

# Studio dell'associazione

## tra due variabili qualitative

Task Description

Due Date

Status

	Task Description	Due Date	Status
<input type="checkbox"/>	Studio dell'associazione tra due variabili qualitative: l'indipendenza assoluta; confronto tra distribuzioni condizionate e		
<input type="checkbox"/>	distribuzione marginale; condizione necessaria e sufficiente per la sussistenza di indipendenza assoluta; frequenze teoriche		
<input type="checkbox"/>	sotto l'ipotesi di indipendenza assoluta; l'indice $\chi^2$ come distanza normalizzata tra la tabella delle frequenze		
<input type="checkbox"/>	osservate e la tabella delle frequenze teoriche; normalizzazione dell'indice $\chi^2$ : l'indice $\Phi^2$ e l'indice $v$ di Grame		
<input type="checkbox"/>			

Prüfmemoria  
(tabelle e doppie entrate)

		Y						
		$y_1$	$y_2$	...	$y_j$	...	$y_c$	
X	$x_1$	$m_{11}$	$m_{12}$	...	$m_{1j}$	...	$m_{1c}$	$M_{1\bullet}$
	$x_2$	$m_{21}$	$m_{22}$	...	$m_{2j}$	...	$m_{2c}$	$M_{2\bullet}$
	$\vdots$	$\vdots$	$\vdots$		$\vdots$		$\vdots$	$\vdots$
	$x_i$	$m_{i1}$	$m_{i2}$	...	$m_{ij}$	...	$m_{ic}$	$M_{i\bullet}$
	$\vdots$	$\vdots$	$\vdots$		$\vdots$		$\vdots$	$\vdots$
$x_R$	$m_{R1}$	$m_{R2}$	...	$m_{Rj}$	...	$m_{Rc}$	$M_{R\bullet}$	
		$M_{\bullet 1}$	$M_{\bullet 2}$	...	$M_{\bullet j}$	...	$M_{\bullet c}$	$M_{\bullet\bullet}$    $n$

freq. congiunta  
 $X = x_i \cap Y = y_j$

freq. marginale  
 $Y = y_j$

freq. marginale  
 $X = x_i$

Sia X che Y sono qualitative  
 (si può lavorare allo stesso modo anche in caso di variabili miste o di due variabili numeriche)

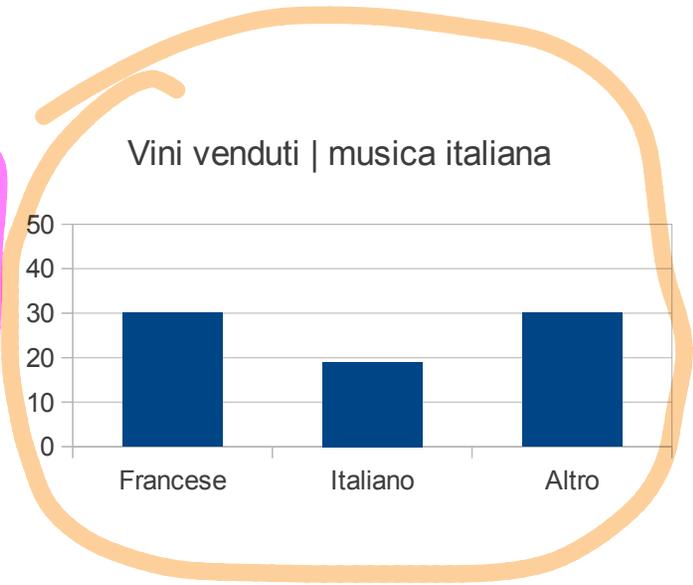
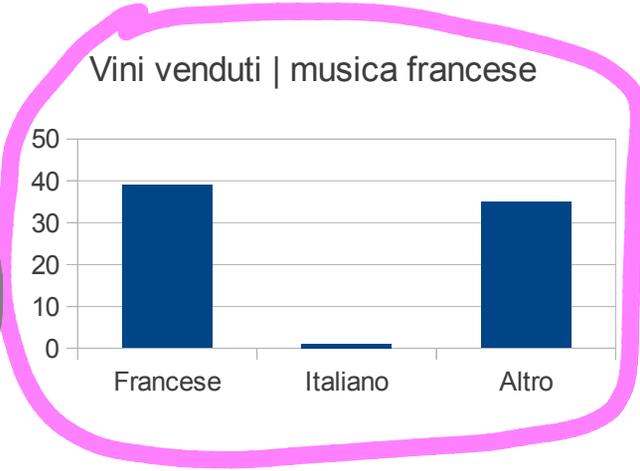
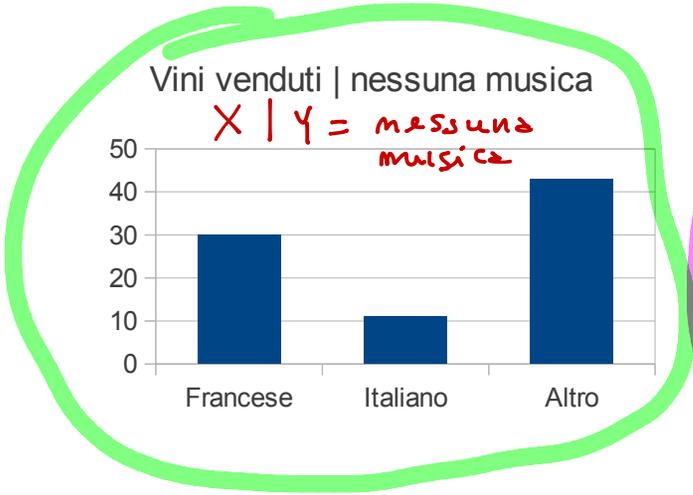
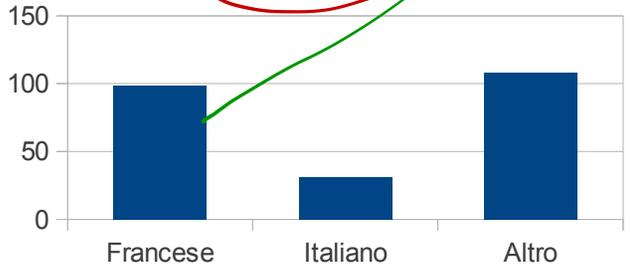
$M_{ij}$

