

## Statistică multivariată

### Lucrarea nr. 10 — Regresia logistică - SPSS

#### A. Noțiuni teoretice

##### Regresia logistică

Regresia logistică modelează relația dintre o mulțime de variabile independente  $x_i$  (catoriale, continue) și o variabilă dependentă dihotomică (nominală, binară)  $Y$ . O astfel de variabilă dependentă apare, de regulă, atunci când reprezintă apartenența la două clase, categorii – prezență/absență, da/nu etc.

Ecuția de regresie obținută, de un tip diferit de celelalte regresii discutate, oferă informații despre:

- importanța variabilelor în diferențierea claselor,
- clasificarea unei observații într-o clasă.

De remarcat că diagrama de împrăștiere a valorilor nu oferă nici un indiciu în privința dependențelor. În asemenea cazuri, regresia liniară clasică nu oferă un model adecvat.

Presupunem că valorile  $y$  (variabilă binară) sunt codificate 0/1, valoarea 1 exprimând în general apariția unui anumit eveniment, astfel încât ceea ce se caută este o estimare a probabilității de producere a respectivului eveniment în funcție de valorile variabilelor independente.

##### Cazul unei singure variabile independente

Modelul este

$$P(y=1|x) = \frac{e^{\alpha+\beta x}}{1 + e^{\alpha+\beta x}}$$

sau

$$\ln\left(\frac{P(y=1|x)}{1 - P(y=1|x)}\right) = \alpha + \beta x$$

Cantitatea din partea stângă este numită (**transformarea**) **logit** a probabilității  $P(y=1|x)$ .

Semnificația expresiei  $P(y=1|x)$  este evidentă: probabilitatea de realizare a valorii  $y=1$  condiționată de valoarea  $x$ . Cu alte cuvinte, probabilitatea de clasare a observației  $x$  în clasa  $y=1$ , sau probabilitatea ca valoarea  $x$  să fie asociată cu producerea evenimentului  $y=1$ . În continuare se notează  $P(y=1|x)$  cu  $p$ , conform notației de la modelul probabilist binomial (probabilitatea de “succes”).

Transformarea logit este necesară pentru a proiecta probabilitatea  $p$  din intervalul  $(0,1)$  în intervalul  $(-\infty, +\infty)$ , fapt necesar în procesul de estimare a parametrilor. Modelul este legat direct de noțiunea de *odds* (raport de șanse), notat OR (*odds report*):

$$OR = \frac{p}{1 - p}$$

care reprezintă raportul dintre probabilitatea de « succes » și probabilitatea de « insucces ».

Modelul se mai poate scrie

$$\frac{p}{1-p} = e^{\alpha + \beta x}$$

de unde interpretarea coeficientului  $\beta$ :

- creșterea cantității logit atunci când  $x$  crește cu o unitate sau
- OR crește de  $e^\beta$  ori atunci când  $x$  crește cu o unitate.

Testarea ipotezei  $\beta = 0$  se realizează prin testul Wald, corespunzător testului  $t$  de la regresia liniară, statistica testului fiind

$$\chi^2 = \frac{b^2}{\text{Var}(b)}$$

care este repartizată  $\chi^2$  cu un singur grad de libertate.

Intervalul de încredere pentru  $\beta$  este, potrivit rezultatelor de la analiza ecuației de regresie,

$$\left( e^{b - z_{1-\frac{\alpha}{2}} SE(b)}, e^{b + z_{1-\frac{\alpha}{2}} SE(b)} \right),$$

unde  $b$  este estimația lui  $\beta$  (din ecuația de regresie estimată) iar  $SE(b)$  este abaterea standard a repartiției de sondaj a lui  $b$ .

Se observă imediat că, pentru o observație, dacă  $p > 0,5$ , atunci este mai probabil ca observația să aparțină grupului caracterizat de  $y=1$ . Această condiție este echivalentă cu  $OR > 1$ , adică  $\text{logit} > 0$ .

### Cazul mai multor variabile independente

Modelul general este

$$\ln \left( \frac{p}{1-p} \right) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k,$$

unde  $p$  este  $P(y = 1 | x_1, x_2, \dots, x_k)$ . Se poate obține imediat și forma exponențială echivalentă.

