

Folosirea pachetului MathCad în analiza statistică

George Daniel Mateescu*
Corina Saman*
Mihai Buneci*

Rezumat

MathCad-ul ne oferă posibilitatea de a face o analiză a datelor și de a testa modele posibile într-un mod simplu și ușor de folosit dacă cunoaștem matematica care se afla în spatele modelului.

Putem verifica ipotezele modelului regresiei cum ar fi dispersia constantă a erorilor, testarea ipotezelor statistice față de coeficienți, corelația erorilor, normalitatea erorilor ca variabile aleatoare.

Cuvinte cheie: regresie, ipoteze statistice

JEL: C22(Time-Series Models), C51(Model Construction and Estimation)

Pachetul **Mathcad** intră în categoria mediilor soft care asistă utilizatorul specialist într-un anumit domeniu. Permite realizarea de calcule matematice foarte complexe precum și activități conexe cum ar fi reprezentări grafice, realizarea unei documentații, precum și posibilitatea de export spre alte medii Windows.

Comenzile principale:

- **File:** operații uzuale privind fișierul ca entitate, deschidere, salvare, tiparire, transmitere prin fax sau poștă electronică
- **Edit:** comenzi de “dactilografie inteligentă”: ștergere, copiere, etc.
- **View:** se referă în special la subseturile de comenzi dedicate pentru diferite operații matematice și care sunt grupate în seturi (palette) ce pot fi activate prin *paleta matematică*
- **Insert:** se referă la obiecte matematice complexe cum sunt grafice de funcții, tabele (matrici) funcții (sistemul conține în formă predefinită cele mai cunoscute funcții matematice)
- **Format:** se referă atât la indicații de “înfățișare” a lucrării cât și la formate specifice calculelor matematice (cum este de exemplu numărul de zecimale care se afișează)
- **Math:** privește comanda de începere a calculelor dar și indicații referitoare la unitățile de măsură
- **Symbolics:** reprezintă un element de noutate în raport cu calculele numerice și conține elemente de inteligență artificială care permit evaluări simbolice (primitive sau derivare simbolică)

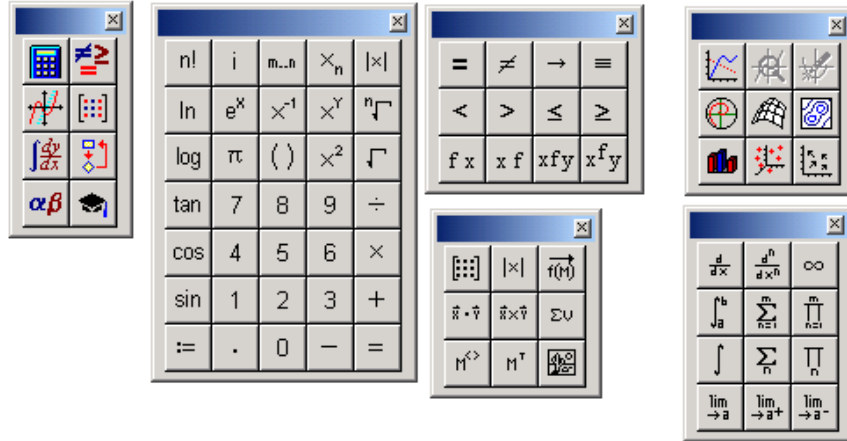
* Institutul de Prognoză Economică, Academia Română

- **Window:** permite gestiunea mai multor lucrări
- **Help:** mediu complex care conține indicații de utilizare

Detalierea comenzilor. Tot ceea ce ține de tratarea fișierului ca entitate, vizualizare, tipărite editare, ferestre, help, sunt foarte asemănătoare cu alte pachete din mediul Windows (Word, Excel, etc.)

Paleta matematică, conține cele mai semnificative comenzi pe următoarele grupe

- ⇒ aritmetică
- ⇒ logică
- ⇒ reprezentări grafice
- ⇒ matrici și vectori
- ⇒ calcul diferențial și integral
- ⇒ elemente de programare
- ⇒ alfabetul grec
- ⇒ evaluarea simbolică



Principiul de utilizare este al unei “foi de hârtie inteligente” în care utilizatorul scrie expresiile matematice în mod obișnuit iar “foaia” efectuează calculele și pune rezultatele exact în continuarea expresiei. Utilizatorul completează liniile și coloanele matricii a iar apoi scrie pur și simplu a⁻¹ iar “foaia de calcul” pune rezultatul. Se pot face calcule oricât de complicate așa cum se vede alăturat. Modul de introducere al expresiilor este “free hand”

$$a := \begin{bmatrix} 1 & 3 & 1 \\ 2 & 1 & -1 \\ 4 & 1 & 2 \end{bmatrix} \quad b := \begin{bmatrix} 4 & 2 & 1 \\ 1 & 2 & 2 \\ 8 & 5 & 0 \end{bmatrix}$$

$$|a| \cdot b - a^T \cdot b^{-1} = \begin{bmatrix} -93.158 & -47.421 & -22.368 \\ -24.316 & -45.842 & -45.737 \\ -186.526 & -114.737 & 1.105 \end{bmatrix}$$

Calcul simbolic

$$\frac{d}{dx} \frac{\sin(x)}{x^3 - \cos(1 + \sqrt{x})} \rightarrow \frac{\cos(x)}{(x^3 - \cos(1 + \sqrt{x}))} - \frac{\sin(x)}{(x^3 - \cos(1 + \sqrt{x}))^2} \cdot \left(3 \cdot x^2 + \frac{1}{2} \cdot \frac{\sin(1 + \sqrt{x})}{\sqrt{x}} \right)$$

