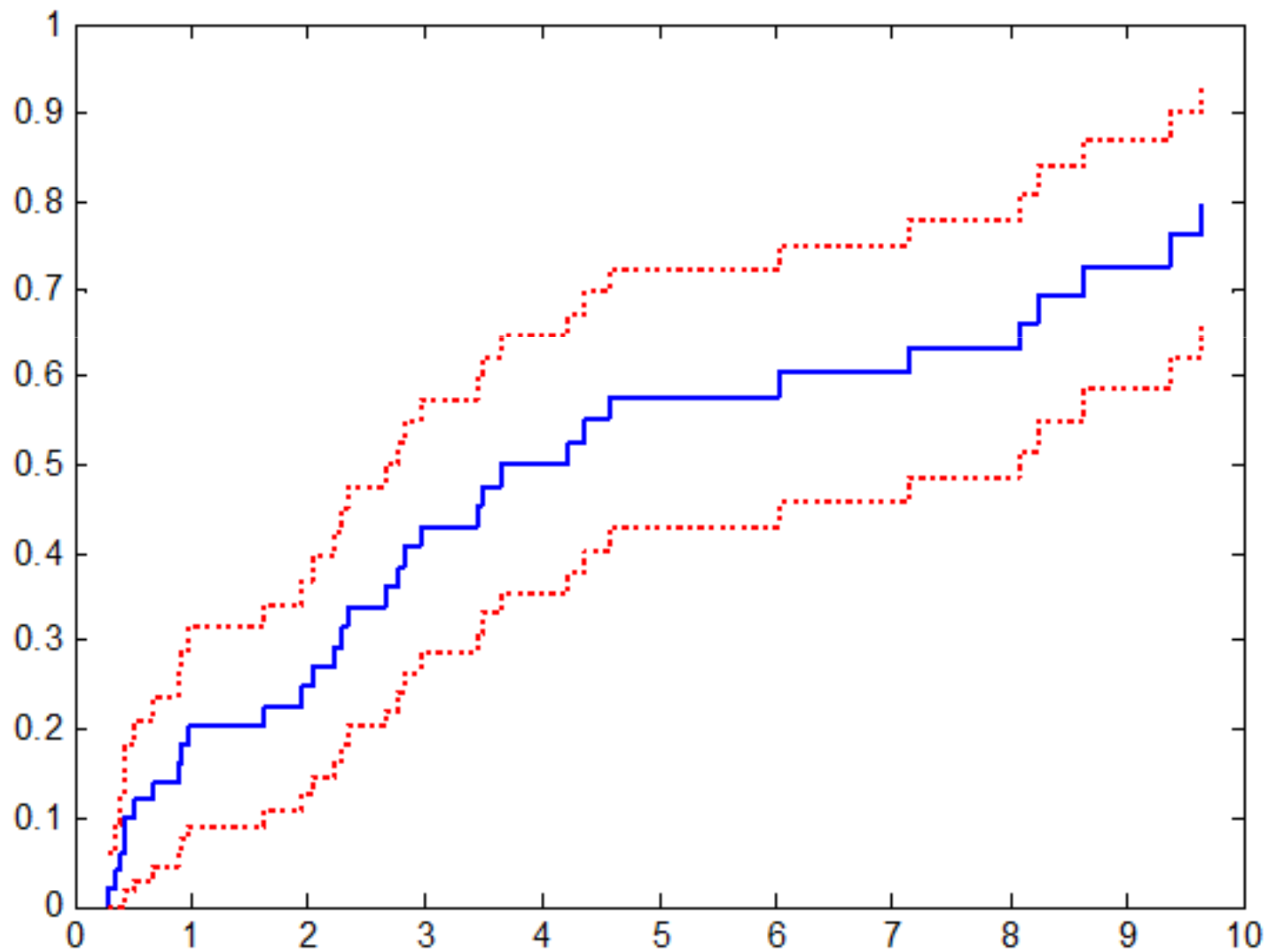
```
>> [f,x,flo,fup] = ecdf(t,'censoring',censored);  
>> stairs(x,f,'LineWidth',2)  
>> hold on
```



Vom calcula și desena limitele intervalului de încredere

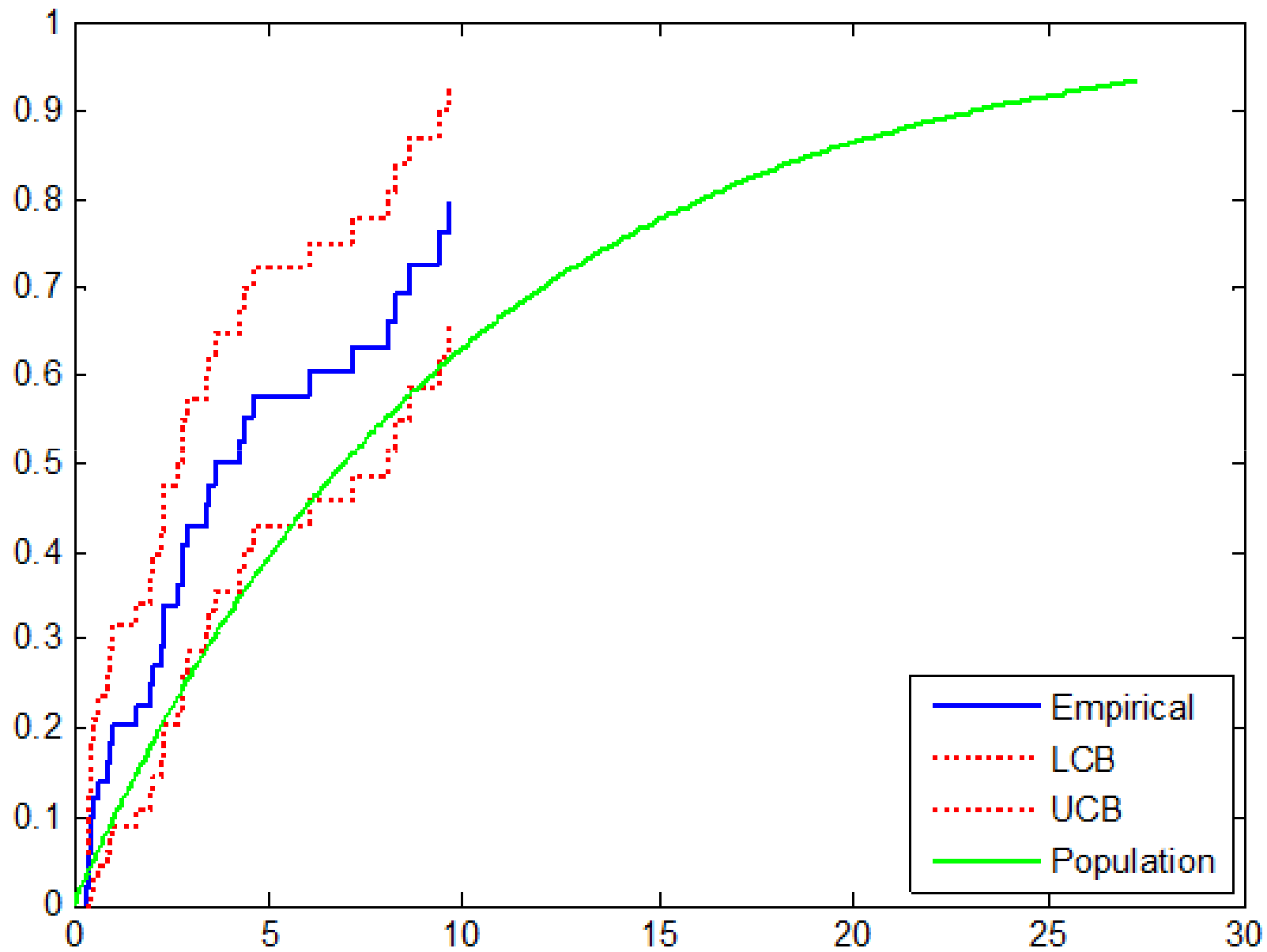
```
>> stairs(x, flo, 'r:', 'LineWidth', 2)
```

```
>> stairs(x, fup, 'r:', 'LineWidth', 2)
```