Interpretarea coeficienților  $\beta_i$  este evidentă: creșterea cantității logit (logaritm din OR) atunci când  $x_i$  crește cu o unitate (celelalte variabile  $x$  rămânând constante). Pentru interpretări mai sofisticate rescriem modelul sub forma:

$$P(y = 1 | x_1, x_2, \dots, x_k) = \frac{\exp(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k)}{1 + \exp(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k)}$$

Se obține atunci, după calcule imediate,

$$\begin{aligned} \exp(\beta_0) &= \frac{P(y = 1 | x_1 = x_2 = \dots = x_k = 0)}{1 - P(y = 1 | x_1 = x_2 = \dots = x_k = 0)} = \\ &= \frac{P(y = 1 | x_1 = x_2 = \dots = x_k = 0)}{P(y = 0 | x_1 = x_2 = \dots = x_k = 0)} \end{aligned}$$

adică OR în situația de bază  $x_1 = x_2 = \dots = x_k = 0$ .

Pentru coeficientul  $\beta_i$  se obține :

$$\begin{aligned} \exp(\beta_i) &= \frac{P(y=1 | x_i=1, x_j=0 \text{ pentru } j \neq i)}{1 - P(y=1 | x_i=1, x_j=0 \text{ pentru } j \neq i)} \times \frac{1}{OR_{\text{baza}}} = \\ &= \frac{OR_{x_i=1, x_j=0 \text{ pentru } j \neq i}}{OR_{\text{baza}}}. \end{aligned}$$

Se ajunge astfel, din caracterul multiplicativ al modelului logistic,

$$OR_{x_1, x_2, \dots, x_k} = \exp(\beta_0) \times \exp(\beta_1 x_1) \times \dots \times \exp(\beta_k x_k),$$

la interpretarea utilă că fiecare  $\beta_i$  exprimă contribuția factorului  $x_i$  la explicarea probabilității (sub forma OR) de producere a evenimentului  $y=1$ . Astfel, fixând  $x_i=1$ ,  $\exp(\beta_i)$  va reprezenta factorul multiplicativ constant indiferent de valorile celorlalte variabile independente.

Dacă  $\beta_i = 0$ , factorul corespunzător nu are nici un efect, (înmulțirea cu 1). Dacă  $\beta_i < 0$  prezența factorului reduce probabilitatea evenimentului  $y=1$ ,  $\beta_i > 0$  măbind această probabilitate.

Construirea modelului se poate realiza și prin metode forward sau backward, testarea semnificației coeficienților realizându-se prin testul Wald sau prin testul raportului de verosimilitate (LR, *likelihood-ratio*).

Testul Wald este prezentat la modelul logistic cu un singur factor. Testul LR se bazează pe statistica obținută ca raport între maximul funcției de verosimilitate sub ipoteza nulă și maximul funcției de verosimilitate în condiții mai largi. Lema Neyman-Pearson arată că acesta este cel mai puternic test la un prag  $\alpha$  fixat. Pentru cazul regresiei logistice, se calculează raportul între valoarea maximă a funcției de verosimilitate pentru modelul complet ( $L_1$ ) și cea pentru modelul mai simplu ( $L_0$ ). Statistica LR este  $-2\log(L_0/L_1)$ , repartizată  $\chi^2$ . Testul LR este recomandat în cazul construirii modelului pas cu pas, verificând dacă variabila eliminată din model este semnificativă, deci dacă modelul poate fi simplificat.

**Observație.** O mai bună imagine intuitivă asupra raportului de verosimilitate este dată în continuare. presupunem că se dorește distingerea între două ipoteze  $H_0$  și  $H_1$  (o contrară a lui  $H_0$ ). Fie  $p_0$  probabilitatea ca datele observate să apară în ipoteza  $H_0$  adevărată și  $p_1$  probabilitatea ca datele observate să apară în ipoteza  $H_1$  adevărată. Raportul  $p_1/p_0$  este raportul de verosimilitate (LR) și măsoară OR (odds report) ca  $H_1$  să fie adevărat ca opusă lui  $H_0$  adevărată.

Deoarece unele simulări arată că datorită datelor "rare" (*sparse*) statistica prin care se compară două modele nu este repartizată  $\chi^2$  și, din acest motiv, s-a dezvoltat testul Hosmer-Lemeshow. De notat că testul este recomandat pentru variabile independente continue și mai mult de 400 de observații. Testul constă în clasificarea în decile a probabilităților prognozate (10 grupuri bazate pe rangul percentilic) și calcularea statisticii  $\chi^2$  care compară frecvențele observate cu cele prognozate (în tabelul  $2 \times 10$ ). Valori mici ale statisticii (deci acceptarea nediferențierii dintre cele două șiruri de frecvențe) arată o bună potrivire a datelor prognozate, deci o adecvanța modelului.