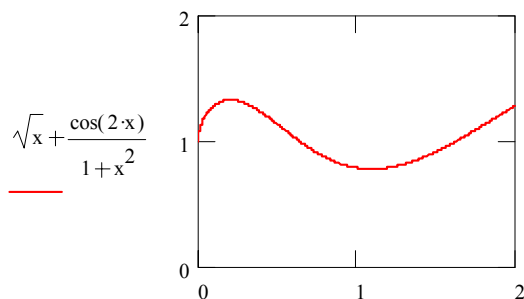
Reprezentări grafice

Se plasează:

pe orizontală, jos: variabila și intervalul (domeniul e definiție)

pe verticală, în stânga, legea de corespondență și codomeniul

Mediul Windows este atât de puternic integrat încât permite, de exemplu în Word, editarea direct în pagină prin lansarea implicită în background a pachetului Mathcad!



Calcul numerice

Integrala definită a unei funcții a cărei primitivă nu este exprimabilă prin funcții elementare. În cazul funcțiilor ale căror primitive se exprimă prin funcții elementare se poate apela și la calculul simbolic

$$\int_1^2 e^{-x^2} dx = 0.135$$

Sisteme de ecuații

Sistemele de ecuații liniare se rezolvă direct prin calcule cu matrici dar se pot rezolva și ecuații sau sisteme de ecuații neliniare. Este necesară introducerea unor valori inițiale pentru x și y (semnificația acestei introduceri poate fi explicată numai după ce vor fi însușite cunoștințe de analiză numerică). Este posibil ca, uneori, să nu fie obținută soluția sistemului dacă valorile de start nu satisfac anumite condiții.

În exemplul alăturat semnul = se obține prin selectare din tabelul de operatori de tip logic. Soluția se obține prin apelul funcției find care are exact acest rol.

$$\begin{aligned} x &:= 4 & y &:= 2 \\ \text{given} & & & \\ x^2 + 3y + x &= 5 \\ 2 + x \cdot y &= 3 \\ \text{find}(x, y) &= \begin{bmatrix} 1 \\ 1 \end{bmatrix} \end{aligned}$$

Analiză statistică

MathCad-ul ne oferă posibilitatea de a face o analiză a datelor și de a testa modele posibile într-un mod simplu și ușor de folosit dacă cunoaștem matematica care se afla în spatele modelului.

Estimarea, deducția și predicția sunt posibile pentru orice model simulat.

Cînd experimentăm modele posibile asociate fenomenelor economice trebuie luate decizii în funcție de rezultatele obținute în testarea ipotezelor statistice. În MathCad este suficient să formulăm ipoteza statistică, apoi măsurarea coincidenței dintre valorile observate și cele pe care ipoteza statistică le validează ca adevărate se face simplu calculînd statistica asociată ipotezei. Este necesar numai să traducem în formula potrivită ipoteza statistică. Avem un test statistic care este o funcție a datelor observate. Tinînd cont că ipoteza nulă determină distribuția de probabilitate a celui test statistic care apoi determină probabilitate este simplu de construit orice test statistic.

În cele ce urmează vom explora posibilitățile de analiză a validității unui model de regresie liniară pe care ni le oferă MathCad-ul. Si acestea sunt: estimatori după metoda celor mai mici pătrate, verificare ipotezelor modelului liniar, deducții asupra parametrilor funcției liniare, măsura potrivirii modelului pentru datele observate, analiza tabelii varianței pentru regresie.

Datele folosite în exemplele de mai jos sunt: producția industrială p_i (variație lunară,%), somajul s (rată lunară,%), și inflația inf (rată lunară,%), pentru perioada ianuarie 2000 - ianuarie 2002.

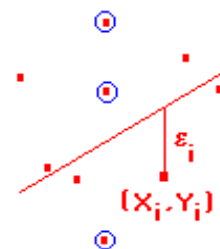
Condițiile impuse de Modelul regresiei Liniare

1. Variabila dependentă Y este o funcție liniară de X .
2. Variabila eroare ε este componenta aleatoare a modelului linear. Valorile lui X sunt presupuse fixe.
3. Termenii eroare corespunzători observațiilor sunt necorelați. În plus, pentru orice valori date pentru X , erorile sunt variabile aleatoare normal distribuite cu media zero și dispersie constantă.

Căutăm cei mai buni estimatori pentru intersecția cu axa Ox (β_0) și panta dreptei (β_1), în cazul modelului regresiei liniare simple dat prin ecuația:

$$Y_i = \beta_0 + \beta_1 \cdot X_i + \varepsilon_i$$

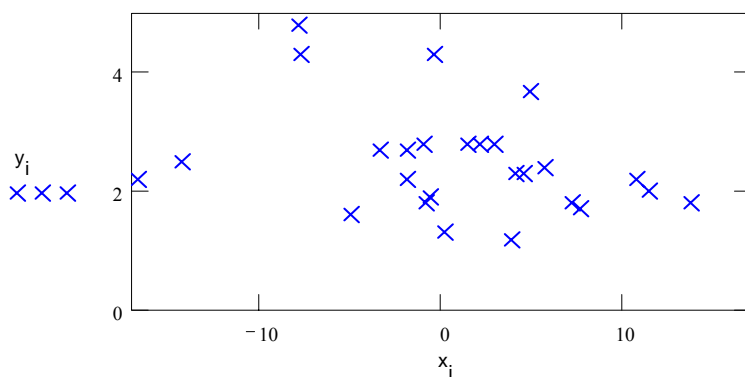
Termenul ε_i pentru observația, i , este diferența măsurată pe verticală, dintre punctele observate (X_i, Y_i) și linia de regresie. Valoarea lui ε va fi pozitivă când valorile observate se găsesc deasupra liniei de regresie și negativă când se găsesc sub linie.