		Musica sottofondo			TOTALE
		Nessuna	Francese	Italia	
Vino	Francese	30	39	30	99
	Italiano	11	1	19	31
	Altro	43	35	30	108
TOTALE		84	75	79	238

$\rightarrow N$

$X$   
 paese provenienza

$Y$   
 Vini Venduti



		Musica sottofondo			TOTALE
		Nessuna	Francese	Italia	
Vino	Francese	30	39	30	99
	Italiano	11	1	19	31
	Altro	43	35	30	108
TOTALE		84	75	79	238

profili colonne →

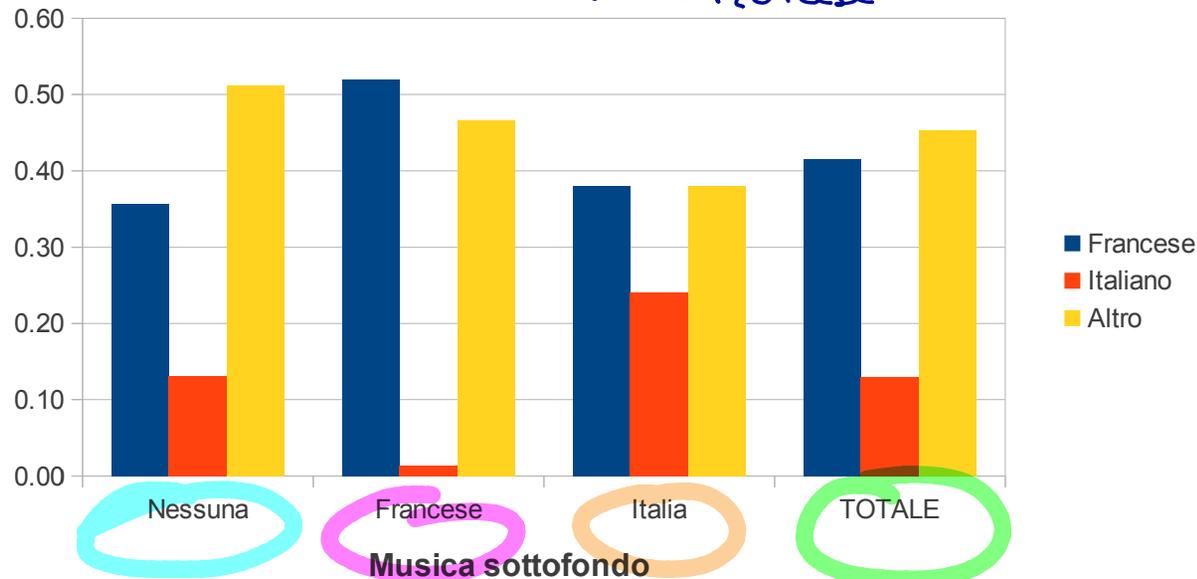
		Musica sottofondo			TOTALE
		Nessuna	Francese	Italia	
Vino	Francese	$\frac{30}{84} = 0.36$	$\frac{11}{84} = 0.13$	$\frac{43}{84} = 0.51$	0.42
	Italiano	0.36	0.01	0.38	0.13
	Altro	0.51	0.47	0.38	0.45
TOTALE		1	1	1	1

