

Capitolul 6

TRATAREA SIMULTANĂ A DOUĂ VARIABLE CALITATIVE

Dependența în planul datelor a două variabile calitative, se numește **asociere**. Și în acest caz funcționează **postulatul epistemologic**:

“O legătură în plan fenomenologic implică asocierea celor două variabile calitative corespunzătoare în planul datelor”.

Reciproca nu este valabilă și deci, ca și în cazul variabilelor cantitative, există **asocieri aparente sau false asocieri**.

Pentru măsurarea asocierii ori pentru susținerea absenței sale, șirurile bivariate de variante ale celor două variabile calitative se grupează în tabele cu dublă intrare, numite *tabele de contingență*.

Definiție: Se numește **tabel de contingență** rezultatul "ventilării unei populații după variantele a două caracteristici calitative" [3].

Exemplul 6.

Să presupunem că avem următorul șir bivariat de specii, respectiv parcele în care s-au întâlnit speciile respective:

Parcela:	1	1	1	1	1	1	1	1	2	2	2	2	2	2	2	2
Specia:	1	1	1	1	2	2	3	3	1	1	2	2	2	2	3	3

Tabelul de contingență corespunzător va fi format din 2 linii corespunzătoare celor 2 parcele și 3 coloane corespunzătoare celor 3 specii și va conține, în celula de pe linia i și coloana j , numărul de perechi (i, j) din șirul bivariat. În final sunt adăugate coloana și linia „totalurilor marginale” (coloana totalurilor pe linii și linia totalurilor pe coloane) precum și totalul general:

	specia 1	specia 2	specia3	Totaluri pe linii:
<i>parcela 1</i>	4	2	2	8
<i>parcela 2</i>	2	4	2	8
<i>Totaluri pe coloane:</i>	6	6	4	Total general: 16

§ 6.1. Sinteza grafică

Un tabel de contingență poate fi reprezentat prin **areale în dreptunghiuri** (în cadrul unui pătrat [16]).

Alături este reprezentat tabelul de contingență anterior.

		1		
		S_1	S_2	S_3
parcela 1		$\frac{4}{16}$	$\frac{2}{16}$	$\frac{2}{16}$
		S_1	S_2	S_3
parcela 2		$\frac{2}{16}$	$\frac{4}{16}$	$\frac{2}{16}$
		0		
		1		

§ 6.2. Sinteza numerică bivariată

Pentru perechile de variabile cantitative am văzut că dependența statistică apare plasată între independență și dependența funcțională. În mod analog, aici vom plasa asocierea între independență și cea ce se numește [32] *asocierea completă*.

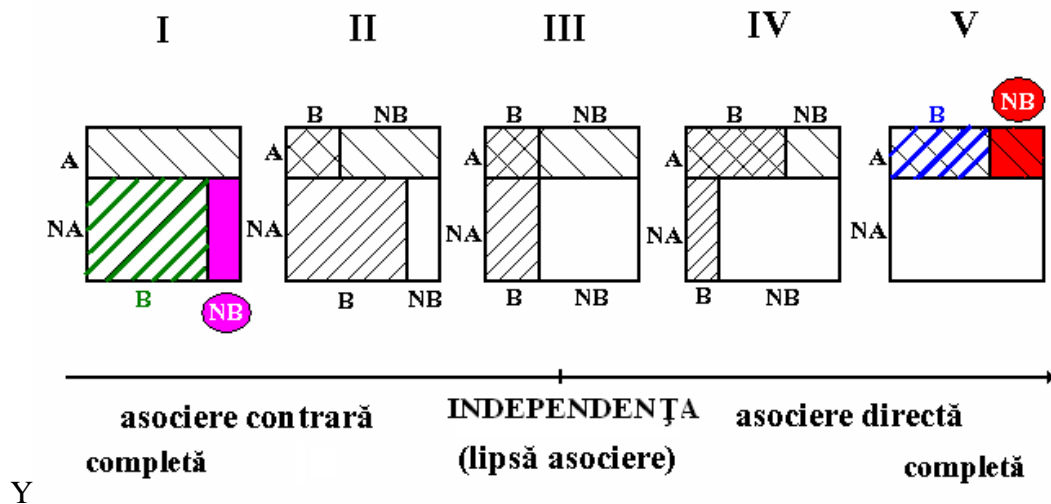
6.2.1. Tabele 2×2

În cazul variabilelor binare, care produc tabele de contingență 2×2 , analogia merge mai departe și anume cu dependențele liniare, unde vorbim despre *corelație directă*, respectiv *inversă*, aici propunem să vorbim despre *asociere directă* respectiv *contrară*.