Un exemplu: regresia inflației (inf) față de producția industrială (pi)
 Reprezentarea grafică:

```

y := inf          x := pi
n := rows(pi)    i := 0.. n - 1
    
```



Cum găsim această dreaptă?
 Prin metoda celor mai mici pătrate.

Suma patratelor Residualelor, SSE

Vom folosi **suma patratelor rezidualelor**

$$SSE(b_0, b_1) := \sum_i [y_i - (b_0 + b_1 \cdot x_i)]^2$$

Putem obține o soluție folosind metoda de rezolvare a sistemelor de ecuații în Mathcad

Valori initiale: $b_0 := 1$ $b_1 := 1$

Blocul ecuatiilor

Given

$$n \cdot b_0 + b_1 \cdot (\sum x) = \sum y$$

$$b_0 \cdot (\sum x) + b_1 \cdot (\sum x^2) = \sum (x \cdot y)$$

$$\begin{bmatrix} b_0 \\ b_1 \end{bmatrix} := \text{Find}(b_0, b_1)$$

produce estimatorii pentru β_0 si β_1

$$b_0 = 2.525 \quad b_1 = -0.04$$

pentru cele mai mici posibile sume de patrate,

$$\text{SSE}(b_0, b_1) = 18.58$$

$$\text{SSE}(b_0 - 0.01, b_1 - 0.01) = 18.717$$

Metoda celor mai mici patrate produce estimatorii cei mai buni - the **best linear unbiased estimators** (BLUE) - pentru coeficienti. Asta inseamna ca au dispersia minima.

In Mathcad sunt dati de functiile: $\text{intercept}(x, y)$ $\text{slope}(x, y)$

$$\text{intercept}(x, y) = 2.525 \quad \text{slope}(x, y) = -0.04$$

unde primul argument reprezinta variabila independenta. Spre deosebire de functia de corelatie ordinea argumentelor este esentiala !

$$\text{intercept}(x, y) \neq \text{intercept}(y, x)$$

$$\text{slope}(x, y) \neq \text{slope}(y, x)$$

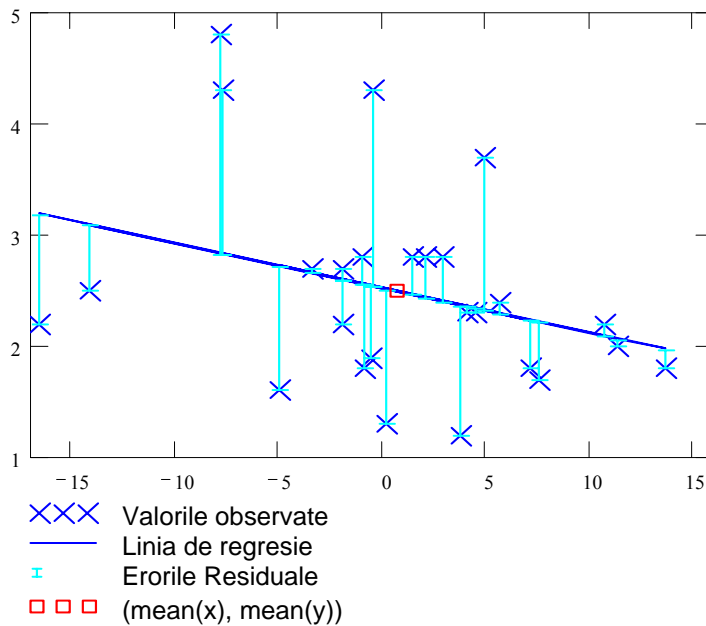
Linia de regresie pentru exemplu dat are ecuatia

$$y_{\text{lin}}(x) := b_0 + b_1 \cdot x$$

se vede pe grafic si de asemenea punctele observate si rezidualele. Punctul dat de valorile medii ale variabilelor aleatoare:

$$(\text{mean}(x), \text{mean}(y))$$

(indicat de un patrat) este de asemenea pe linia de regresie.



Estimatori pentru Dispersia Erorilor

Suma patratelor rezidualelor SSE poate fi folosita in estimarea dispersiei erorilor σ^2 . Aceasta dispersie masoara imprastierea datelor observate fata de linia de regresie.

Un estimator pentru σ^2 in cazul unui esantion de lungime n este **mean square error** (MSE), care se calculeaza impartind valoarea lui SSE prin $(n - 2)$

$$MSE := \frac{SSE(b_0, b_1)}{n - 2} \quad MSE = 0.774$$

Aceasta definitie urmareste ideea ca dispersia este media sumei patratelor deviatiei fata de linie. Impartim prin $(n - 2)$ deoarece am estimat deja cei doi parametrii β_0 si β_1 . Un estimator pentru abaterea standard a erorilor σ este **standard error of estimate**

$$se_e := \sqrt{MSE} \quad se_e = 0.88$$

$$x := pi \quad y := inf$$

Erorile reziduale, e , sunt date de:

$$y_{lin}(x, y) := intercept(x, y) + slope(x, y) \cdot x$$

$$e := y - y_{lin}(x, y)$$

Si estimatorul pentru abaterea standard este deci

$$se_e(e) := \sqrt{\frac{\left(\sum e^2\right)}{rows(e) - 2}} \quad se_e(e) = 0.88$$

