



UNIVERSITÀ
DEGLI STUDI DELLA
TUSCIA



SAPIENZA
UNIVERSITÀ DI ROMA

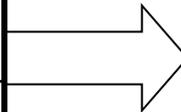
Statistics for business and decision making course

Dr. Ilaria Benedetti

06. Association among variables

Two way table

zone	on line selling
center	yes
outskirt	yes
Semi central	no
outskirt	no
center	no
center	no
outskirt	no
Semi central	no
centre	yes



		on line selling		Tot
		yes	no	
Zone	Center	2	2	4
	Semi central	0	2	2
	outsk	1	2	3
Tot		3	6	9

N=9 selling points

Two way table

		on line selling		Tot
		yes	no	
Zone	Center	2	2	4
	Semi central	0	2	2
	outsk	1	2	3
Tot		3	6	9

- Identify the number of selling point located in the center
- Identify the number of selling point who sell on line
- Considering the selling point who sell on-line, identify the number of selling point located in the center
- Considering the selling point in the outskirts, identify the number of selling point who sell on line

Two way table

Conditional distribution of X for $Y=y_j$

marginal distributions of variable X

Y

		y_1	...	y_j	...	y_k	Totale
X	X_1	n_{11}	...	n_{1j}	...	n_{1k}	$n_{1.}$

	X_i	n_{i1}	...	n_{ij}	...	n_{ik}	$n_{i.}$

	X_H	n_{H1}	...	n_{Hj}	...	n_{HK}	$n_{H.}$
Totale	$n_{.1}$...	$n_{.j}$...	$n_{.k}$	n	

Conditional distribution of Y for $X=x_i$

marginal distributions of variable Y

Two way table

$$\rightarrow n_{i.} = \sum_{j=1}^K n_{ij} \quad \text{per } i=1,2,\dots,H$$

$$\rightarrow n_{.j} = \sum_{i=1}^H n_{ij} \quad \text{per } j=1,2,\dots,K$$

$$\rightarrow n = \sum_{i=1}^H \sum_{j=1}^K n_{ij} = \sum_{i=1}^H n_{i.} = \sum_{j=1}^K n_{.j}$$

The distribution of a variable that is restricted to cases satisfying a condition is called a conditional distribution

Dependence and interdependence

Dependence: studies how the modalities of one character depend on those of another character according to a unidirectional link

Interdependence: It is assumed that the two characters have the same role and that the link is bidirectional

Two way table

Freq. Observed: n_{ij}

Freq. Expected: n_{ij}' in case of independence

$$n_{ij}' = \frac{n_{i.} \times n_{.j}}{n}$$

		Y					Tot
		y ₁	...	y _j	...	y _K	
X	X ₁	n ₁₁		n _{1j}		n _{1k}	n _{1.}
	...						
	X _i	n _{i1}		n _{ij}		n _{ik}	n _{i.}
	...						
	X _H	n _{H1}		n _{Hj}		n _{HK}	n _{H.}
Tot		n _{.1}		n _{.j}		n _{.K}	n

 Freq. obs

 Freq. Used to find the exp. freq.

$$n'_{ij} = \frac{n_{i.} \times n_{.j}}{n}$$

Obs. Freq.

		On-line sell		Tot
		yes	no	
Zone	Center	2	2	4
	Semi - Central	0	2	2
	outskirt	1	2	3
Tot		3	6	9

Exp. Freq

		on line selling		Tot
		yes	no	
Zone	Centre	$\frac{4 \times 3}{9}$	$\frac{4 \times 6}{9}$	4
	Semi central	$\frac{2 \times 3}{9}$	$\frac{2 \times 6}{9}$	2
	outsk.	$\frac{3 \times 3}{9}$	$\frac{3 \times 6}{9}$	3
Tot		3	6	9

In case of statistical independence

Observed. Freq. vs expected freq

f. obs		on line sell		Tot
		yes	no	
zone	Centre	2	2	4
	Semi centr	0	2	2
	outsk	1	2	3
Tot		3	6	9

	Teoriche	on line selling		Tot
		si	no	
zone	Center	1,33	2,67	4
	Semi centr	0,67	1,33	2
	outsk.	1	2	3
Tot		3	6	9

➡ **What is it the degree of association among these two var?**

compare the observed data (original table) with the «artificial» table (expected frequencies in case of no association).

Chi-squared compares the observed contingency table with an artificial – expected - table with the same marginal totals.