În regresia logistică nu există un indicator absolut similar coeficientului  $R^2$  din regresia liniară. S-au dezvoltat însă indicatori similari. Astfel în SPSS există Cox & Snell Pseudo- $R^2$  definit prin

$$R^2 = 1 - \left[ \frac{-2LL_{null}}{-2LL_k} \right]^{2/n}$$

unde  $LL_{null}$  este logaritmul din maximul funcției de verosimilitate pentru modelul constant, iar  $LL_k$  este logaritmul din maximul funcției de verosimilitate pentru modelul cu variabile independente incluse. Se poate astfel observa că se merge pe varianta de comparare a cantităților  $-2LL$  prin intermediul raportului lor și nu a împărțirii lor (ca la LR). Acest  $R^2$  nu atinge 1 și a fost introdusă de Nagelkerke o modificare prin care se atinge 1. Formula pentru Nagelkerke Pseudo- $R^2$  este

$$R^2 = \frac{1 - \left[ \frac{-2LL_{null}}{-2LL_k} \right]^{2/n}}{1 - (-2LL_{null})^{2/n}}$$

Alți indicatori sunt:

- AIC (*Akaike's Information Criterion*) definit ca  $-2LL_k + 2k$ , unde  $k$  este numărul de parametri estimați.
- BIC (*Bayesian Information Criterion*) definit ca  $-2LL_k + k \cdot \log(n)$  unde  $k$  este numărul de parametri estimați iar  $n$  este numărul de observații. BIC mai este referit și drept criteriul Schwartz (care l-a argumentat).

Vor fi preferate modelele pentru care criteriile (AIC sau BIC) au valori mai mici. Se observă că ambele criterii "recompensează" buna potrivire a modelului dar și "penalizează" numărul de parametri estimați, astfel încât să se obțină un model bun dar cu un număr minim de parametri. În BIC, penalizarea lui  $k$  este mai puternică decât în AIC. Ambii indicatori necesită condiția ca erorile (reziduurile) să fie normal distribuite.

### Regresia logistică multinomială

Modelul regresional logistic multinomial (cunoscut și ca regresia logistică polinomială – *polytomous logistic regression* – sau ca model de alegere discretă – *discrete choice model* – în econometrie) este o generalizare a modelului logistic acceptând ca variabila dependentă  $Y$  să aibă mai mult de două valori.

Să presupunem că variabila  $Y$  are ca valori posibile elementele mulțimii neordonate  $\{1, \dots, g\}$ . Modelul logistic multinomial presupune că probabilitatea ca  $Y$  să fie egal cu  $s$  în observația  $i$  depinde de valorile variabilelor  $x_{i1}, \dots, x_{ip}$  prin

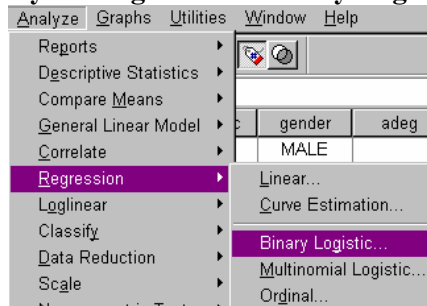
$$P(Y_i = s) = \frac{e^{\eta_{is}}}{\sum_{t=1}^g e^{\eta_{it}}}$$

unde  $\eta_{is} = \sum_{k=1}^p x_{ik} \beta_{ks}$  este o funcție liniară. În această formulare a modelului, este de remarcat că există coeficienți de regresie  $\beta_{ks}$  diferiți pentru fiecare  $k$  și, mai ales,  $s$ . Prin urmare, fiecare valoare posibilă  $Y$  are un model asociat.