Reziduale Standardizate

$$\text{standard_e}(e) := \frac{e}{\text{se_e}(e)} \quad \text{standard_e}(e)_1 = 0.119$$

Reziduale Studentizate

Un alt estimator pentru abaterea standard s a unei variabile X este:

$$i := 0..n-1 \quad \text{leverage}_i := \frac{1}{\text{rows}(x)-1} \cdot \frac{(x_i - \text{mean}(x))^2}{\text{Var}(x)} \quad \text{leverage}_0 = 0.055$$

Ajustind rezidualele standardizate obtinem **reziduale studentizate**

$$\text{student_e}_i := \frac{e_i}{\text{se_e}(e)} \sqrt{\left[1 - \left(\frac{1}{\text{rows}(x)} + \text{leverage}_i \right) \right]}$$

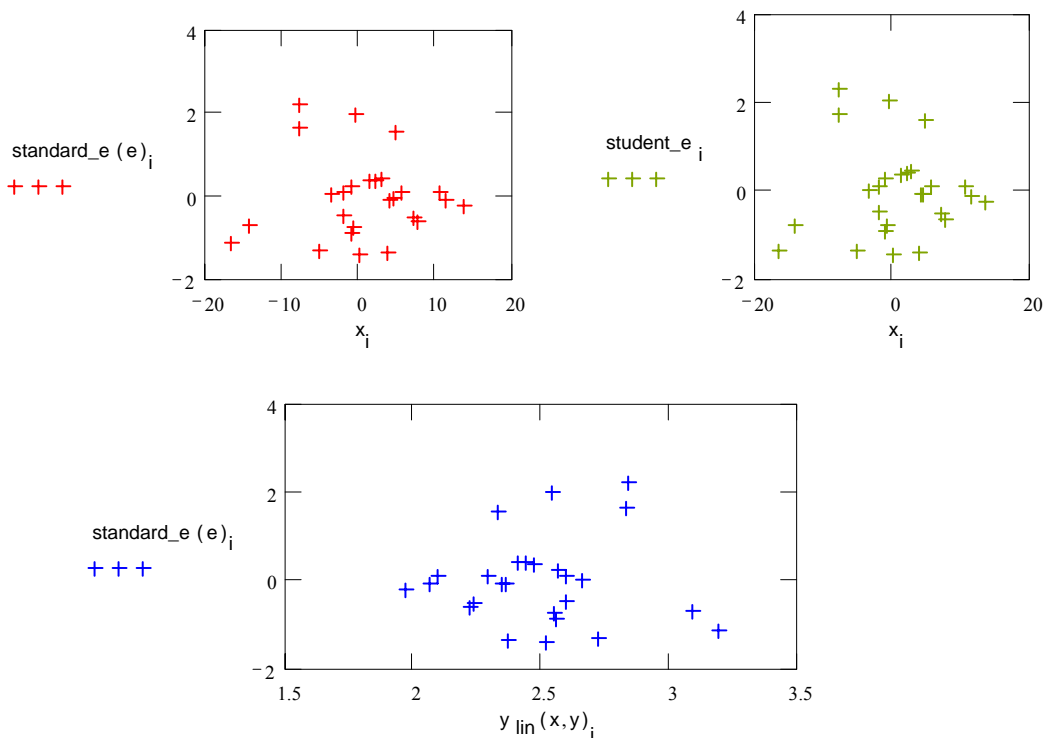
unde membrul drepteste un estimator pentru abaterea standard a rezidualelor e_i .

Rezidualele studentizate sunt mai precise in sensul ca indica mai exact diferentele in dispersia erorilor .

$$\text{standard_e}(e)_1 = 0.119 \quad \text{student_e}_1 = 0.127$$

Graficul Rezidualelor

Daca datele urmeaza o relatie liniara atunci graficul rezidualelor standardizate (sau studentizate) axa (y-axis) fata de valorile X sau fata de valorile previzionate de model nu vor prezenta nici un model evident ci vor fi repartizate aleator.



Daca datele nu respecta modelul liniar reziduale vor prezenta grafic un model(pattern) . De exemplu daca generam aleator un esantion de 100 valori dintr-o o serie de date care urmeaza o functie exponentiala:

```
n_esantion := 100      nr := 20
k := 0.. n_esantion - 1   x_exp_k := rnd(nr)
```

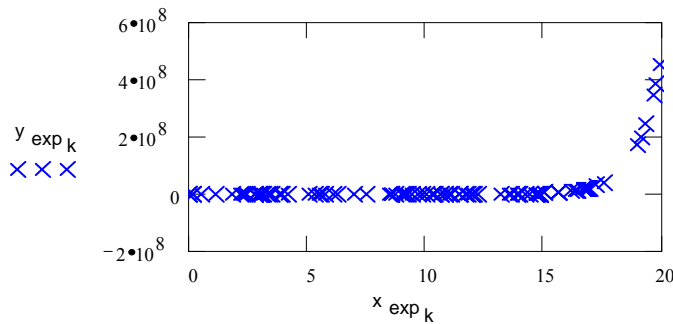
Si de asemenea erori normale cu dispersia σ :

```
 $\sigma := 10$        $\varepsilon := \text{rnorm}(n\_esantion, 0, \sigma)$ 
```