If the tables are similar, then there's not association.

The greater the difference between the tables becomes, the greater the association.

Summary: steps to compute chi-square metrics

1. From the table of observed freq. distribution, prepare the table of **expected frequency distribution** (with the same marginal distributions);
2. Subtract the cells of the expected table from those of the observed table and then square the differences.
3. Divide each squared differences in each cell by expected freq. Distribution.
4. Sum the normalized values obtained in point 3. over all of the cells.

Chi-square index

$$\chi^2 = \sum_{i=1}^H \sum_{j=1}^K \frac{c_{ij}^2}{n'_{ij}} \longrightarrow c_{ij} = n_{ij} - n'_{ij}$$

$\chi^2 = 0 \Rightarrow$ Statistical Independence

$\chi^2 > 0 \Rightarrow$ Statistical interdependence

• Calculating The Chi-Square Statistic

The tables below show the observed and expected counts for the wine and music experiment. Calculate the chi-square statistic.

Observed Counts					Expected Counts				
	Music					Music			
Wine	None	French	Italian	Total	Wine	None	French	Italian	Total
French	30	39	30	99	French	34.22	30.56	34.22	99
Italian	11	1	19	31	Italian	10.72	9.57	10.72	31
Other	43	35	35	113	Other	39.06	34.88	39.06	113
Total	84	75	84	243	Total	84	75	84	243

For the French wine with no music, the observed count is 30 bottles and the expected count is 34.22. The contribution to the χ^2 statistic for this cell is

$$\frac{(\text{Observed} - \text{Expected})^2}{\text{Expected}} = \frac{(30 - 34.22)^2}{34.22} = 0.52$$

The χ^2 statistic is the sum of nine such terms:

$$\begin{aligned} \chi^2 &= \sum \frac{(\text{Observed} - \text{Expected})^2}{\text{Expected}} = \frac{(30 - 34.22)^2}{34.22} + \frac{(39 - 30.56)^2}{30.56} + \dots + \frac{(35 - 39.06)^2}{39.06} \\ &= 0.52 + 2.33 + \dots + 0.42 = 18.28 \end{aligned}$$

V Cramer index

Relative index of association

$$V = \sqrt{\frac{\chi^2 / n}{\min[(H - 1), (K - 1)]}} \quad 0 \leq V \leq 1$$

$V=0$ \Rightarrow statistical independence – no assoc

$V=1$ \Rightarrow perfect association

If $V=0.25$, we will say that the association is weak.

If $V=0.75$, we will say that the association is strong.

In between, we will say there is moderate association

χ^2 and V

$$\begin{aligned} \chi^2 &= \frac{(2 - 1,33)^2}{1,33} + \frac{(2 - 2,67)^2}{2,67} + \\ &+ \frac{(0 - 0,67)^2}{0,67} + \frac{(2 - 1,33)^2}{1,33} + \\ &+ \frac{(1 - 1)^2}{1} + \frac{(2 - 2)^2}{2} \end{aligned} \quad \Rightarrow \quad \begin{aligned} \chi^2 &= 0,33 + 0,17 + \\ &+ 0,33 + 0,67 = 1,5 \end{aligned}$$

H=3, K=2 min among H-1
and K-1 in equal to 1

$$\Rightarrow \quad V = \sqrt{\frac{1,50}{9}} = 0,41$$

A character Y perfectly depends on X when only one mode of Y is associated with each mode of X, that is, when in a double table for each i there is only one j for which $n_{ij} \neq 0$

		on line		Tot
		yes	no	
	Centre	4	0	4
	Semi central	0	2	2
	Outsk	0	3	3
Tot		5	5	9

In this table the online sale (Y) depends perfectly on the location (X) Attention X does not depend on Y

Between two characters X and Y there is perfect interdependence if each modality of one of the two characters corresponds to one and only one modality of the other character and vice versa.

		Y			Tot
		Y ₁	Y ₂	Y ₃	
X	X ₁	4	0	0	4
	X ₂	0	2	0	2
	X ₃	0	0	3	3
Tot		4	2	3	9

ONLY for
Squared table!

Checklist:

Chi-squared and Cramer's V. Chi-squared and Cramer's V measure association between two categorical variables that define a contingency table.

Before you use these, verify that your data meet this prerequisite:

✓ Categorical variables. If a variable is numerical, there are better ways to measure association.