Vom desena în același cadran și funcția de distribuție cumulativă a populației:

```
>>xx = 0:1:max(t);  
>>yy = 1-exp(-xx/10);  
>>plot(xx,yy,'g-','LineWidth',2)  
>>legend('Empirical','LCB','UCB','Population',...  
         'Location','SE')  
>>hold off
```



Dacă suntem interesați în a compara experiența de supraviețuire la două grupuri de subiecți, putem calcula curba Kaplan-Meier pentru fiecare grup. O asemenea abordare dă o comparație în anumite valori de timp, arbitrare, iar dacă cele două curbe sunt semnificativ diferite, numai în anumite valori ale timpului.

10. Exemplu

Vom considera un al doilea grup de subiecți testați pentru răul de mare, în același simulator, doar că frecvența și accelerația mișcării s-a dublat față de primul experiment. Tabelul următor datele obținute:

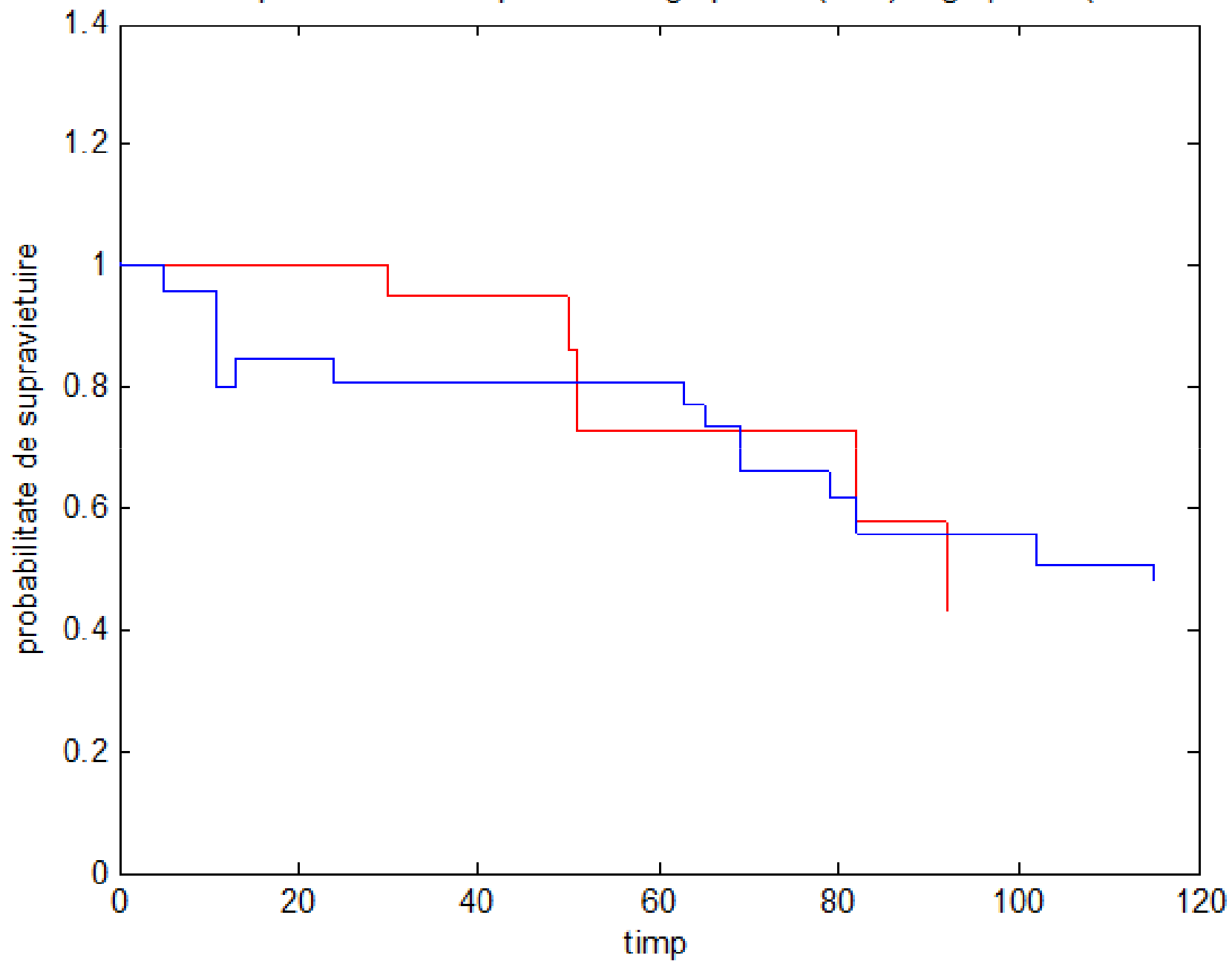
nr	1	2	3	4	5	6	7	8	9	10	11	12
T	5	6*	11	11	13	24	63	65	69	69	79	82
p	0.96			0.89	0.85	0.81	0.77	0.74		0,66	0,62	

nr	13	14	15	16	17	28
T	82	102	115	120*	120*	120*
p	0.56	0.51	0,48				

Vom desena curbele Kaplan-Meier în același sistem de axe:

```
>>x=[0 0.01 30 50 51 82 92];y=[ 1.01 1 0.95 0.86 0.73 0.58 0.43];  
>>x1=[ 92 120];y1=[0 0];  
>> [xb, yb]=stairs(x, y); [x1b y1b]=stairs(x1,y1);  
>>x2=[0 0.01 5 11 13 24 63 65 69 79 82 102 115];  
>> y2=[1.01 1 0.96 0.80 0.85 0.81 0.77 0.74 0.66 0.62 0.56 0.51 0.48];  
>> [x2b, y2b]=stairs(x2, y2); [x3b y3b]=stairs(x3,y3);  
>> plot(xb,yb,'r',x1b,y1b,'k', x2b,y2b,'b',x3b,y3b,'k')
```

curbele Kaplan Meier corespunzatoare grupului 1 (rosu) si grupului 2(albastru)