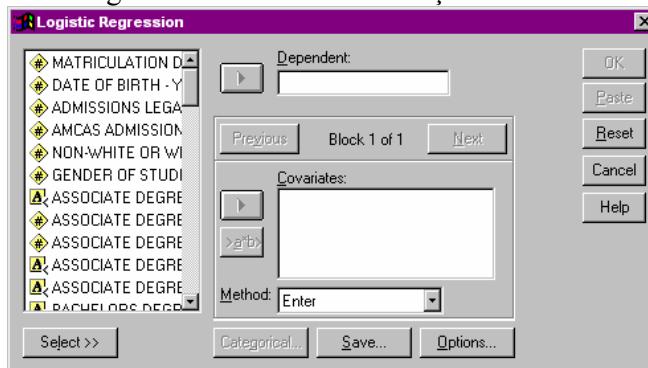
Modelul astfel definit este supraparametrizat, ceea ce impune o reducere prin fixarea unei valori  $Y$ , de exemplu  $Y = 1$ , drept categorie de referință (adică  $\beta_{11}, \dots, \beta_{p1}$  sunt egali cu zero). Alegerea categoriei de referință poate facilita interpretarea.

## B. Instrumente SPSS

Comanda este **Analyse - Regression - Binary Logistic**.

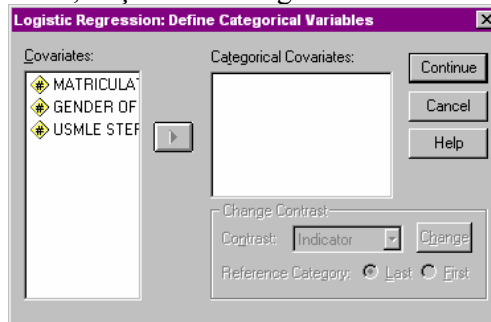


Se afișează dialogul de fixare a variabilelor și statisticilor.

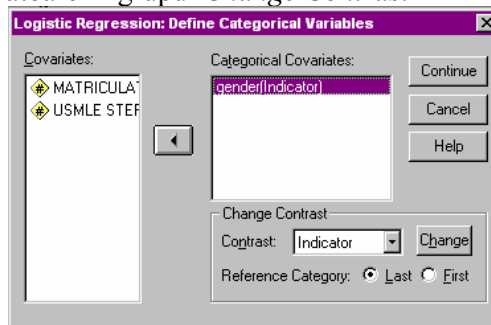


Se mută variabila dependentă (binară) în *Dependent*. Variabila independentă sau variabilele independente (în cazul multivariat) sunt mutate în lista *Covariates*.

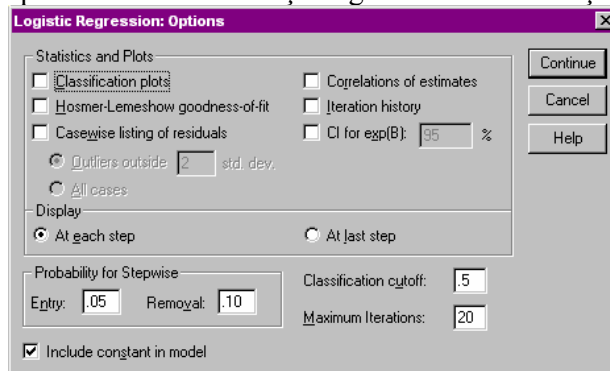
Pentru a indica variabilele independente care este categoriale (discrete), se va acționa butonul *Categorical*, afișându-se dialogul



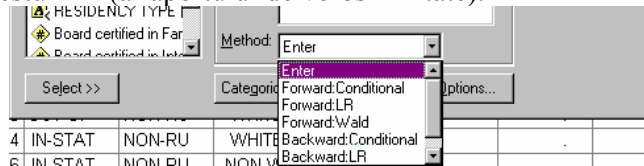
Fiecare variabilă trecută în lista *Categorical Covariates* poate fi caracterizată prin selecții corespunzătoare în grupul *Change Contrast*



Acționând butonul *Options* din dialogul principal, se deschide dialogul sinonim în care se precizează statisticile și diagramele dorite în ieșire.



În dialogul principal *Logistic Regression*, se poate alege metoda utilizată pentru introducerea variabilelor la estimarea regresiei. De reținut metoda *Enter* în care variabilele sunt introduse în bloc (se estimează o singură ecuație) sau metode de selectare pas cu pas (ca la regresia liniară multiplă), cum ar fi *Forward: LR*. Aceasta înseamnă că modelul este construit ascendent, criteriul de introducere a unei noi variabile fiind testul LR (a raportului de verosimilitate).