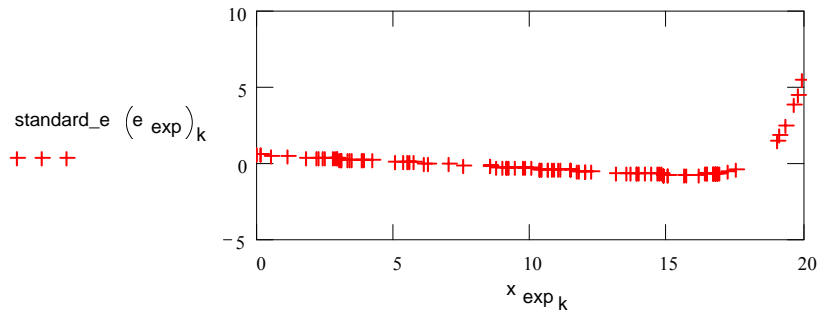
$$y_{exp} := \overrightarrow{\left(\exp(x_{exp})\right)} + \varepsilon$$

$$e_{exp} := y_{exp} - y_{lin}(x_{exp}, y_{exp})$$

obtinem urmatorul grafic



si rezidualele



Dispersia Erorilor Constanta

Daca dispersia erorilor nu este constanta de la o valoare a lui X la alta , graficul rezidualelor va arata o repartizare crescatoare sau descrescatoare de la stanga la dreapta.

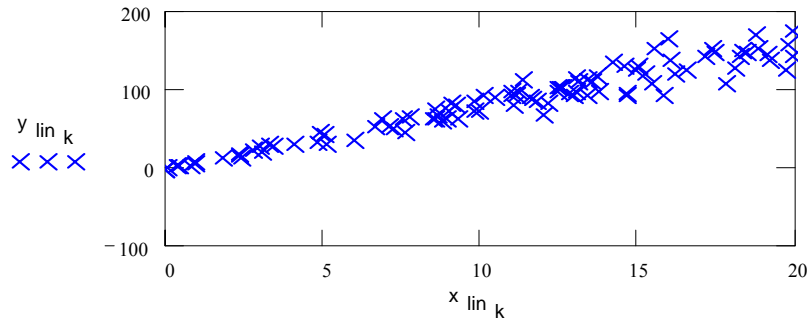
Sa vedem un exemplu de care urmeaza o relatie liniara asa cum arata graficul lui Y fata de X si cum arata si coeficientul de corelatie:

```
n_esantion := 100      nr := 20
k := 0.. n_esantion - 1   x_lin_k := rnd(nr)
```


$$\sigma := x_{lin} \quad \varepsilon_k := \text{norm} (n_{esantion}, 0, \sigma_k)_k$$

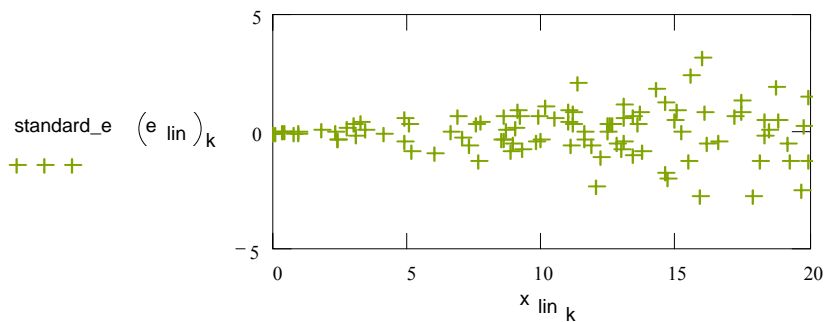
$$m := 8 \quad b := -2 \quad y_{lin} := m \cdot x_{lin} + b + \varepsilon$$

$$e_{lin} := y_{lin} - y_{lin}(x_{lin}, y_{lin})$$



$$\text{corr}(x_{lin}, y_{lin}) = 0.966$$

dispersia erorilor prezinta o tendinta de crestere de la stinga la dreapta .



Corelatia Erorilor

Putem verifica daca exista corelatie intre termenii alaturati in seria erorilor folosind statistica **Durbin-Watson (DW)**.

$$e := y - y_{lin}(x, y)$$

$$DW(e) := \frac{\sum_{d=1}^{\text{rows}(e)-1} (e_d - e_{d-1})^2}{\sum_{d=0}^{\text{rows}(e)-1} (e_d)^2}$$

$$DW(e) = 1.867$$

Valori pentru statistica Durbin-Watson mai mici decit 2 indica corelatie pozitiva pentru erori si valori mai mari decit 2 o corelatie negativa. Totusi aceasta statistica nu poate da un raspuns decit pentru corelatia termenilor alaturati.

R² masura pentru dependenta cauzala intre variabila independenta si cea dependenta

Dependenta dintre variabile se poate exprima prin covarianta

$$E[(X - \mu_x) \cdot (Y - \mu_y)]$$

care de obicei se calculeaza normalizat prin coeficientul de corelatie ρ ,

$$\rho = \frac{E[(X - \mu_x) \cdot (Y - \mu_y)]}{\sigma_x \cdot \sigma_y}$$

$$\rho_{pi.inf} := \frac{\sum_{i=0}^{rows(pi) - 1} [(pi_i - \mu_{pi}) \cdot (inf_i - \mu_{inf})]}{\sigma_{pi} \cdot \sigma_{inf}} \cdot \frac{1}{n}$$

$$\rho_{pi.inf} = -0.318 \quad \text{corr}(pi, inf) = -0.318 \quad \text{rows}(pi) = 26$$

$$y := inf \quad x := pi$$

Erorile reziduale, e , sunt date de:

$$y_{lin}(x, y) := \text{intercept}(x, y) + \text{slope}(x, y) \cdot x$$

$$e := y - y_{lin}(x, y)$$

Si estimatorul pentru abaterea standard este deci $se_e(e) := \sqrt{\frac{\sum e^2}{rows(e) - 2}}$ $se_e(e) = 0.88$

$$R2_{pi.inf} := 1 - \frac{\sum(e)^2}{\sum(y - \text{mean}(y))^2} \quad R2_{pi.inf} = 0.101$$