media ponderata di 0,36, 0,52 e 0,38  

$$\frac{0,36 \times 84 + 0,52 \times 75 + 0,38 \times 79}{238}$$

il 42% dei clienti  
 di acquista una bottiglia  
 di vino francese

il profilo marginale  
 può essere interpretato  
 come profilo medio



Se tra le due variabili X e Y non vi fosse associazione (connessione o legame) dovei aspettarmi distribuzioni condizionate uguali tra loro (uguali anche alla distribuzione marginale)

↳ condizione di indipendenza assoluta

Frequenze osservate

$M_{ij}$		Musica			
		N	F	I	
VINO	F	30	39	30	99
	I	11	1	19	31
	A	43	35	30	108
		84	75	79	238

$\frac{M_{ij}}{M_{.j}}$

		Musica			
		N	F	I	
VINO	F	0,36	0,52	0,38	0,42
	I	0,13	0,01	0,24	0,13
	A	0,51	0,47	0,38	0,45
		1	1	1	1

## CONDIZIONE DI INDIPENDENZA ASSOLUTA

Se non vi fosse alcuna relazione scelta di acquisto e scenario di acquisto, ciascun profilo colonna dovrebbe essere uguale al profilo marginale

Frequenze osservate

$M_{ij}$		Musica			
		N	F	I	
VINO	F	30	39	30	99
	I	11	1	19	31
	A	43	35	30	108
		84	75	79	238

$\frac{M_{ij}}{M_{0j}}$

		Musica			
		N	F	I	
VINO	F	0,36	0,52	0,38	0,42
	I	0,13	0,01	0,24	0,13
	A	0,51	0,47	0,38	0,45
		1	1	1	1

frequenze teoriche o attese  
sotto l'ipotesi di indipendenza

$$l_{ij} = \hat{M}_{ij} = M_{ij}^* = M_{ij}^!$$

		Musica			
		N	F	I	
VINO	F	0,42 × 84	0,42 × 75	0,42 × 79	99
	I	0,13 × 84	0,13 × 75	0,13 × 79	31
	A	0,45 × 84	0,45 × 75	0,45 × 79	108
		84	75	79	238

questi devono essere uguali alla tabella  $M_{ij}$

frequenze teoriche o attese  
 sotto l'ipotesi di indipendenza

$$e_{ij} = \hat{M}_{ij} = M_{ij}^* = M_{ij}^1$$

		Musica			
		N	F	I	
VINO	F	0,42 × 84	0,42 × 75	0,42 × 79	99
	I	0,13 × 84	0,13 × 75	0,13 × 79	31
	A	0,45 × 84	0,45 × 75	0,45 × 79	108
		84	75	79	238

$$0,42 = \frac{99}{238} = \frac{m_{10}}{n}$$

$$0,13 = \frac{31}{238} = \frac{m_{20}}{n}$$

$$0,45 = \frac{108}{238} = \frac{m_{30}}{n}$$

$$\rightarrow 0,42 \times 84 = \frac{m_{10}}{n} \cdot m_{01} = \frac{m_{10} \cdot m_{01}}{n} = \hat{M}_{11}$$

$$\rightarrow 0,13 \times 79 = \frac{m_{20}}{n} \cdot m_{03} = \frac{m_{20} \cdot m_{03}}{n} = \hat{M}_{23}$$

$$\rightarrow 0,45 \times 75 = \frac{m_{30}}{n} \cdot m_{02} = \frac{m_{30} \cdot m_{02}}{n} = \hat{M}_{32}$$