Testul log-rank

Metoda cea mai cunoscută de comparare a timpurilor de supraviețuire pentru două grupuri independente de subiecți este *testul log-rank*.

Testul log-rank este o metodă non-parametrică de testarea ipotezei nule. H_0 , în sensul că verificăm dacă nu există diferențe semnificative între curbele de supraviețuire.

Curbele Kaplan-Meier și testul logrank se folosesc când variabila predictor este categorială (de exemplu medicament vs placebo) sau când ia un număr mic de valori, care pot fi considerate categoriale.

Testul logrank se utilizează pentru a detecta diferențele între curbele de supraviețuire în condițiile în care rata moratalității (ieșirii din funcțiune) într-un grup este semnificativ mai mare decât rata corespunzătoare în al doilea grup și raportul acestor rate este constant în timp.

Acest test dă ponderi egale tuturor evenimentelor (decese, ieșiri din funcțiune)

$$T = \sum_{j=1}^n \sum_{i=1}^m \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

Principiul testului log-rank constă în divizarea scalei timpilor de supraviețuire, ignorând timpii cenzurați.

13. Exemplu

Considerăm două grupuri de 10 și respectiv 8 pacienți ce au suferit un transplant de ficat, cărora li se administrează medicație diferită. Prezentăm timpul de supraviețuire pentru cele două grupuri, în luni.

G1	1* 3 4* 5 5 6* 7 7 7* 8
G2	2 2 3* 4 6* 6* 7 10

Reamintim că * înseamnă dată cenzurată, în cazul nostru pacient retras în luna t_i din studiu.

În tabelele următoare vom folosi notațiile:

- t – timpul în care are loc evenimentul (în luni);
- n - numărul de subiecți care sunt încă sub observație în luna t
- n_1 - numărul de subiecți din grupul 1 care sunt încă sub observație în luna t
- n_2 - numărul de subiecți din grupul 2 care sunt încă sub observație în luna t
- r – numărul de decese raportat în luna t .
- c - numărul de valori cenzurate raportat în luna t .
- σ_1 – numărul de decese din grupul 1, raportat în luna t .
- σ_2 – numărul de decese din grupul 2, raportat în luna t .
- e_1 – numărul de decese din grupul 1, așteptat în luna t .
- e_2 – numărul de decese din grupul 2, așteptat în luna t .

Luna nr 1

t	n	n_1	n_2	r	c	σ_1	σ_2	e_1	e_2
1	18	10	8	0	1	0	0	0	0

Luna nr 2

t	n	n_1	n_2	r	c	σ_1	σ_2	e_1	e_2
1	18	10	8	0	1	0	0	0	0
2	17	9	8	2	0	0	2	1.06	0.94

Raportul de subiecți sub observație în grupul 1 este $\frac{9}{17}$ iar în grupul 2, $\frac{8}{17}$.

Astfel evenimentele așteptate sunt $e_1 = 2 \cdot \frac{9}{17} = 1.06$ și $e_2 = 2 \cdot \frac{8}{17} = 0.94$.

Luna 3

t	n	n_1	n_2	r	c	σ_1	σ_2	e_1	e_2
1	18	10	8	0	1	0	0	0	0
2	17	9	8	0	0	0	2	1.06	0.94
3	15	9	6	1	1	1	0	0.6	0.4

Raportul de subiecți sub observație în grupul 1 este $\frac{9}{15}$ iar în grupul 2, $\frac{6}{15}$,
evenimentele așteptate sunt $e_1 = \frac{9}{15} = 0.6$ și $e_2 = \frac{6}{15} = 0.4$.

Tabelul următor prezintă situația cu toți timpii de supraviețuire:

t	n	n_1	n_2	r	c	σ_1	σ_2	e_1	e_2
1	18	10	8	0	1	0	0	0	0
2	17	9	8	0	0	0	2	1.06	0.94
3	15	9	6	1	1	1	0	0.6	0.4
4	13	8	5	1	1	0	1	0.62	0.38
5	11	7	4	2	0	2	0	1.27	0.73
6	9	5	4	0	3	0	0	0	0
7	6	4	2	3	1	2	1	2	1
8	2	1	1	1	0	1	0	0.5	0.5
Total						6	4	6.05	3.95

Aplicăm testul χ^2 cu un grad de libertate (având două grupuri)

$$T = 1.06 + \frac{(1.06)^2}{0.94} + \frac{(0.4)^2}{0.6} + 0.4 + 0.62 + \frac{(0.62)^2}{0.38} + \frac{(0.73)^2}{1.27} + 0.73 + 0.5 + 0.5 = 6.7032$$

	p value											
df	0.25	0.20	0.15	0.10	0.05	0.025	0.02	0.01	0.005	0.0025	0.001	0.0005
1	1.32	1.64	2.07	2.71	3.84	5.02	5.41	6.63	7.88	9.14	10.83	12.12
2	2.77	3.22	3.79	4.61	5.99	7.38	7.82	9.21	10.60	11.98	13.82	15.20
3	4.11	4.64	5.32	6.25	7.81	9.35	9.84	11.34	12.84	14.32	16.27	17.73
4	5.39	5.59	6.74	7.78	9.49	11.14	11.67	13.23	14.86	16.42	18.47	20.00
5	6.63	7.29	8.12	9.24	11.07	12.83	13.33	15.09	16.75	18.39	20.51	22.11
6	7.84	8.56	9.45	10.64	12.53	14.45	15.03	16.81	18.55	20.25	22.46	24.10
7	9.04	9.80	10.75	12.02	14.07	16.01	16.62	18.48	20.28	22.04	24.32	26.02
8	10.22	11.03	12.03	13.36	15.51	17.53	18.17	20.09	21.95	23.77	26.12	27.87

Din tabelul testului χ^2 obținem că $P \approx 0.01 < 0.05$ și astfel ipoteza nulă nu este valabilă, ceea ce se traduce că există deosebiri între timpii de supraviețuire ale grupurilor 1 și 2.