Ipoteza statistica $H_0: \beta_1 = 0$ se poate testa folosind statistica t

$t = \frac{b_1}{se_{b1}(x, y)}$ urmeaza o distributie t student cu $(n - 2)$ grade de libertate. Valuarea $se_{b1}(x, y)$ este abaterea standard a erorii pentru parametrul $b1$

$$se_{b1}(x, y) := \frac{se_e(e)}{\sqrt{\sum(x - \text{mean}(x))^2}} \quad se_{b1}(x, y) = 0.024$$

$$t(x, y, \beta_1) := \frac{\text{slope}(x, y) - \beta_1}{se_{b1}(x, y)} \quad t(x, y, 0) = -1.643$$

Să încercăm și o regresie multiplă. Datele folosite sunt producția industrială pi (variație lunară, %), șomajul șomaj (rată lunară, %) și inflația inf (rată lunară, %) pentru perioada ianuarie 2000 - ianuarie 2002. In ordinea dată ele se regăsesc și in matricea DATE.

Variabila dependenta inflatia, adica :

$$y := \text{DATE}^{<2>}$$

marimea esantionului: $n := \text{rows}(\text{DATE})$

Numarul coloanelor ce reprezinta variabilele indepedente(productia industriala si somajul):

$$x_col := \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad p := \text{rows}(x_col) \quad p = 2$$

Matricea X a variabilelor indepedente:

$$X := \begin{array}{l} \text{for } i \in 0..n-1 \\ \quad X_{i,0} \leftarrow 1 \\ \quad \text{for } j \in x_col \\ \quad \quad X^{<j+1>} \leftarrow \text{DATE}^{<j>} \\ X \end{array}$$

Vectorul parametrilor estimati:

$$b := (X^T \cdot X)^{-1} \cdot (X^T \cdot y)$$

Ecuatia regresiei si erorile reziduale

$$ylin := X \cdot b \quad e := y - ylin$$

Suma patratelor erorilor si abaterilor:

$$SSE := e^T \cdot e \quad \text{suma patratelor erorilor}$$

suma patratelor abaterii date de regresie:

$$SSR := (ylin - \text{mean}(y))^T \cdot (ylin - \text{mean}(y))$$

suma patratelor abaterii/deviatiei totale

$$SST := SSE + SSR$$

Gradele de libertate

Cum avem $(p + 1)$ parametri de estimat avem

$$\text{grL_erori} := n - (p + 1) \quad \text{grL_erori} = 23$$

grade de libertate asociate cu suma patratelor erorilor. Suma patratelor abaterii/deviatiei date de regresie are

$$\text{grL_regresie} := p \quad \text{grL_regresie} = 2$$

grade de libertate deoarece in formula sunt p variabile independente .

Suma patratelor abaterii/deviatiei totale are

$$\text{grL_total} := n - 1 \qquad \text{grL_total} = 25$$

grade de libertate si $\text{grL_erori} + \text{grL_regresie} = \text{grL_total}$

Pentru a estima dispersia impartim suma patratelor la nr. de grade de libertate. Media abaterilor patratice pentru erorile reziduale

$$\text{MSE} := \frac{\text{SSE}_0}{\text{grL_erori}}$$

este o estimatie pentru dispersia erorilor σ^2 , sau abaterea neexplicata de regresie. (Vom identifica SSE ca fiind primul element din vector.) Alti doi estimatori pentru dispersia care este explicata de modelul regresiei respectiv dispersia totala sunt:

$$\text{MSR} := \frac{\text{SSR}_0}{\text{grL_regresie}} \qquad \text{MST} := \frac{\text{SST}_0}{\text{grL_total}}$$

Testul F

In cazul ipotezei nule

H_0 : nu avem model de tip regresie multipla

statistica $F := \frac{\text{MSR}}{\text{MSE}}$ are o distributie F cu $n1 := \text{grL_regresie}$ si $n2 := \text{grL_erori}$

grade de libertate.

Valuarea p-value a testului este data de :

$$\text{p_val} := 1 - \text{pF}(F, n1, n2) \qquad \text{p_val} = 0.063$$

$R2 := \frac{\text{SSR}_0}{\text{SST}_0}$ si $R2_ajustat := 1 - \frac{\text{MSE}}{\text{MST}}$ dau masura potrivirii modelului de regresie liniara.

Analiza Abaterilor

GradeLibertate	SS	MS	F	R2
grL_regresie = 2	SSR = 4.408	MSR = 2.204	F = 3.117	R2 = 0.213
grL_erori = 23	SSE = 16.262	MSE = 0.707	p-value	R2_ajustat = 0.145
grL_total = 25	SST = 20.67	MST = 0.827	p_val = 0.063	

Abaterea explicata de regresia liniara (MSR) este mai mare decit cea datorata erorilor reziduale (MSE). Diferenta este destul de mare (valoarea p-value este suficient de mica) pentru a respinge ipoteza nula.