Să presupunem că dorim să exprimăm dependența între o anumită specie și altitudinea la care este găsită. De exemplu, se știe că floarea de colț crește numai la altitudine mare. Spunem că există o asociere completă directă între altitudine și floarea de colț. Analog, vom spune că există o asociere completă CONTRARĂ²⁴ între altitudine și grâu, deoarece grâul crește numai acolo unde NU exista altitudine (mare). Prin urmare, două variabile calitative binare a (cu variantele A și NA - citește *Non A*) și b (cu variantele B și NB) se numesc **complet asociate în mod direct**, dacă B are loc doar în prezența lui A sau invers. Se vor numi **complet asociate în mod contrar**, dacă B are loc doar în absența lui A sau invers. În exemplele anterioare a este variabila binară "altitudine", cu variantele A = "altitudine" și NA = "Non altitudine", iar b este variabila specie cu variantele B = "floare de colț" și NB = "Non floare de colț", în primul caz, respectiv B = "grâu" și NB = "Non grâu", în al doilea caz.

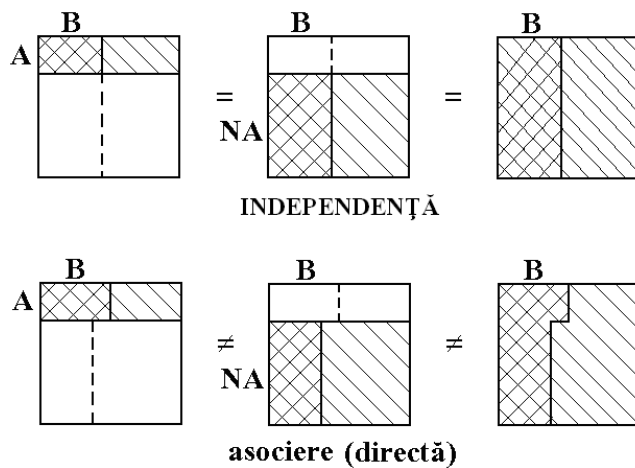
Exemplul cu floarea de colț este reprezentat în desenul V din figura următoare, iar cel cu grâul în desenul I. Desenul IV arată o situație în care B are loc în prezența lui A , dar și în absența sa (adică prezența lui NA). Proporția prezenței lui B în NA este însă mai mică decât proporția prezenței în A (vezi aria hașurată exclusiv ascendent raportată la dreptunghiul orizontal inferior în care este inclusă, respectiv aria dublu hașurată raportată la dreptunghiul superior din care face parte). Este deci cazul unei asocieri directe între A și B dar nu complete. În mod simetric, desenul II exprimă o asociere contrară dar nu completă.

²⁴ În [32] se vorbește despre caracteristici calitative complet asociate, respectiv complet neasociate. Ultima denumire se poate, însă, confunda cu independența. De aceea am propus această nouă terminologie.



În fine, desenul III exprimă **independența** între două va-riabile calitative: proporția lui B din A este egală cu proporția lui B din NA (ceea ce implică și egalitatea cu proporția lui B).

În caz contrar, cele două variabile sunt asociate. (Vezi figura alăturată. O prezentare mai ri-guroasă a independenței se află în volumul următor și în [12].)



Exemplele 6.2.1.

a. Să presupunem că am împărțit în două parcele o anumită suprafață de teren și am numărat pentru fiecare parcelă exemplarele care aparțin, respectiv nu aparțin unei anumite specii S, obținând următorul tabel de contingență la care am adăugat totalurile pe linii, pe coloane (denumite totaluri marginale) și totalul general.