10. Exemplu

Reluăm experimentele cu răul de mare, pentru a vedea dacă există diferențe semnificative între timpii de supraviețuire a celor două grupuri:

În primul experiment avem 5 evenimente(rău de mare) în minutele: 30,50,51,82,92.

În al doilea experiment avem 14 evenimente(rău de mare): câte unul în minutele 5,13,24,63,65,78,102,115 și câte două în minutele 11,69,82.

Pentru cele două evenimente combinate avem 16 valori distincte de timp, și vom obține 16 intervale de timp:

$(0,5']$, $(5',11']$, $(11',13']$, $(13',24']$, $(24',30']$, $(30',50']$, $(50',51']$, $(51',63']$,
 $(63',65']$, $(65',78']$, $(78',82']$, $(82',92']$, $(92',102']$, $(102',115']$, $(115',120']$

Aplicăm algoritmul de calcul prezentat în exemplul nr 12, utilizând aceleași notații:

t	n	n_1	n_2	r	c	σ_1	σ_2	e_1	e_2
t_1	42	21	21	1	0	0	1	0.5	0.5
t_2	41	21	20	2	1	0	2	1.03	0.97
t_3	38	21	17	1	0	0	1	0.55	0.45
t_4	37	21	16	1	0	0	1	0.55	0.45
t_5	36	21	16	1	0	0	1	0.58	0.42
t_6	35	21	15	1	1	1	0	0.57	0.43
t_7	33	18	15	1	0	1	0	0.56	0.44
t_8	32	17	15	1	0	0	1	0.53	0.47
t_9	31	17	14	1	0	0	1	0.55	0.45
t_{10}	30	17	13	2	1	0	2	1.86	1.14
t_{11}	27	16	11	1	0	0	1	0.59	0.41

t_{11}	27	16	11	1	0	0	1	0.59	0.41
t_{12}	26	16	10	3	0	1	2	1.86	1.14
t_{13}	23	15	8	1	0	1	0	0.65	0.35
t_{14}	22	14	8	1	0	0	1	0.64	0.36
t_{15}	21	14	7	1	0	0	1	0.67	0.33
t_{16}	20	14	6	0	20	0	0	0	0

□

$$T = \sum_{j=1}^{16} \sum_{i=1}^2 \frac{(O_{ij} - E_{ij})^2}{E_{ij}} = 20.3261$$

Din tabelul testului χ^2 obținem $P \approx 0$ și astfel respingem ipoteza nulă.

Modelul Cox

Metoda clasică de analiza supraviețuirii nu poate fi utilizată pentru a explora efectul simultan asupra supraviețuirii, a mai multor variabile.

În acest caz, nu se poate utiliza direct metoda regresiei multiple din cel puțin două motive:

- (a) variabila ce descrie timpul de supraviețuire nu este de cele mai multe ori normal repartizată (de obicei este exponențială sau Weibull) și
- (b) analiza supraviețuirii utilizează așa-numitele *date cenzurate*, adică unele observații sunt incomplete.

Analiza supraviețuirii utilizează mai multe metode regresive specializate, din care amintim:

- Analiza regresivă a hazardurilor proporționale –modelul Cox (*Cox proportional hazards regression analysis*), pe scurt -modelul Cox al hazardului proporțional (*Cox proportional hazard model*);
- Modelul Cox al hazardului proporțional cu covarianțe dependente de timp (*Cox proportional hazard model with time-dependent covariates*);
- Modelul regresiv exponențial (*Exponential regression model*);
- Modelul liniar regresiv log-normal (*Log-normal linear regression model*);
- Modelul liniar regresiv normal (*Normal linear regression model*).

Modelul Cox al hazardului proporțional

Funcția de supraviețuire (survival function) este definită prin probabilitatea:

$$S(t) = P\{T > t\},$$

unde t reprezintă timpul în general, iar T este timpul până la deces.

Repartiția duratei de viață (*lifetime*) este dată de

$$F(t) = 1 - S(t),$$

unde $f(t) = \frac{d}{dt} F(t)$ reprezintă rata de deces pe unitatea de timp.

Funcția hazard (*hazard function*) este definită de formula:

$$\lambda(t) = P(t < T < t + dt) = \frac{f(t)dt}{S(t)} = -\frac{S'(t)dt}{S(t)}$$

Funcția hazard reprezintă deci riscul de a deceda într-un interval foarte scurt de timp dt după un timp dat T , presupunând evident supraviețuirea până la acel moment.

Curba de supraviețuire reprezintă grafic supraviețuirea cumulată în funcție de timp, în timp ce derivata curbei de supraviețuire este rata survenirii decesului într-un interval scurt de timp, iar acesta este hazardul.

De exemplu, dacă ne așteptăm ca 20% dintre pacienții cu un anumit tip de cancer să moară în acest an, atunci hazardul este 20% pe an.

Așadar, hazardul este un risc, iar raportul hazardelor (*hazard ratio*) este un risc relativ.

Dacă raportul hazardelor este 0,5, atunci riscul de a muri într-un grup (cel tratat) este jumătate din riscul de a muri în celălalt grup (placebo).

Regresia hazardului proporțional (modelul lui Cox) utilizează metodele de regresie pentru a prezice riscul relativ pe baza uneia sau mai multor variabile X.

Hazardul (riscul) se poate modifica în timp:

- el poate crește în timp; *de exemplu* la pacienții cu cancer, cu cât trece mai mult timp de la momentul diagnosticului, cu atât riscul de deces ca urmare a cancerului este mai mare.
- hazardul poate scădea în timp: tot la pacienții cu cancer, la care s-a obținut o remisie prin tratament, cu cât trece mai mult timp fără o recădere, cu atât riscul de recădere este mai mic.

Se numește "hazard proporțional" deoarece se bazează pe presupunerea că, dacă hazardul poate să varieze în timp, raportul hazardelor rămâne constant (sau altfel spus, cele două funcții ale hazardului sunt proporționale una cu cealaltă).

Presupunerea "hazardului proporțional" este rezonabilă în multe situații clinice dar nu întotdeauna; *de exemplu* atunci când comparăm un tratament medical cu unul chirurgical, al doilea are de obicei un hazard mai mare perioperator, pe când la cel medical hazardul crește după un timp mai lung, motiv pentru care regresia hazardului proporțional nu este indicată.