Prin acționarea butonului *Save* în dialogul principal se pot preciza noile variabile care pot fi create din ieșirea procedurii, ca și la regresia multiplă.

Informațiile care apar în fișierul de ieșire SPSS sunt explicate în continuare.

Un prim tabel cu informațiile generale (număr de observații valide etc.).

### Logistic Regression

Case Processing Summary

Unweighted Cases <sup>a</sup>		N	Percent
Selected Cases	Included in Analysis	370	33.1
	Missing Cases	747	66.9
Total		1117	100.0
Unselected Cases		0	.0
Total		1117	100.0

a. If weight is in effect, see classification table for the total number of cases.

Un tabel în care se precizează codificările variabilelor categoricale (inclusiv cea dependentă). Pentru variabilele categoricale independente are loc o recodificare cu considerarea categoriei de referință: aceasta este recodificată 0.

Dependent Variable Encoding

Original Value	Internal Value
0 no	0
1 yes	1

Categorical Variables Codings

		Frequency	Parameter (1)
GENDER	1 FEMALE	109	1.000
OF STUDENT	2 MALE	261	.000

Ieșirea diferă ca structură după metoda de selectare a variabilelor, dar conține un prim bloc de informații care se referă la modelul simplu (doar cu termenul constant). De remarcat structura: clasificare, variabile în ecuație, variabile candidate.

**Block 0: Beginning Block**

Classification Table<sup>a,b</sup>

Observed		Predicted		
		Board certified in FM, IM, Peds or Ob-Gyn		Percentage Correct
		no	yes	
Step 0	Board certified in FM, IM, Peds or Ob-Gyn	no	yes	
		309	0	100.0
		61	0	.0
Overall Percentage				83.5

a. Constant is included in the model.

b. The cut value is .500

Tabelul de clasificare este construit prin considerarea probabilității de clasificare prognozate de modelul curent pentru fiecare observație. după principiul că  $OR > 1$  clasează observația în grupul codificat 1. Un model bun trebuie să numere cele mai multe observații pe diagonala principală a tabelului.

Tabelul care urmează, referitor la model, este explicat și se interpretează potrivit celor spuse la tabelul similar dintr-un pas intermediar afișat ceva mai departe în lucrare.

Variables in the Equation

	B	S.E.	Wald	df	Sig.	Exp(B)	
Step 0	Constant	-1.622	.140	134.090	1	.000	.197

Variables not in the Equation

Step	Variables	Score	df	Sig.
0	MATYR	78.766	1	.000
	GENDER(1)	.865	1	.352
	USMLE2	18.236	1	.000
Overall Statistics		91.356	3	.000

Informațiile oferite pentru faza finală sunt după structura

**Block 1: Method = Forward Stepwise (Likelihood Ratio)**

Omnibus Tests of Model Coefficients

	Chi-square	df	Sig.	
Step 1	Step	70.688	1	.000
	Block	70.688	1	.000
	Model	70.688	1	.000
Step 2	Step	10.984	1	.001
	Block	81.672	2	.000
	Model	81.672	2	.000
Step 3	Step	6.862	1	.009
	Block	88.534	3	.000
	Model	88.534	3	.000

Se observă că în fiecare pas al estimării modelului se testează dacă trecerea de la precedent este semnificativă (se respinge ipoteza nulității variabilei sau variabilelor adăugate).

Indicatorii similari coeficientului de determinare din regresia multiplă sunt în tabelul care urmează.

**Model Summary**

Step	-2 Log likelihood	Cox & Snell R Square	Nagelkerke R Square
1	260.573	.174	.294
2	249.589	.198	.335
3	242.727	.213	.360

**Hosmer and Lemeshow Test**

Step	Chi-square	df	Sig.
1	23.800	7	.001
2	11.019	8	.201
3	7.160	8	.519

Testul Hosmer & Lemeshow este explicat pentru fiecare pas prin raportarea celor 10 frecvențe observate/așteptate (statistica și semnificația sunt raportate în tabelul precedent).