	"specia S"=B	"specie ≠ S"=NB	
"parcela 1"=A	8	12	20
"parcela 2"=NA	22	33	55
	30	45	75

Se observă că prin introducerea parcelelor, proporția speciei S nu se modifică, adică proporția speciei S în parcela 1 ($8 / 20 = 2 / 5$) este egală cu proporția în parcela 2 ($22 / 55 = 2 / 5$) și cu proporția în tot arealul ($30 / 75 = 2 / 5$). Prin urmare, cele două variabile sunt independente. Altfel spus cele două parcele împart același biotop, doar în mod convențional.

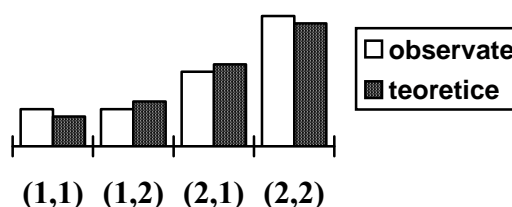
- b. Să repetăm problema anterioară schimbând datele din interiorul tabelului, dar păstrând totalurile marginale pentru comparabilitate.

	"specia S"=B	"specie ≠ S"=NB	
"parcela 1"=A	10	10	20
"parcela 2"=NA	20	35	55
	30	45	75

În acest caz, observăm deosebiri între proporțiile speciei S în cele două parcele: $10 / 20 (= 1 / 2)$ este diferit de $20 / 55 (= 4 / 11)$, proporțiile fiind diferite, evident și de proporția speciei S în tot arealul ($30 / 75 = 2 / 5$). Există, deci, o asociere directă (vezi fig. IV de mai sus) între A și B ($1 / 2$ fiind mai mare decât $4 / 11$). În consecință, judecând în cadrul statisticii descriptive, cele două parcele delimitează două biotopuri ce conțin biocenoze diferite.

- ✓ Deoarece totalurile marginale sunt aceleași, iar proporțiile din primul tabel (a) respectă condiția de independență, acesta poate fi considerat pentru cel de-al doilea (b), tabel teoretic în ipoteza de independență. Modul în care se calculează, în general, pentru un tabel dat, tabelul cu frecvențele teoretice în ipoteza de independență, este explicat la subpunctul următor.

Judecând astfel, putem măsura asocierea prin depărtarea de la independență. Aceasta, deoarece reperul central al problemei asocierii este independența, așa cum se observă și din figurile I-V. Prin urmare, va trebui să apelăm la o măsură a depărtării dintre cele două distribuții, cea de frecvențe observate (o) și cea de frecvențe teoretice în ipoteza de independență (t).



Măsura clasică de comparare a două distribuții a fost prezentată în 3.7.7. și este:

$$\chi^2 = \frac{\sum (o-t)^2}{t} \text{ (formula teoretică)} = \sum \frac{o^2}{t} - N \text{ (formula de calcul),}$$

unde N este totalul general (volumul seriei bivariate).

Exemplul 6.2.1'.

În cazul exemplului 5.2.1. vom avea:

$$\chi^2 = \sum o^2 / t - N = 10^2 / 8 + 10^2 / 12 + 20^2 / 22 + 35^2 / 33 - 75 \cong 12,5 + 8,3 + 18,2 + 37,1 - 75 = 76,1 - 75 = 1,1.$$

Problema rămasă nerezolvată este cum calculăm tabelul în ipoteza de independență. Vom prezenta acest lucru în continuare, pe cazul general al unui tabel de contingență cu p linii și q coloane.