Modelul lui Cox ne furnizează o estimare a efectului tratamentului asupra supraviețuirii după ajustarea pentru celelalte variabile independente și ne permite să estimăm hazardul (sau riscul) de deces sau alt eveniment de interes la indivizi, date fiind variabilele lor prognostice.

Chiar dacă grupurile (tratament și martor) sunt similare în privința variabilelor cunoscute ca afectând supraviețuirea, folosirea unui model al lui Cox pentru aceste variabile poate produce o estimare mai precisă a efectului tratamentului.

De exemplu, într-o analiză multivariată trebuie să includem atât de mulți pacienți, încât să avem cel puțin câte 10 efecte (în acest caz decese) pentru fiecare variabilă independentă.

Modelul Cox al hazardului proporțional este foarte general între modelele regresive deoarece nu se bazează pe nicio ipoteză prealabilă privind repartiția supraviețuirii. El se bazează doar pe presupunerea că hazardul este o funcție doar de variabilele independente (predictive, covarianțe) Z_1, Z_2, \dots, Z_k , adică:

$$h(t; Z_1, Z_2, \dots, Z_k) = h_0(t) \cdot \exp(b_1 \cdot Z_1 + b_2 \cdot Z_2 + \dots + b_k \cdot Z_k),$$

sau, logaritmând,

$$\ln\left(\frac{h(t, Z_1, \dots, Z_k)}{h_0(t)}\right) = b_1 \cdot Z_1 + \dots + b_k \cdot Z_k$$

fiind deci un model semi-parametric.

$h(t, Z)$ poate fi considerată o densitatea de repartiție Weibull deoarece:

$$h(t, Z) = \lambda \cdot p \cdot t^{p-1},$$

unde

$$\lambda = \exp(b_1 \cdot Z_1 + \dots + b_k \cdot Z_k)$$

$$h_0(t) = p \cdot t^{p-1}$$

Termenul $h_0(t)$ se numește hazardul de bază (*baseline hazard, underlying hazard function*) reprezentând hazardul pentru un anumit individ atunci când toate variabilele independente sunt egale cu zero. $h_0(t)$ depinde doar de timp $-t$.

Exponențiala se referă doar la $Z = (Z_1, \dots, Z_k)$ și variabilele Z_i sunt independente de timp.

Exemple de variabile independente de timp (în domeniul medical):

- variabile care nu se schimbă în timp, de exemplu: sexul. Valorile sunt setate la $t = 0$;
- variabile ce nu par a se schimba în timp, *de exemplu* statutul de fumător;
- variabile ce sunt considerate că nu depind de timp: vârsta, greutatea.

Modelul Cox este foarte utilizat, fiind o alegere sigură în multe situații.

Hazardurile estimate sunt nenegative. în plus chiar dacă $h_0(t)$ nu este specificat, putem estima coeficienții β_i și apoi calcula rata hazardului (rata de risc - *hazard ratio*), notată HR .

În analiza supraviețuirii este preferat modelul Cox celui logistic, deoarece modelul Cox lucrează și cu date cenzurate.

hazard ratio- HR este definită de formula:

$$HR = \frac{h(t, Z^*)}{h(t, Z)},$$

unde $Z = (Z_1^*, \dots, Z_k^*)$ și $Z = (Z_1, \dots, Z_k)$

Pentru simplificarea interpretării luăm $HR \geq 1$, adică $h(t, Z^*) \geq h(t, Z)$ considerând Z^* un grup cu risc (hazard) mai mare, de exemplu grupul placebo și Z un grup cu risc mai mic, cum ar fi grupul cărui i se aplică un anumit tratament.

Astfel :

$$HR = \frac{h_0(t) \cdot \exp\left(\sum_{i=1}^k \beta_i \cdot Z_i^*\right)}{h_0(t) \cdot \exp\left(\sum_{i=1}^k \beta_i \cdot Z_i\right)} = \exp\left(\sum_{i=1}^k \beta_i \cdot (Z_i^* - Z_i)\right)$$

Să menționăm că, totuși, trebuie luate în considerație două condiții:

- (a) trebuie să existe o relație multiplicativă între $h_0(t)$ și funcția log-liniară a covarianțelor – *ipoteza proporționalității*, prin prisma hazardului
- (b) trebuie să existe o relație log-liniară între hazard și variabilele independente.

Remarcă

1. Selecția variabilelor explicative se face după aceeași regulă ca și în cazul regresiei liniare multiple.
2. Semnul coeficienților de regresie b_i trebuie interpretat după cum urmează. Un semn pozitiv indică un hazard ridicat deci, în consecință, un prognostic negativ pentru individul cu o valoare ridicată a acelei variabile. În schimb, un semn negativ pentru un anumit coeficient indică un hazard scăzut relativ la acea variabilă.
3. La fel ca și în cazul regresiei liniare multiple/logistice și în acest caz putem utiliza valoarea predictivă a modelului, folosind indicele prognostic definit de $IP = b_1 \cdot Z_1 + b_2 \cdot Z_2 + \dots + b_k \cdot Z_k$. Se poate calcula astfel funcția de supraviețuire $S(t) = \exp[-H_0(t)]^{\exp(IP)}$, unde $H_0(t)$, numit și *hazardul de bază cumulat*, este o funcție scară în variabila timp.

1. În cazul unei singure variabile $Z \in \mathbf{R}$, cu parametrul $\beta \in \mathbf{R}$, rata hazardului este:

$$h(t, Z) = h_0(t) \exp(\beta Z).$$

Rata hazardului pentru doi subiecți corespunzând lui Z_1 și Z_2 este:

$$\exp(\beta(Z_1 - Z_2)).$$

Exemplu

Vârsta dependenților de droguri.

Intr-un program de dezintoxicare ce a cuprins 600 subiecți, se notează cu a_i vârsta subiectului i la intrarea în program.

Rata hazardului a fost calculată ca fiind:

$$h(t, a) = h_0(t) \cdot \exp(-0.013a).$$

Se estimează că fiecare an din vârsta unui dependent multiplică riscul de a relua drogurile cu $e^{-0.013} = 0.99$. Nivelul de semnificație fiind $p = 0.07$ nu avem motiv să respingem ipoteza că vârsta nu are efect asupra reluării consumului de droguri.

2. În cazul a două variabile $(Z_1, Z_2) \in \mathbf{R}^2$, cu parametrii $(\beta_1, \beta_2) \in \mathbf{R}^2$, rata hazardului este:

$$h(t, Z_1, Z_2) = h_0(t) \exp(\beta_1 Z_1 + \beta_2 Z_2)$$

Dacă există o interacțiune între Z_1 și Z_2 , parametrii sunt $(\beta_1, \beta_2, \beta_{12}) \in \mathbf{R}^3$ și rata hazardului este:

$$h(t, Z_1, Z_2) = h_0(t) \exp(\beta_1 Z_1 + \beta_2 Z_2 + \beta_{12} Z_1 \cdot Z_2)$$

Exemplu

Numărul de tratamente de dezintoxicare (droguri) și numărul de tratamente antidepresive.

