

Introduction to data science & artificial intelligence (INF7100)

Arthur Charpentier

#261 Bivariate Statistics

été 2020

Pearson's Chi-Square Test

Consider two factor variables,

$$x_1 \in \{a_1, \dots, a_I\} \text{ and } x_2 \in \{b_1, \dots, b_J\}$$

convert the dataframe into a contingency table

$$n_{i,j} = \sum_k \mathbf{1}(x_{1,k} = a_i, x_{2,k} = b_j)$$

```
1 > loc = "http://freakonometrics.free.fr/titanic.RData"
2 > download.file(loc, "titanic.RData")
3 > load("titanic.RData")
4 > base = base[,1:7]
5 > table(base$Survived, base$Pclass)
```

	1	2	3
0	64	90	270
1	120	83	85

Pearson's Chi-Square Test

Test : $H_0 : X_1 \perp\!\!\!\perp X_2$, i.e. (cf [definition](#)), $\forall i, j$

$$\mathbb{P}[X_1 = a_i, X_2 = b_j] = \mathbb{P}[X_1 = a_i] \cdot \mathbb{P}[X_2 = b_j]$$

i.e. under H_0 , we wish we had

$$\frac{n_{i,j}}{n} \approx \frac{n_{i,\cdot}}{n} \cdot \frac{n_{\cdot,j}}{n} = \frac{n_{i,j}^\perp}{n} \text{ où } n_{i,\cdot} = \sum_{j=1}^J n_{i,j}, \quad n_{\cdot,j} = \sum_{i=1}^I n_{i,j}$$

Pearson's chi-square statistics is

$$Q = \sum_{i,j} \frac{(n_{i,j} - n_{i,j}^\perp)^2}{n_{i,j}^\perp} \sim \chi^2((I-1)(J-1))$$

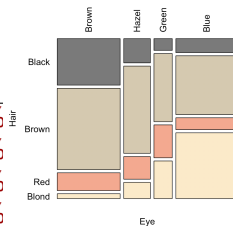
where $u_{i,j} = \frac{n_{i,j} - n_{i,j}^\perp}{\sqrt{n_{i,j}^\perp}}$ is the contribution of pair (i, j) .

Pearson's Chi-Square Test

	brown	hazel	green	blue	
black	63.0%	13.9%	4.6%	18.5%	100.0%
brown	41.6%	18.9%	10.1%	29.4%	100.0%
red	36.6%	19.7%	19.7%	23.9%	100.0%
blond	5.5%	7.9%	12.6%	74.0%	100.0%
	37.2%	15.7%	10.8%	36.3%	



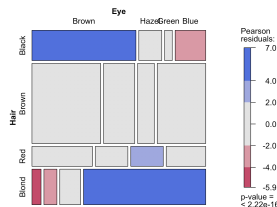
	brown	hazel	green	blue	
black	30.9%	16.1%	7.8%	9.3%	18.2%
brown	54.1%	58.1%	45.3%	39.1%	48.3%
red	11.8%	15.1%	21.9%	7.9%	12.0%
blond	3.2%	10.8%	25.0%	43.7%	21.5%
	100.0%	100.0%	100.0%	100.0%	



Pearson's Chi-Square Test

	brown	hazel	green	blue	
black	68	15	5	20	108
brown	119	54	29	84	286
red	26	14	14	17	71
blond	7	10	16	94	127
	220	93	64	215	

	brown	hazel	green	blue	
black	40	17	12	39	108
brown	106	45	31	104	286
red	26	11	8	26	71
blond	47	20	14	46	127
	220	93	64	215	



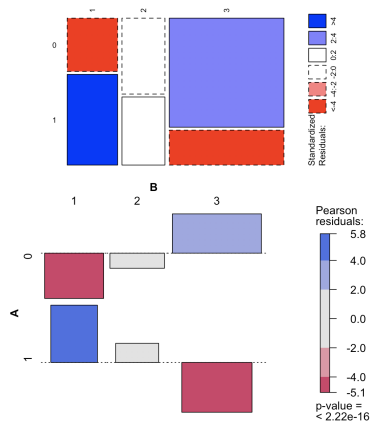
Compare $n_{i,j}$ and $n_{i,j}^\perp$

$$n_{i,j}^\perp = \frac{n_{i,\cdot} \cdot n_{\cdot,j}}{n}$$

Pearson's Chi-Square Test

$x_1 \in \{0, 1\}$, $\mathbf{x}_1 = (\mathbf{1}_0, \mathbf{1}_1)$ and $x_2 \in \{A, B, C\}$, $\mathbf{x}_2 = (\mathbf{1}_A, \mathbf{1}_B, \mathbf{1}_C)$

```
1 > (T=table(base$Survived,
2         base$Pclass))
3         Pclass
4 Survived    1    2    3
5           0   80   97  372
6           1  136   87  119
7
8 > chisq.test(T)
9
10      Pearson's Chi-squared test
11
12 X-squared = 91.081, df = 2,
13    p-value < 2.2e-16
14
15 > library("graphics")
16 > mosaicplot(T)
17 > library("vcd")
18 > assoc(T)
```



Correspondance Analysis

One can also look at **correspondance analysis** to visualize the associations of various categories (see #271).