	$y_1$	$y_2$	...	$y_j$	...	$y_c$
$x_1$						
$x_2$						
⋮						
$x_i$				$m_{ij}$		$m_{i\cdot}$
⋮						
$x_r$						
						$m_{\cdot j}$
						$n$

	$y_1$	$y_2$	...	$y_j$	...	$y_c$
$x_1$						
$x_2$						
⋮						
$x_i$				$\hat{m}_{ij} = \frac{m_{i\cdot} \cdot m_{\cdot j}}{n}$		$m_{i\cdot}$
⋮						
$x_r$						
						$m_{\cdot j}$
						$n$

### CONDIZIONE DI INDIPENDENZA ASSOLUTA

↳ i profili colonna ( $y_j$ ) devono essere tutti uguali

$$\frac{m_{ij}}{m_{\cdot j}} = \frac{m_{i\cdot}}{n} \quad \forall i, j$$

$$\frac{m_{ij}}{m_{i\cdot}} = \frac{m_{\cdot j}}{n} \quad \forall i, j$$

CONDIZIONE INDIPENDENZA  
SUI PROFILI COLONNA

$$\frac{m_{ij}}{m_{\cdot j}} = \frac{m_{i\cdot}}{n} \quad \forall i, j$$

CONDIZIONE INDIPENDENZA  
SUI PROFILI RIGA

$$\frac{m_{ij}}{m_{i\cdot}} = \frac{m_{\cdot j}}{n} \quad \forall i, j$$

$$\hat{m}_{ij} = \frac{m_{i\cdot} \cdot m_{\cdot j}}{n}$$

CONDIZIONE  
NECESSARIA E  
SUFFICIENTE

frequenze attese o teoriche sotto l'ipotesi di indipendenza  
(per evitare confusione le indichiamo con  $\hat{m}_{ij}$ )


$m_{ij}$


$\hat{m}_{ij}$

l'indice di associazione  $\bar{e}$  una distanza tra le tabelle delle  $m_{ij}$  e la tabelle delle  $\hat{m}_{ij}$

$m_{ij}$

		Musica sottofondo			TOTALE
		Nessuna	Francese	Italia	
Vino	Francese	30	39	30	99
	Italiano	11	1	19	31
	Altro	43	35	30	108
TOTALE		84	75	79	238

$\frac{m_{i \cdot} \cdot m_{\cdot j}}{n} = \hat{m}_{ij}$

$34.94 = \frac{99 \times 84}{238} = \frac{m_{i \cdot} \cdot n_{\cdot 1}}{N}$

		Musica sottofondo			TOTALE
		Nessuna	Francese	Italia	
Vino	Francese	34.94	31.20	32.86	99
	Italiano	10.94	9.77	10.29	31
	Altro	38.12	34.03	35.85	108
TOTALE		84	75	79	238

$m_{ij} - \hat{m}_{ij}$

		Musica sottofondo			TOTALE
		Nessuna	Francese	Italia	
Vino	Francese	-4.94	7.80	-2.86	0
	Italiano	0.06	-8.77	8.71	0
	Altro	4.88	0.97	-5.85	0
TOTALE		0	0	1.0658E-014	0

$\frac{(m_{ij} - \hat{m}_{ij})^2}{\hat{m}_{ij}}$

		Musica sottofondo			TOTALE
		Nessuna	Francese	Italia	
Vino	Francese	0.70	1.95	0.25	
	Italiano	0.00	7.87	7.37	
	Altro	0.63	0.03	0.95	
TOTALE					19.75

$\sum_{i=1}^n \sum_{j=1}^c \frac{(m_{ij} - \hat{m}_{ij})^2}{\hat{m}_{ij}}$

somma di tutte le distanze (normalizzate)

indice del chi-quadro ( $\chi^2$ )

la divisione per  $\hat{m}_{ij}$  ha l'effetto di normalizzare tutte le celle

trattare tutte le celle allo stesso modo: senza questa operazione le celle con  $m_{i \cdot}$  e/o  $m_{\cdot j}$  avrebbero un peso maggiore

queste sono le due celle che contribuiscono di più a spiegare l'associazione tra X e Y

contingenze

## Ricapitoliamo ---

Lo studio dell'indipendenza si basa sul confronto tra i profili

↓  
nell'es. sul confronto tra i profili colonna e il profilo marginale

$$\left[ \begin{array}{l} \text{distribuz. condizionate} \\ X | Y = y_j \\ j = 1, \dots, c \end{array} \right] \quad \left[ \times \right]$$

lo studio è di tipo simmetrico

↓  
confrontare i profili colonna equivale a confrontare i profili riga  
da un punto di vista numerico