**Contingency Table for Hosmer and Lemeshow Test**

		PRIMCARE Board certified in FM, IM, Peds or Cb-Gyn = 0 no		PRIMCARE Board certified in FM, IM, Peds or Ob-Gyn = 1 yes		Total
		Observed	Expected	Observed	Expected	
		Step 3	1	37	36.817	
	2	37	36.535	0	.465	37
	3	37	36.102	0	.898	37
	4	37	35.507	0	1.493	37
	5	35	34.348	2	2.652	37
	6	31	32.848	6	4.152	37
	7	28	30.777	9	6.223	37
	8	26	28.110	11	8.890	37
	9	26	23.459	11	13.541	37
	10	15	14.496	22	22.504	37

Se raportează de asemenea tabelul de clasificare pentru fiecare pas al procedurii.

**Classification Table<sup>a</sup>**

Observed		Predicted		Percentage Correct
		Board certified in FM, IM, Peds or Ob-Gyn		
		no	yes	
Step 1	Board certified in FM, IM, Peds or Ob-Gyn	no	yes	93.9
	Overall Percentage	46	15	24.6
				82.4
Step 2	Board certified in FM, IM, Peds or Ob-Gyn	no	yes	96.8
	Overall Percentage	299	10	23.0
		47	14	84.6
Step 3	Board certified in FM, IM, Peds or Ob-Gyn	no	yes	95.8
	Overall Percentage	296	13	32.8
		41	20	85.4

a. The cutvalue is .500

În tabelul referitor la variabilele din model se raportează:

- coeficienții B
- Exp (B) cu interpretarea, dată în partea teoretică, că reprezintă modificare OR a variabilei dependente la modificarea cu o unitate a variabilei independente, deci  $Exp(B) \approx 1$  pentru variabilele nesemnificative.
- informații asociate testul Wald de semnificație a fiecărui coeficient.



Variables in the Equation

		B	S.E.	Wald	df	Sig.	Exp(B)
Step 1 <sup>a</sup>	MATYR	-.568	.083	46.714	1	.000	.566
	Constant	1129.854	165.478	46.619	1	.000	.
Step 2 <sup>b</sup>	MATYR	-.561	.086	42.232	1	.000	.571
	USMLE2	-.026	.008	10.387	1	.001	.974
	Constant	1120.518	171.865	42.507	1	.000	.
Step 3 <sup>c</sup>	MATYR	-.598	.090	44.344	1	.000	.550
	GENDER(1)	.957	.366	6.855	1	.009	2.605
	USMLE2	-.029	.008	12.294	1	.000	.971
	Constant	1193.796	178.716	44.621	1	.000	.

a. Variable(s) entered on step 1: MATYR.

b. Variable(s) entered on step 2: USMLE2.

c. Variable(s) entered on step 3: GENDER.

În tabelul următor (apare doar pentru anumite metode de selectare a variabilelor) se prezintă informațiile necesare pentru a testa ce s-ar întâmpla dacă o variabilă din model este exclusă. Pentru un model care se construiește ascendent, acestea pot sugera prezența unor variabile care au devenit ne semnificative prin includerea altor variabile.

Model if Term Removed

Variable	Model Log Likelihood	Change in -2 Log Likelihood	df	Sig. of the Change
Step 1 MATYR	-165.630	70.688	1	.000
Step 2 MATYR	-156.540	63.491	1	.000
USMLE2	-130.287	10.984	1	.001
Step 3 MATYR	-155.605	68.483	1	.000
GENDER	-124.795	6.862	1	.009
USMLE2	-127.933	13.138	1	.000

Pentru variabilele care nu sunt în model, se prezintă testele care decid necesitatea prezenței lor. La pasul următor, va fi introdusă în model variabila cu scorul cel mai mare (scor calculat potrivit metodei selectate).

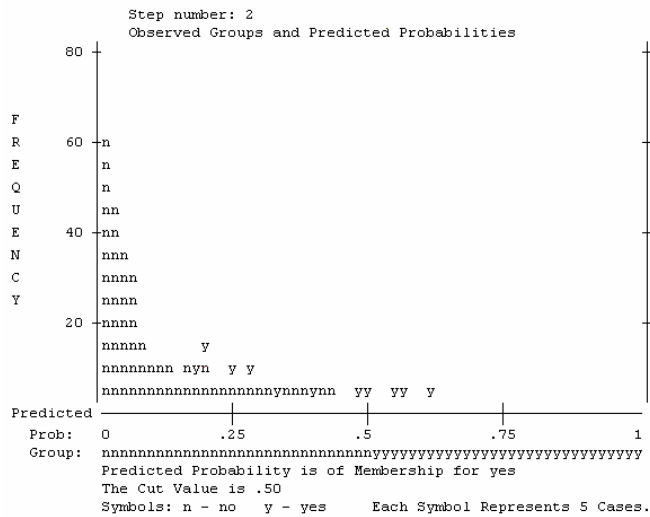
Variables not in the Equation

		Score	df	Sig.
Step 1 Variables	GENDER(1)	4.885	1	.027
	USMLE2	10.948	1	.001
Overall Statistics		17.728	2	.000
Step 2 Variables	GENDER(1)	7.140	1	.008
	Overall Statistics	7.140	1	.008