6.2.2. χ^2 pentru independență versus asociere, în tablele de contingență $p \times q$

Forma generală a unei tablele de contingență este:

$y:$ $x =$	Variantele lui x	Sume pe linii
	$x_1 \ x_2 \ \dots \ x_j \ \dots \ x_q$	$l(i):$
y_1	.	$l(i) = \sum_{j=1}^q o(i, j)$
y_2	.	
.	.	
.	.	
y_i	. . . $o(i, j)$. . .	
.	.	
y_p	.	
Sume pe coloane $c(j) :$	$c(j) = \sum_{i=1}^p o(i, j)$	$N = \sum_{i=1}^p \sum_{j=1}^q o(i, j) = \sum_{i=1}^p l(i) = \sum_{j=1}^q c(j)$

Prin $l(i)$ s-a notat totalul marginal al elementelor de pe linia i , iar prin $c(j)$ totalul marginal al elementelor de pe coloana j . N este totalul general.

Tabelul de frecvențe teoretice în ipoteza de independență va fi alcătuit din elementele:

$$t(i, j) = l(i) \cdot c(j) / N.$$

Formula de calcul rezultă din caracterizarea independenței de mai sus²⁵.

Cititorul va face două exerciții utile calculând frecvențele teoretice²⁶ pentru tabelul din exemplul 5.2.1.b. (obținând tabelul 5.2.1.a) și pentru tabelul 2×3 din exemplul 5 (pentru care va obține tabelul din desenul alăturat).

		1		
		S_1	S_2	S_3
parc ela	1	$\frac{3}{16}$	$\frac{3}{16}$	$\frac{2}{16}$
	2	$\frac{3}{16}$	$\frac{3}{16}$	$\frac{2}{16}$
		0		

În cazul general formula teoretică, respectiv cea de calcul pentru χ^2 vor fi:

$$\chi^2 = \sum_{i=1}^p \sum_{j=1}^q \frac{(o(i, j) - t(i, j))^2}{t(i, j)} = \sum_{i=1}^p \sum_{j=1}^q \frac{o(i, j)^2}{t(i, j)} - N.$$

Reținem, deci, că χ^2 măsoară depărtarea de la independență, respectiv intensitatea absolută a asocierii.

²⁵ în care înlocuim A cu varianta i a lui y și B cu varianta j a lui x . Atunci independență înseamnă că proporția variantei j a lui x din întreaga variantă i a lui y ($t(i, j) / l(i)$) este egală cu proporția întregii variante j ($c(j) / N$), adică: $t(i, j) / l(i) = c(j) / N$, de unde rezultă egalitatea de demonstrat.

²⁶ Frecvențele teoretice pot fi și numere cu zecimale, nu numai numere întregi, ca în exemplele de aici.

1° Coeficientul de contingență al lui Ciuprov

Pentru că χ^2 prezintă dezavantajul creșterii o dată cu numărul de linii p , de coloane q și numărul de unități statistice N , fără a exprima în același timp o creștere a asocierii, s-a introdus o măsură relativă a asocierii și anume *coeficientul de contingență al lui Ciuprov*, notat T , coeficient al cărui pătrat este dat de formula:

$$T^2 = \frac{\chi^2}{N \cdot \sqrt{(p-1) \cdot (q-1)}}.$$

T se obține extrăgându-se radicalul de ordinul doi din T^2 . Se consideră numai rezultatul pozitiv²⁷, deoarece în cazul acestui coeficient nu mai există avantajul oferit în cazul variabilelor cantitative, de coeficientul de corelație liniară al lui Pearson, care prin semnul său indică sensul legăturii (corelației). Este natural să fie așa, deoarece în cazul variabilelor calitative nu mai dispunem de o ordine pentru fiecare variabilă și, deci, nu putem vorbi de un sens al asocierii decât în cazul variabilelor binare, așa cum am arătat mai sus în 5.2.1.

Exemplul 5.2.1''.

Continuând exemplul numeric de mai sus (5.2.1'),

$$T^2 = \frac{\chi^2}{N \cdot \sqrt{(p-1) \cdot (q-1)}} = \frac{1,1}{75 \cdot \sqrt{(2-1) \cdot (2-1)}} = \frac{1,1}{75} \cong 0,01,$$

de unde rezultă că $T = \sqrt{0,01} = 0,1$. Deoarece T este diferit de zero, există o asociere între cele două variabile, de intensitate slabă însă, T fiind apropiat de 0.

²⁷ T variază între 0 și 1 dacă $p = q$ [32].