↓  
non è così rispetto all'interpretazione dei risultati

$y_1$  $y_2$  $y_j$  $y_c$  $x_1$ 

$$\frac{m_{11}}{m_{\cdot 1}} = \frac{m_{12}}{m_{\cdot 2}} = \dots = \frac{m_{1j}}{m_{\cdot j}} = \dots = \frac{m_{1c}}{m_{\cdot c}} = \frac{m_{1\cdot}}{N}$$

 $x_i$ 

$$\frac{m_{i1}}{m_{\cdot 1}} = \frac{m_{i2}}{m_{\cdot 2}} = \dots = \frac{m_{ij}}{m_{\cdot j}} = \dots = \frac{m_{ic}}{m_{\cdot c}} = \frac{m_{i\cdot}}{N}$$

 $x_2$ 

$$\frac{m_{21}}{m_{\cdot 1}} = \frac{m_{22}}{m_{\cdot 2}} = \dots = \frac{m_{2j}}{m_{\cdot j}} = \dots = \frac{m_{2c}}{m_{\cdot c}} = \frac{m_{2\cdot}}{N}$$

 $(1)$  $(1)$  $(1)$  $(1)$  $(1)$

formulazione  
compatte  
condizione di  
indipendenza

$$\frac{n_{ij}}{n_{\cdot j}} = \frac{n_{i\cdot}}{N}$$

$\forall i=1, \dots, r$   
 $\forall j=1, \dots, c$

freq. che osservi  
osservare in ciascuna  
cella  $\uparrow$  se  $r_i$  fosse  
indipendenza

moltiplicando  
ambo i membri  
per  $n_{\cdot j}$

$$n_{ij} = \frac{n_{i\cdot} \cdot n_{\cdot j}}{N}$$

formulaz. equivalente

per distinguere queste freq. teoriche o attese sotto  
l'ipotesi di indipendenza di solito le si indica  
con  $\hat{n}_{ij}$  oppure  $m_{ij}^*$  oppure  $e_{ij}$

	Y
X	$m_{ij}$ freq. osservate

	Y
X	$\hat{m}_{ij}$ freq. teoriche o attese sotto l'ipotesi di indipendenza

distanza tra le  $m_{ij}$  e le  $\hat{m}_{ij}$  come  
 misura dell'indipendenza



se vi fosse indipendenza le  $\hat{m}_{ij}$  dovrebbero  
 essere uguali alle  $m_{ij}$  (le due tabelle dovrebbero  
 essere identiche)

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(m_{ij} - \hat{m}_{ij})^2}{\hat{m}_{ij}}$$

$0$   
 $\max$

se solo se  $m_{ij} = \hat{m}_{ij}$   
 $\forall i, j$

dipende dalla struttura  
 delle tabelle  $(r, c)$   
 ma anche dal n. di osservati  $(N)$

Si definisce:  $\Phi^2 = \frac{\chi^2}{N} \rightarrow$  elimina la dipendenza da  $N$   
 indice di contingenza  
 media

$$v = \sqrt{\frac{\Phi^2}{\min(r, c) - 1}}$$

(Cramer)

$\rightarrow$  alcuni autori estraggono la radice

elimina l'effetto anche della struttura

$\rightarrow$  è un indice normalizzato tra 0 e 1

$x^e$   $\left\{ \begin{array}{l} 0 \text{ minimo} \\ N \times [\min(r, c) - 1] \text{ massimo} \end{array} \right.$

$$\frac{x^2 - \min}{\max - \min} = \frac{x^2}{N[\min(r, c) - 1]} = \frac{\Phi^e}{\min(r, c) - 1}$$

$\downarrow$   
 $\Phi^e$

quando si verifica  $x^2 = \max$  (?)