Se notează b_i numărul de tratamente de dezintoxicare ale subiectului i . b_i ia valori între 0 și 40.

Scorul de depresie Beck atribuit subiectului i se notează cu e_i . Acesta ia valori între 0 și 54 și anume: între 0 și 9, depresie inexistentă, între 10 și 18 depresie ușor moderată, între 19 și 29 depresie sever moderată și între 30 și 54 depresie severă.

Rata hazardului este:

$$h(t, b, e) = h_0(t) \cdot \exp(0.030b + 0.010e)$$

Astfel raportul hazardului pentru un dependent tipic cu depresie severă - 40 și un dependent cu depresie ușoară - 15 este $e^{25 \cdot 0.01} = 1.27$

Nivelul de semnificație p fiind mai mic de 0.05 indică un efect slab, în sensul că dependenti de droguri sever depresivi au mai multe șanse de recidivă.

Raportul hazardului pentru un dependent ce a urmat 20 tratamente de dezintoxicare în raport cu un dependent ce nu a urmat niciun tratament este $e^{20 \cdot 0.03} = 1.8$, ceea ce înseamnă că de 1.8 ori este mai probabil ca primul subiect să se drogheze din nou față de primul. Acest rezultat este foarte semnificativ, deoarece $p = 0.01$

Rata hazardului în cazul interacțiunii dintre numărul de tratamente de dezintoxicare și depresie este:

$$h(t, b, e) = h_0(t) \cdot \exp(0.079b + 0.043e + 0.0012b \cdot e)$$

S-a constatat că este totuși preferabil modelul care nu ia în seamă interacțiunea.

3. Dacă variabila $Z \in \{0,1\}$, (binară) și parametrul $\beta \in \mathbf{R}$, rata hazardului este

$$h(t, Z) = h_0(t) \exp(\beta Z).$$

În particular:

$$h(t, 0) = h_0(t) \text{ și } h(t, 1) = h_0(t) \cdot e^\beta$$

Rata hazardului pentru grupul 1 este de e^β ori mai mare decât rata hazardului la grupul 0.

Exemplu

Efectul rasei asupra tratamentului de dezintoxicare.

Subiecții au fost clasificați în albi și alții. S-au codificat cu albi cu 0, alții cu 1 și s-a notat cu f_i rasa subiectului i .

Rata hazardului este

$$h(t, f) = h_0(t) \exp(-0.29 f)$$

și rata hazardului pentru alții este 0.75 din rata hazardului pentru albi.

Acest rezultat este foarte semnificativ, deoarece $p < 0.01$, ceea ce înseamnă că tratamentul are mai mult succes la alții decât la albi.

4. Dacă variabilele sunt $Z_1 \in \{0,1\}$ și $Z_2 \in \mathbf{R}$, cu parametrii $(\beta_1, \beta_2) \in \mathbf{R}^2$, și respectiv $(\beta_1, \beta_2, \beta_{12}) \in \mathbf{R}^3$ în cazul în care există o interacțiune între Z_1 și Z_2 .

Pentru modelul fără interacțiuni, rata hazardului este:

$$h(t, Z_1, Z_2) = h_0(t) \exp(\beta_1 Z_1 + \beta_2 Z_2)$$

În particular

$$h(t, Z_1 = 0, Z_2) = h_0(t) \exp(\beta_2 Z_2)$$

$$h(t, Z_1 = 1, Z_2) = h_0(t) \exp(\beta_1 + \beta_2 Z_2)$$

Pentru modelul cu o interacțiune, rata hazardului este:

$$h(t, Z_1, Z_2) = h_0(t) \exp(\beta_1 Z_1 + \beta_2 Z_2 + \beta_{12} Z_1 \cdot Z_2)$$

În particular

$$h(t, Z_1 = 0, Z_2) = h_0(t) \exp(\beta_2 Z_2 + \beta_{12} Z_1 \cdot Z_2)$$

$$h(t, Z_1 = 1, Z_2) = h_0(t) \exp(\beta_1 + \beta_2 Z_2 + \beta_{12} Z_1 \cdot Z_2)$$

Exemplu

rasa și numărul de tratamente de dezintoxicare.

Fiecare individ este codificat $Z_1 = 0$, alb și $Z_1 = 1$, pentru non-alb.

Pentru modelul fără interacțiune, funcția hazard este:

$$h(t, f, Z_2) = h_0(t) \exp(-0.26Z_1 + 0.027Z_2)$$

în timp ce pentru modelul cu interacțiune, funcția hazard este:

$$h(t, f, Z_2) = h_0(t) \exp(-0.25Z_1 + 0.027Z_2 - 0.001Z_1 \cdot Z_2)$$

11. Exemplu

Datele următoare reprezintă supraviețuirea, în zile, a unor pacienți cu limfom histiocitar difuz, pacienți ce fac parte din două grupuri distincte. Vor fi comparate două grupuri distincte: pacienți aflați în stadiul al 3-lea, respectiv în stadiul al 4-lea.

Stadiul 3: 6, 19, 32, 42, 42, 43*, 94, 126*, 169*, 207, 211*, 227*, 253, 255*, 270*, 310*, 316*, 335*, 346*

Stadiul 4: 4, 6, 10, 11, 11, 11, 13, 17, 20, 20, 21, 22, 24, 24, 29, 30, 30, 31, 33, 34, 35, 39, 40, 41*, 43*, 45, 46, 50, 56, 61*, 61*, 63, 68, 82, 85, 88, 89, 90, 93, 104, 110, 134, 137, 160*, 169, 171, 173, 175, 184, 201, 222, 235*, 247*, 260*, 284*, 290*, 291*, 302*, 304*, 341*, 345*

