

Basic Inference

K. Gibert

*Department of Statistics and Operation Research
Knowledge Engineering and Machine Learning group*

Universitat Politècnica de Catalunya, Barcelona

karina.gibert@upc.edu

www.eio.upc.edu/homepages/karina

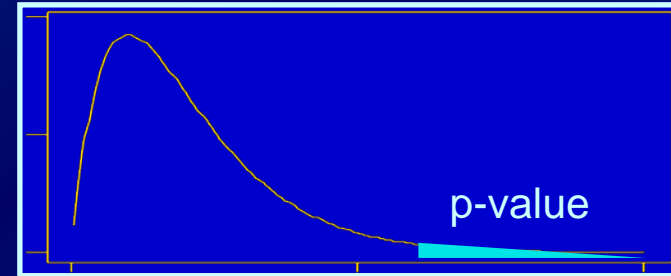
Statistical Significance test

1. Hypothesis

H_0 : (null hypothesis)

H_1 : (alternative hypothesis)

$f(s)$, assuming H_0 is true



2. Statistics (S):

Measure of proximity between sample and H_0

s_0

Known probability distribution under H_0 (f)

CARE: Technical conditions might be required

3. Observed value of statistics (s_0):

Compute statistics over sample data

4. p-value: $P_f(|S| > s_0)$

5. Decision rule: if p-value < alpha then Reject H_0

**p-value
depends on H_1**

Association between numerical variables

Correlation coefficient

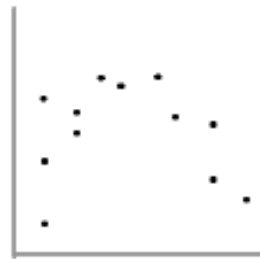
$$r_{x_j y} = \frac{\sum_i (x_{ij} - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_i (x_{ij} - \bar{x})^2 \sum_i (y_i - \bar{y})^2}} = \frac{s_{x_j y}}{s_x s_y}$$



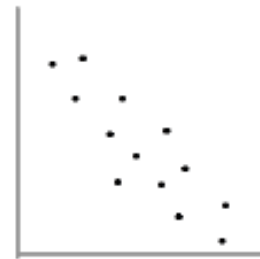
correlación 0,6



correlación 0



correlación 0
(relación cuadrática)



correlación -0,9

Test :

$$H_0: \text{cor}(X, Y) = 0$$

$$H_1: \text{cor}(X, Y) \neq 0$$

Only linear relationships

$$\sqrt{(n-2)} \frac{R}{\sqrt{(1-R^2)}} \sim t_{(n-2)}$$

Sheffer
Generalized
coefficient

Assessing association between categorical variables

The chi2 independence Test

Test: H_0 : X,Y are independent ($n_{kj}=np_kp_j \forall kj$)
 H_1 : X,Y are associated

Missperformance
if $n_{kj} < 5$

Statistics:

$$X^2 = \sum_{k=1}^p \sum_{j=1}^q \frac{(n_{kj} - \frac{n_k n_j}{n})^2}{\frac{n_k n_j}{n}} \sim \chi^2_{(p-1)(q-1)}$$

	1...	j	...	q
1	<div><div></div></div>			
k				
p				
		n_j		n

Care with
Simpson's
Paradox

Assessing association between categorical variables

The Simpson's Paradox

Apparently independent

or

Apparently dependent



**Lurking
Variables**

The Simpson's Paradox

Apparently independent

1978: Warren McClesky (black man) sentenced to death in Georgia for killing a police (white man)

He appealed to the US Supreme Court arguing racial bias of death penalty in Georgia

326 defendants in homicide in 20 Florida counties in between 1976-1977
[Radlet 1981] [Agresti 1990]

		YES	NO	Total	(%Yes)
Race of Suspect	White	39	308	347	(11.2%)
	Black	32	345	377	(8.5%)
	Total	71	653	724	(9.8%)

**Apparently
Independent**

When color of victim considered.....

White victim

		YES	NO	Total	(%Yes)
Race of Suspect	White	39	279	318	(12.3%)
	Black	29	121	150	(19.3%)
	Total	68	400	468	(14.5%)

Black victim

		YES	NO	Total	(%Yes)
Race of Suspect	White	0	29	29	(0.0%)
	Black	3	224	227	(1.3%)
	Total	3	253	256	(1.2%)

Assessing association between one categorical variable and one numerical

The F Test

Test: $H_0: \mu_{Y|X=x_1} = \mu_{Y|X=x_2} = \dots = \mu_{Y|X=x_s} = \mu(X, Y \text{ independent})$

$H_1: \exists x \in \{x_1, \dots, x_s\}: \mu_{Y|X=x} \neq \mu$

Requires
Normality

(X,Y associated)

Statistics:

$$F = \frac{S_B^2 / (q-1)}{S_W^2 / (n-q)} \sim F_{q-1, n-q}$$

$$S_W^2 = \sum_{k=1}^q \sum_{i=1}^{n_k} (x_{ki} - \bar{x}_k)^2$$

$$S_B^2 = \sum_{k=1}^q n_k (\bar{x}_k - \bar{x})^2$$

Kruskal-Wallis
test

Basic Inference

Karina Gibert

Dpt. Statistics and Operation Research

Knowledge Engineering and Machine Learning Research group

Universitat Politècnica de Catalunya-BarcelonaTech (Spain)

karina.gibert@upc.edu

www.eio.upc.edu/homepages/karina



Are there any questions?...