Corelatia dintre variabilele modelului este data in matricea CORR

$$j := 0..p \quad k := 0..p \quad \text{CORR}_{j,k} := \text{corr}(\text{DATE}^{<j>}, \text{DATE}^{<k>})$$

$$\text{CORR} = \begin{matrix} & \begin{matrix} x1 & x2 & x3 \end{matrix} \\ \begin{matrix} x1 \\ x2 \\ x3 \end{matrix} & \begin{bmatrix} 1 & 0.268 & -0.318 \\ 0.268 & 1 & 0.237 \\ -0.318 & 0.237 & 1 \end{bmatrix} \end{matrix}$$

Cea mai mare valoare a corelatiei dintre variabilele independente este cea intre variabilele ($x1$ and $x2$)

$$\text{CORR}_{1,0} = 0.268$$

Teste T

$$b = \begin{bmatrix} 0.522 \\ -0.052 \\ 0.196 \end{bmatrix} \quad \text{unde} \quad b = \begin{bmatrix} b_0 \\ b_1 \\ b_2 \end{bmatrix}$$

Forma statisticii t este aceeași ca și în cazul regresiei simple. În ipoteza nulă, $t = \frac{b_1}{se_{b1}}$

urmează o distribuție t student cu $(n - 2)$ grade de libertate. Valoarea se_{b1} este abaterea standard a erorii pentru parametrul $b1$

Matricea Varianta-Covariantă pentru Parametrii Estimati

Fiecare estimatie are o dispersie/varianța dar și o covarianță cu fiecare dintre celelalte estimatii. Acestea se pot reprezenta într-o matrice Var_Covar_b data de $\sigma^2 \cdot (X^T \cdot X)^{-1}$

Putem estima varianța necunoscută σ^2 prin abaterea medie pătratică "mean square error: (MSE),

$$\text{Var_Covar_b} := (X^T \cdot X)^{-1} \cdot \text{MSE}$$

$$\text{Var_Covar_b} = \begin{matrix} & \begin{matrix} b_0 & b_1 & b_2 \end{matrix} \\ \begin{matrix} b_0 \\ b_1 \\ b_2 \end{matrix} & \begin{bmatrix} 1.251524 & 0.006805 & -0.119521 \\ 0.006805 & 0.000589 & -0.000703 \\ -0.119521 & -0.000703 & 0.01167 \end{bmatrix} \end{matrix}$$

$\text{Var_Covar_b}_{0,1} = 6.805 \cdot 10^{-3}$ este covarianța dintre b_0 și b_1 . Pe diagonală matricii este varianța/dispersia estimată a coeficienților.

$\text{Var_Covar_b}_{1,1} = 5.891 \cdot 10^{-4}$ este o estimatie a varianței/dispersiei lui b_1 .

Erorile Standard ale Parametrilor Estimati

Prin definitie vectorul cu erorilor standard ale parametrilor estimati poate fi obtinuta astfel:

$$k := 0..p$$

$$se_b_k := \sqrt{\text{Var_Covar_b}_{k,k}} \quad se_b = \begin{bmatrix} 1.119 \\ 0.024 \\ 0.108 \end{bmatrix}$$

Statistica T

$$t := \frac{\vec{b}}{se_b} \quad t = \begin{bmatrix} 0.467 \\ -2.142 \\ 1.81 \end{bmatrix}$$

Considerind ipoteza nula,

$$H_0: \beta_k = 0$$

$k = 0, 1, 2, \dots, p$, fiecare test statistic urmeaza o distributie t student cu $n - (p + 1) = 23$ grade de libertate, egal cu nr. de grade de libertate rezidualelor $grL_erori = 23$

Valuarea corespunzatoare p-value poate fi calculata astfel:

$$p_t_k := \begin{cases} 2 \cdot (1 - pt(t_k, grL_erori)) & \text{if } t_k > 0 \\ 2 \cdot pt(t_k, grL_erori) & \text{otherwise} \end{cases}$$

Tabelul Coeficientilor Regresiei

$$\begin{bmatrix} b_0 \\ b_1 \\ b_2 \end{bmatrix} \quad b = \begin{bmatrix} 0.522 \\ -0.052 \\ 0.196 \end{bmatrix} \quad se_b = \begin{bmatrix} 1.119 \\ 0.024 \\ 0.108 \end{bmatrix} \quad t = \begin{bmatrix} 0.467 \\ -2.142 \\ 1.81 \end{bmatrix} \quad p_t = \begin{bmatrix} 0.645 \\ 0.043 \\ 0.083 \end{bmatrix}$$

Pentru un nivel de semnificatie $\alpha_j := 0.05$

vom respinge ipoteza nula privind neimportanta pentru model a variabilei pentru care conditia este 1

$$cond_j := (\alpha_j \geq p_t_j) \quad cond = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \quad \begin{bmatrix} intercept \\ x1 \\ x2 \end{bmatrix}$$

Daca conditiile impuse de model sunt indeplinite si nu avem multicolaritate, cu un nivel de semnificatie $\alpha = 0.05$ vom pastra in model variabila $x1$ si vom exclude $x2$.

Daca incercam un model de tip $y = \beta_0 + \beta_1 \cdot x_1 + (\beta_1)^2 \cdot x_2 + \varepsilon$

$x_1 := \text{DATE}^{<0>}$ $n := \text{rows}(\text{DATE}^{<0>})$

$x_2 := \text{DATE}^{<1>}$ $i := 0..n-1$

$y := \text{DATE}^{<2>}$

Modelul este $y = [\beta_0 + \beta_1 \cdot x_1 + (\beta_1)^2 \cdot x_2 + \varepsilon]$

este intrinsec nelinear deoarece parametrii apar ca puteri, deci neliniar.