Diagrama de clasificare (afișată în continuare) este alcătuită:

- Axa X este probabilitatea prognozată (de la 0 la 1) de a fi clasificat în grupul codat "1". Sub axă sunt diferențiate zonele de clasificare prin simbolurile care codifică grupul 1 (Yes) și grupul 2 (No). Se observă pragul de 0.5 care schimbă clasificarea.
- Axa Y este frecvență (număr de cazuri).
- Coloanele care apar în diagramă sunt alcătuite din marcaje (fiecare reprezintă un număr de cazuri pentru simplificare) care reprezintă clasificarea observată a cazurilor.
- Examinarea diagramei constă în analiza faptului dacă marcajele corespund la același semn (Yes/No) situat sub axa X. Prin urmare
  - semnele Y care corespund la valori Y de pe axa OX (și semnele N care corespund la valori N de pe axa OX) reprezintă clasificări prognozate corect de model.
  - celelalte marcaje (semnele Y care corespund la valori N de pe axa OX, precum și semnele N care corespund la valori Y de pe

axa OX) reprezintă cazuri clasate eronat, deci observații pentru care modelul estimat nu funcționează.



### C. Lucrarea practică

1. Un studiu care urmărește de cine depinde gustul brânzeturilor de tip cheddar a prelevat probe și a determinat concentrația unor compuși chimici. Fiecare probă a fost supusă unui proces de degustare și a primit o notă. Unele valori au fost transformate în prealabil (*Acetic* și *H2S* sunt obținute prin logaritmarea valorilor măsurate). Fișierul de date este [www.infoiasi.ro/~val/statistica/CheeseData.txt](http://www.infoiasi.ro/~val/statistica/CheeseData.txt)  
Variabilele sunt
  - i. *Taste*: nota obținută în urma combinării notelor acordate de mai mulți degustători
  - ii. *Acetic*: logaritm natural din concentrația de acid acetic
  - iii. *H2S*: logaritm natural din concentrația de H<sub>2</sub>S.
  - iv. *Lactic*: concentrația de acid lactic
  - Să se modeleze variabila *Taste* cu ajutorul celorlalte trei variabile.
  - Să se analizeze modelul obținut.
2. Date privind un număr de companii au fost selectate din lista Forbes 500 pentru anul 1986 (printr-un sondaj sistematic 1/10 din lista alfabetică a companiilor). Studiul urmărește volumul de vânzări al companiei.  
Fișierul de date este [www.infoiasi.ro/~val/statistica/ForbesData.txt](http://www.infoiasi.ro/~val/statistica/ForbesData.txt)  
Variabilele sunt:
  - i. *Company*: numele companiei
  - ii. *Assets*: bunurile companiei (milioane \$)
  - iii. *Sales*: volumul de vânzări (milioane \$)
  - iv. *Market Value*: valoarea de piață a companiei (milioane \$)
  - v. *Profits*: profitul (milioane \$)
  - vi. *Cash Flow*: volumul tranzacțiilor (milioane \$)
  - vii. *Employees*: numărul de angajați (mii persoane)
  - viii. Sector: domeniul de activitate a companiei.

- Să se modeleze volumul de vânzări în funcție de celelalte variabile. Să se analizeze modelul obținut.
  - Să se determine transformările prealabile necesare pentru unele variabile și să se refacă modelarea.
3. Se va deschide fișierul **Employee Data.sav** din setul de fișiere test oferite de SPSS. Să se decidă dacă faptul că un angajat aparține minorității (*minority* = 1) este reflectat de variabilele *educ*, *prevexp*, *jobcat* și *gender*.  
Pentru aceasta se va estima și se va analiza o regresie logistică în care variabila dependentă este *minority*, restul variabilelor fiind considerate independente.