Vino	Musica			
	Nessuna	Francesca	Italiana	
Francesca	0	99	0	99
Italiano	0	0	31	31
Altro	108	0	0	108
	108	99	31	238

$$\chi^2 = 238 \times [\min(3,3) - 1]$$

$$\Phi^2 = \min(3,3) - 1$$

$$\sqrt{\phantom{x}} = 1$$

nel caso di associazione  
1 ad 1 righe - colonne



massima interdipendenza



il max assoluto è osservabile solo nel  
caso di tabelle quadrate ( $r=c$ )

$$\frac{m_{ij}}{m_{\cdot j}} = \frac{m_{i\cdot}}{n}$$

ragion. profili  
colonna

$$\frac{m_{ij}}{m_{i\cdot}} = \frac{m_{\cdot j}}{n}$$

ragionam. profili  
riga

$$\hat{m}_{ij} = \frac{m_{i\cdot} m_{\cdot j}}{n}$$

indice  $\chi^2$  è simmetrico



studio associazione è  
simmetrico



X e Y giocano lo stesso ruolo

INTERDIPENDENZA

dependenza reciproca

Si arriva alla stessa  
condizione di indipendenza  
sia che si parte dai profili  
riga che dai profili colonna

# Schema riepilogativo di calcolo dell'indice $\chi^2$

① tabelle frequenze osservate input

$$M_{ij} \quad i = 1, \dots, r$$
$$j = 1, \dots, c$$

② tabelle frequenze teoriche o attese ( sotto le ipotesi di indipendenza assoluta )

$$\hat{M}_{ij} = \frac{M_{i \cdot} \cdot M_{\cdot j}}{n}$$

③ distanza (normalizzata) tra le  $M_{ij}$  e  $\hat{M}_{ij}$

$$\frac{(M_{ij} - \hat{M}_{ij})^2}{\hat{M}_{ij}}$$

④ sintetizzo la tabella delle distanze con il seguente indice

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(M_{ij} - \hat{M}_{ij})^2}{\hat{M}_{ij}}$$

# Formula alternativa (ridotta) del $\chi^2$

$$\chi^2 = \sum_i \sum_j \frac{(m_{ij} - \hat{m}_{ij})^2}{\hat{m}_{ij}} = \sum_i \sum_j \frac{\left(m_{ij} - \frac{m_{i\cdot} m_{\cdot j}}{n}\right)^2}{\frac{m_{i\cdot} m_{\cdot j}}{n}} =$$

$$= \sum_i \sum_j \frac{m_{ij}^2 + \left(\frac{m_{i\cdot} m_{\cdot j}}{n}\right)^2 - 2 m_{ij} \frac{m_{i\cdot} m_{\cdot j}}{n}}{\frac{m_{i\cdot} m_{\cdot j}}{n}} =$$

$$= \sum_i \sum_j \frac{m_{ij}^2 + \left(\frac{m_{i\cdot} m_{\cdot j}}{n}\right) \left(\frac{m_{ij} m_{\cdot j}}{n}\right) - 2 m_{ij} \frac{m_{i\cdot} m_{\cdot j}}{n}}{\frac{m_{i\cdot} m_{\cdot j}}{n}} =$$

∕

$$= \sum_i \sum_j \left\{ \frac{m_{ij}^2}{m_{i\cdot} m_{\cdot j}} + \frac{\left( \frac{m_{i\cdot} m_{\cdot j}}{m} \right) \left( \frac{m_{ij} m_{\cdot j}}{m} \right)}{\frac{m_{i\cdot} m_{\cdot j}}{m}} - \frac{2 m_{ij} \frac{m_{i\cdot} m_{\cdot j}}{m}}{\frac{m_{i\cdot} m_{\cdot j}}{m}} \right\} =$$

$$= m \sum_i \sum_j \frac{m_{ij}^2}{m_{i\cdot} m_{\cdot j}} + \sum_i \sum_j \frac{m_{i\cdot} m_{\cdot j}}{m} - 2 \sum_i \sum_j m_{ij} =$$

$$= m \sum_i \sum_j \frac{m_{ij}^2}{m_{i\cdot} m_{\cdot j}} + \frac{\overset{(\text{= } m)}{\sum_i m_{i\cdot}} \overset{(\text{= } m)}{\sum_j m_{\cdot j}}}{m} - 2m =$$

$$= m \sum_i \sum_j \frac{m_{ij}^2}{m_{i\cdot} m_{\cdot j}} - m = m \left[ \sum_i \sum_j \frac{m_{ij}^2}{m_{i\cdot} m_{\cdot j}} - 1 \right]$$