$\gg b = \text{coxphfit}(X, y)$ returnează vectorul coloană de dimensiune p ale cărui componente sunt coeficienții regresiei lui Cox a hazardelor (riscurilor) prognozate pentru răspunsurile y , pentru cei p predictorii din matricea X . X este o matrice cu p coloane și n linii (n fiind numărul observațiilor).

```

>> x1=[6, 19, 32, 42, 42, 43, 94, 126, 169, 207, 211, 227, 253, 255, 270,
310, 316, 335, 346];
>> x1cens=[1 1 1 1 1 0 1 0 0 1 0 0 1 0 0 0 0 0 0];
>> x2=[4, 6, 10, 11, 11, 11, 13, 17, 20, 20, 21, 22, 24, 24, 29, 30, 30, 31,
33, 34, 35, 39, 40, 41, 43, 45, 46, 50, 56, 61, 61, 63, 68, 82, 85, 88, 89,
90, 93, 104, 110, 134, 137, 160, 169, 171, 173, 175, 184, 201, 222, 235,
247, 260, 284, 290, 291, 302, 304, 341, 345];
>> x2cens=[1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 0 1 1 1 1 0 0 1
1 1 1 1 1 1 1 1 1 1 1 0 1 1 1 1 1 1 1 0 0 0 0 0 0 0 0 0 0];
>> b1=[1]; for i=2:19 b1=[1 b1];end
>> b2=[2]; for i=2:61 b2=[1 b2];end
>> b=[b1 b2]';x=[x1 x2]';
>> xcens=[x1cens x2cens]';
>> A=[b x xcens];
>> b = coxphfit(A,x)

```

```

b =
    3.8410
   -3.2932
    0.0994

```

Aşadar rata hazardului este:

$$h(t, b, x, xcens) = h_0(t) \exp(3.8410b - 3.2932x + 0.0994xcens)$$

12. Exemplu

Vom genera aleator date cu repartiție Weibull:

```
>> x = 4*rand(100,1);  
>> A = 50*exp(-0.5*x); B = 2;  
>> y = wblrnd(A,B);
```

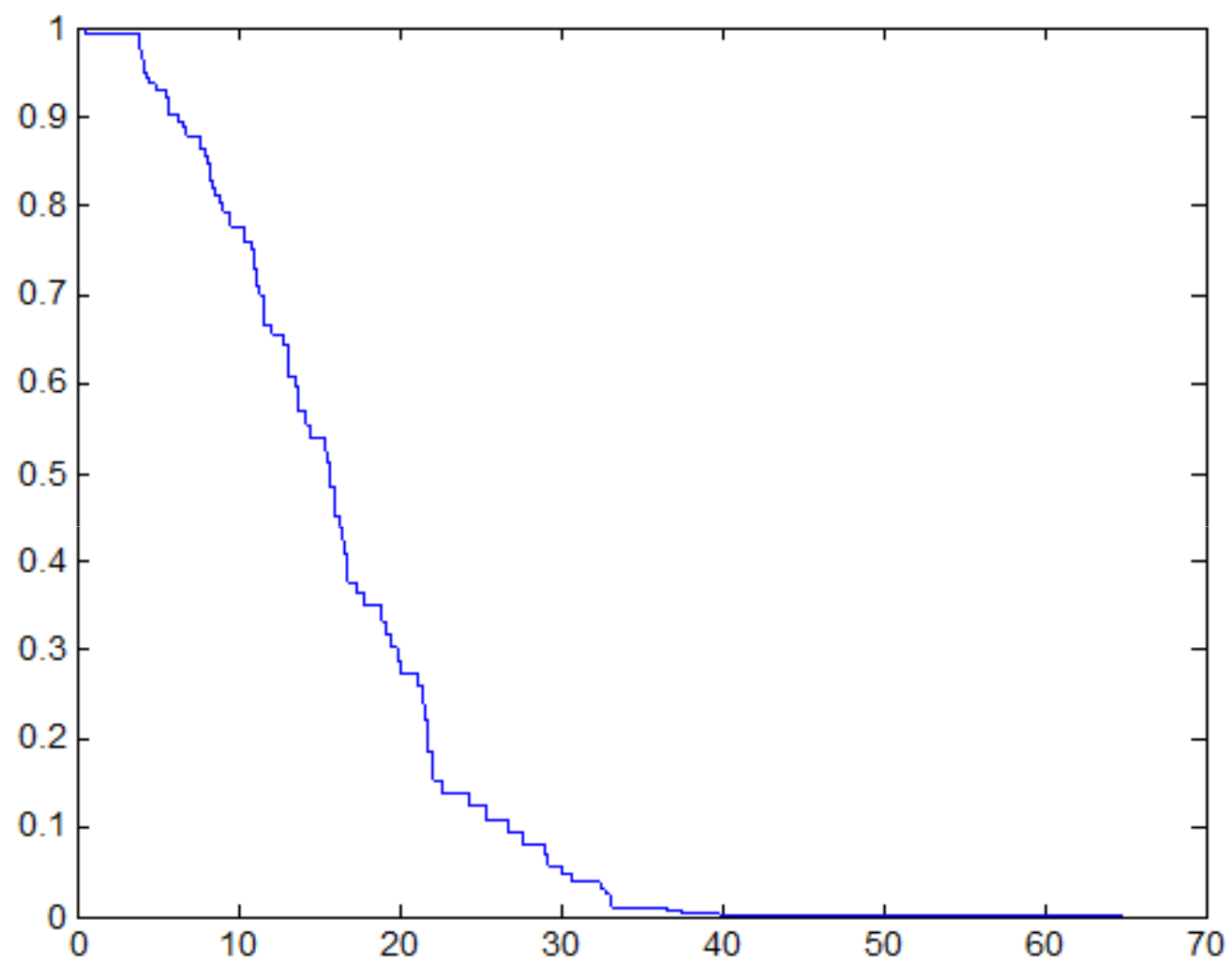
Vom implementa modelul Cox:

```
>> [b,logL,H,stats] = coxphfit(x,y);
```