In aceste cazuri intrinsec neliniare putem aplica metoda celor mai mici patrate.

$b_0 := 0$ $b_1 := 0$

Given

$$\sum_i (-2 \cdot y_i + 2 \cdot b_0 + 2 \cdot b_1 \cdot x_{1i} + 2 \cdot b_1^2 \cdot x_{2i}) = 0$$

$$\sum_i [2 \cdot (y_i - b_0 - b_1 \cdot x_{1i} - b_1^2 \cdot x_{2i}) \cdot ((-x_{1i} - 2 \cdot b_1 \cdot x_{2i}))] = 0$$

$$\begin{bmatrix} b_0 \\ b_1 \end{bmatrix} := \text{Find}(b_0, b_1)$$

$b_0 = 2.508$ $b_1 = -0.041$

si calculam $SSY := \sum y^2$ $SSY = 182.67$ cu n grade de libertate $grL_{total} := n$

$$SSE := \sum \left[y - [b_0 + b_1 \cdot x_1 + (b_1)^2 \cdot x_2] \right]^2$$

$SSE = 18.54$ cu n-p grade de libertate $grL_{erori} := n - p$

Pentru variatia explicata de model avem SSR,

$SSR := SSY - SSE$ $SSR = 164.13$ cu p grade de libertate $grL_{regresie} := p$

$$MSE := \frac{SSE}{grL_{erori}} \quad MSR := \frac{SSR}{grL_{regresie}}$$

Analiza tabelului variatiilor

Grade Libertate	SS	MS
$grL_{regresie} = 2$	$SSR = 164.13$	$MSR = 82.065$
$grL_{erori} = 24$	$SSE = 18.54$	$MSE = 0.773$
$grL_{total} = 26$	$SSY = 182.67$	

Matricea Varianta-Covarianta

Jacobianul:

$$J_{i,0} := -\frac{d}{db_0} [b_0 + b_1 \cdot x1_i + (b_1)^2 \cdot x2_i]$$

$$J_{i,1} := -\frac{d}{db_1} [b_0 + b_1 \cdot x1_i + (b_1)^2 \cdot x2_i]$$

si inmultind cu media patratelor erorilor obtinem matricea

$$\text{Var_Covar} := \text{MSE} \cdot (J^T \cdot J)^{-1}$$

$$\text{Var_Covar} = \begin{bmatrix} 0.026 & -2.869 \cdot 10^{-4} \\ -2.869 \cdot 10^{-4} & 5.633 \cdot 10^{-4} \end{bmatrix}$$

Luam apoi radacina patratica din elementele de pe diagonala pentru a obtine erorile standard pentru parametrii estimati

$$j := 0..p-1$$

$$\text{se}_j := \sqrt{\text{Var_Covar}_{j,j}} \quad \text{se} = \begin{bmatrix} 0.161 \\ 0.024 \end{bmatrix}$$

ca apoi sa formam intervalele de confidenta pentru b_0 ,

$$\alpha_t := 0.05$$

$$t := \text{qt} \left(1 - \frac{\alpha_t}{2}, \text{grL_erori} \right)$$

$$\text{int_stinga } b_0 := b_0 - t \cdot \text{se}_0 \quad \text{int_dreapta } b_0 := b_0 + t \cdot \text{se}_0$$

si apoi pentru b_1 ,

$$\text{int_stinga } b_1 := b_1 - t \cdot \text{se}_1 \quad \text{int_dreapta } b_1 := b_1 + t \cdot \text{se}_1$$

Avem $1 - \alpha_t = 95\%$ siguranta, ca valoarea lui β_0 se afla intre valorile

$$\text{int_stinga } b_0 = 2.176 \quad \text{si} \quad \text{int_dreapta } b_0 = 2.84$$

si ca valoarea lui β_1 se afla intre valorile

$$\text{int_stinga } b_1 = -0.09 \quad \text{si} \quad \text{int_dreapta } b_1 = 7.929 \cdot 10^{-3}$$

Parametrii Estimati	Erorile Standard	Intervalele de incredere
---------------------	------------------	--------------------------

$b_0 = 2.508$	$\text{se}_0 = 0.161$	$\text{int_stinga } b_0 = 2.176$ $\text{int_dreapta } b_0 = 2.84$
---------------	-----------------------	---

$b_1 = -0.041$	$\text{se}_1 = 0.024$	$\text{int_stinga } b_1 = -0.09$ $\text{int_dreapta } b_1 = 7.929 \cdot 10^{-3}$
----------------	-----------------------	--

Cum intervalul pentru b_1 contine valoarea 0 concluzionam ca acest parametru nu este semnificativ pentru model. Deci modelul nu este potrivit.

Bibliografie

- | | |
|---|---|
| <i>Robert S Pindyck, Daniel L. Rubinfeld</i> | MathCad Resource Center |
| <i>McGraw-Hill,Inc., International Edition 1991</i> | Econometric Models and Econometric Forecasts, |
| <i>Edmond Malinvaud</i> | Methodes statistiques de l'econometrie, editura |
| <i>Dunod, 1964</i> | |
| <i>George Daniel Mateescu</i> | Bazele utilizarii calculatoarelor, Editura Donaris |
| <i>2004</i> | |