Vom obține $b = 0.9409$ și $p = 6.9462e-014$

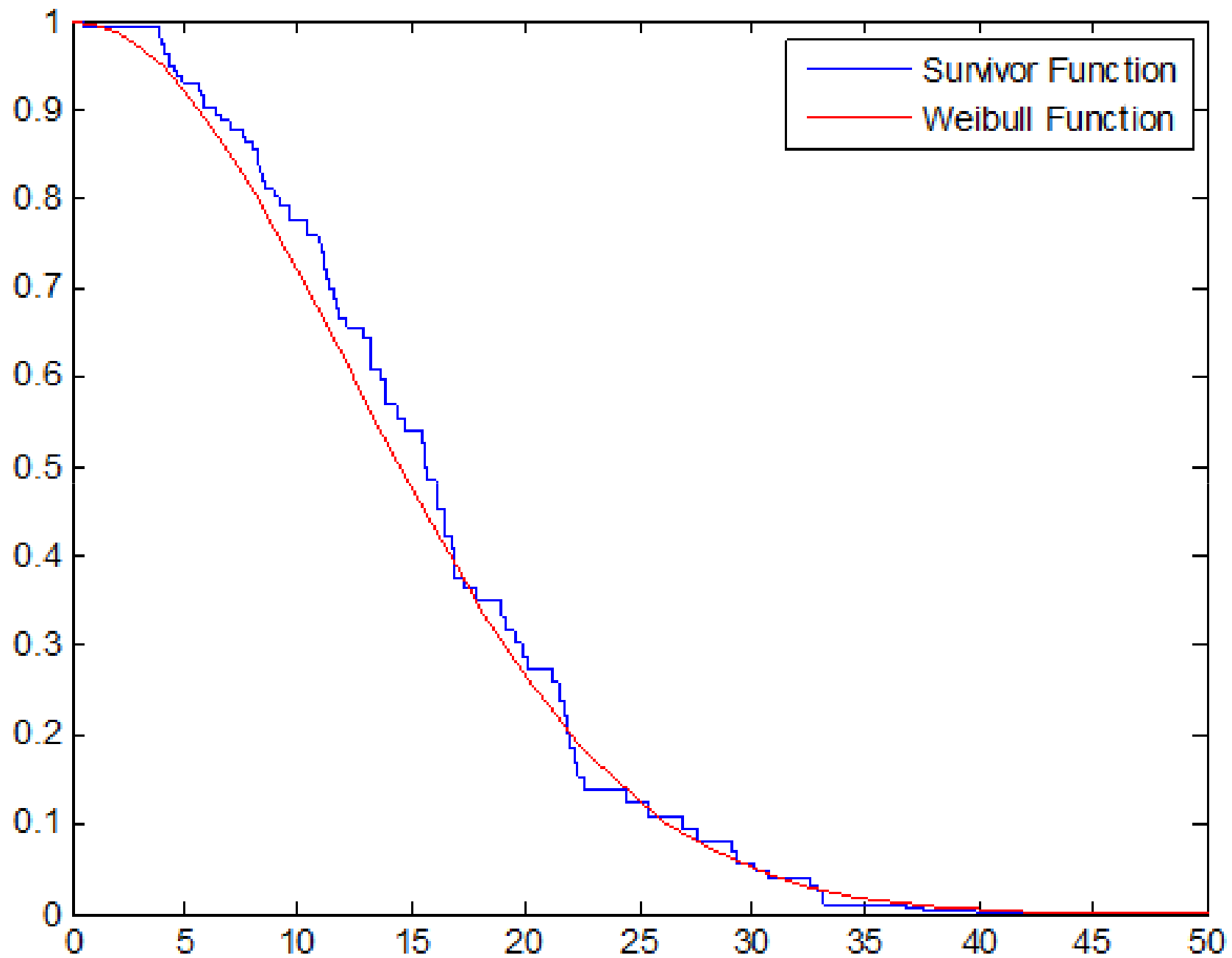
Vom desena funcția Cox estimată:

```
>> stairs(H(:,1),exp(-H(:,2)))
```



Vom desena funcția Cox estimată în același sistem de axe cu repartiția Weibull:

```
>>xx = linspace(0,100);  
>>line(xx,1-wblcdf(xx,50*exp(-0.5*mean(x)),B),'color','r')  
>>xlim([0,50])  
>>legend('Survivor Function','Weibull Function')
```

Modele aditive

Presupunem, la fel ca și în cazul regresiei multiple, că avem o variabilă dependentă Y și k variabile explicative (predictoare) X_1, X_2, \dots, X_k .

Spre deosebire de cazul modelelor liniare, cazul modelului aditiv implică posibila legătură între variabila dependentă și predictorii sub forma:

$$Y = f_1(X_1) + f_2(X_2) + \dots + f_k(X_k) + \varepsilon$$

unde $f_j, j = 1, 2, \dots, k$ sunt, în general, funcții de clasă C^∞ , în unele cazuri și funcții de clasă C^1 , iar ε este o variabilă aleatoare repartizată normal standard $N(0, 1)$.

Este ușor de observat că un model aditiv reprezintă generalizarea modelului regresiei liniare multiple (pentru $\varepsilon = 0$).

Cu alte cuvinte, în loc de un singur coeficient *per* variabilă explicativă, la modelele aditive găsim o funcție nespecificată *per* fiecare predictor, care va trebui să fie estimată în vederea prognozei optime a valorilor variabilei dependente.

Ipoteza aditivității $\sum f_i(X_i)$ este o restricție a cazului general al unui model predictiv de tipul $Y = f(X_1, X_2, \dots, X_n)$.

Funcțiile parametru ale modelului aditiv sunt estimate până la o constantă aditivă.

Modele aditive generalizate

Un *model liniar generalizat* este reprezentat prin ecuația:

$$Y = g(b_0 + b_1 \cdot X_1 + b_2 \cdot X_2 + \dots + b_k \cdot X_k),$$

unde g este o funcție de clasă C^∞ , nedeterminată.

Dacă notăm formal g^{-1} inversa funcției g , funcție numită *funcție legătură* (*link function*), atunci putem scrie ecuația de mai sus sub forma ușor modificată:

$$g^{-1}(E[Y]) = b_0 + b_1 \cdot X_1 + b_2 \cdot X_2 + \dots + b_k \cdot X_k,$$

unde $E[Y]$ reprezintă media variabilei dependente Y .

Combinând un model aditiv cu un model liniar generalizat, vom obține ecuația modelului, cunoscut sub denumirea de *model aditiv generalizat*, sub forma:

$$g^{-1}(E[Y]) = f_1(X_1) + f_2(X_2) + \dots + f_k(X_k) .$$

Problema de bază la aceste modele este estimarea funcțiilor parametri f_i ale modelului.

Cea mai cunoscută metodă de evaluare a funcțiilor f_i este reprezentată de o interpolare pe bază de diagrame de dispersie (*scatterplot smoother*), utilizând funcții spline cubice.

De exemplu, în cazul unui model simplu cu doar două funcții f_1 și f_2 , având forma:

$$Y = f_1(X_1) + f_2(X_2) + \varepsilon ,$$

folosind aproximarea spline, se obține forma celor două funcții:

$$f_1(X) = \delta_1 + X \cdot \delta_2 + \sum_{j=1}^{q_1-2} R(X, X_j^*) \cdot \delta_{j+2} ,$$
$$f_2(X) = \gamma_1 + X \cdot \gamma_2 + \sum_{j=1}^{q_2-2} R(X, X_j^*) \cdot \gamma_{j+2} ,$$

unde δ_j , γ_j sunt parametrii necunoscuți ai funcțiilor f_1 și f_2 , q_1 , q_2 reprezintă numărul parametrilor necunoscuți, iar x_j^* reprezintă nodurile de interpolare pentru cele două funcții.