STATISTICA DESCRITTIVA B

(prof. Fortunato Pesarin)

FACOLTA' DI ASTRONOMIA Corso di Laurea in **Astronomia**

Correzione Prova scritta del 11 Luglio 2006

1) Per valutare l'area di attrazione della Facoltà di Scienze Statistiche su studenti di genere diverso, si è rilevata la residenza delle 180 matricole dell'A.A. 2003-04. I dati sono sintetizzati nella seguente tabella:

X	Y Residenza		
Genere	nella provincia di Padova	nel Veneto, fuori provincia di Padova	fuori del Veneto
maschio	30	48	22
femmina	38	25	17

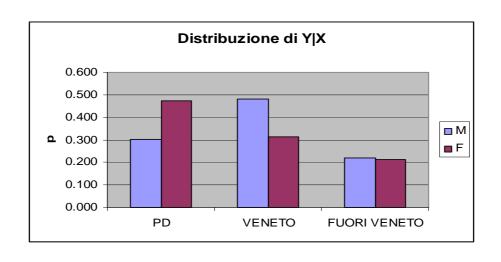
1.1) Qual è l'unità statistica, quale la natura delle variabili e su quale scala sono misurate.

L'unità statistica è la matricola iscritta alla facoltà di Scienze statistiche nell' A.A. 2003-2004. Le variabili rilevate sono entrambe qualitative su scala sconnessa.

1.2) Si ottengano le frequenze relative della variabile **Y**.

Y Residenza					
nella provincia di Padova	nel Veneto, fuori provincia di Padova	fuori del Veneto	Totale		
0.378	0.406	0.217	1		

1.3) Rappresentare graficamente la distribuzione percentuale della variabile Y condizionatamente alle modalità di X.



1.4) Si può affermare che la mutabilità della variabile Y sia la stessa tra i maschi e le femmine?

Indice di Gini:
$$G = 1 - \sum_{i=1}^{3} p_i^2$$
:

$$G_{\rm M} = 0.6312$$
 $G_{\rm F} = 0.6316$

$$G_F = 0.6316$$

Indice di Gini standardizzato:
$$kG/(k-1)$$
 $G_M = 0.9468$ $G_F = 0.9473$

$$G_{\rm M} = 0.9468$$

$$G_F = 0.9473$$

Entropia di Shannon
$$H = \sum_{i=1}^{3} p_i \log(p_i)$$
:

$$H_{\rm M} = 0.4545$$
 $H_{\rm F} = 0.4543$

$$H_F = 0.4543$$

Entropia di Shannon standardizzata:
$$H/log(k)$$
 ' $H_M = 0.9526$ ' $H_F = 0.9523$

$$^{\circ}$$
H_M = 0.9526

$$^{\circ}H_{\rm F} = 0.9523$$

Entrambi gli indici confermano che la mutabilità di Y nei due sessi sia praticamente la stessa.

1.5) Utilizzando un indice opportuno, si può affermare che il genere e la residenza siano dipendenti?

Frequenze attese in ipotesi di indipendenza

37.78	40.56	21.67	100
30.22	32.44	17.33	80
68	73	39	180

Indice χ^2 : 6.689137 > (3-1) × (2-1) = 2 \implies esiste dipendenza tra sesso e residenza.

$$V = \sqrt{\frac{\chi^2}{n \cdot \min(R - 1, C - 1)}} = 0.1927$$

2) Sugli studenti del corso di laurea in SGI si sono osservate tre variabili: X= età (in anni); Y= voto di Algebra (in trentesimi) e Z = voto di Statistica (in trentesimi). La matrice di correlazione tra le tre variabili è risultata uguale a:

	X	Y	Z
X	0	+0,48	+0,75
Y	+0,48	1	+0,81
Z	+0,75	+0,81	1

2.1) La tabella sopra contiene un dato impossibile, quale? Perchè?

Il dato impossibile è nella prima cella della prima riga: ρ (X, X) = 0, perchè questo equivale ad affermare che COV(X, X) = 0, ma per ciascuna variabile non degenere si ha $COV(X, X) = V(X) > 0 e \rho (X, X) = 1.$

2.2) Quale dei seguenti due modelli:

$$I^{\circ}$$
) $Z = \alpha_1 + \beta_1 X + \epsilon$ oppure II°) $Z = \alpha_2 + \beta_2 Y + \epsilon$

II°)
$$Z = \alpha_2 + \beta_2 Y +$$

conviene usare per prevedere Z e perché?

Conviene usare il secondo perchè Z risulta maggiormente correlata a Y di quanto lo sia X, perchè ρ (Z, Y) = 0.81 > 0.75 = ρ (Z, X)

2.3) Sono note le seguenti quantità: M(X) = 27.3; M(Y) = 25.4; M(Z) = 21.6; $M(X^2) = 751.59$; $M(Y^2) = 654.96$; $M(Z^2) = 471.26$, dove:

$$M(X) = \overline{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$$
 e $M(X^2) = \frac{1}{n} \sum_{i=1}^{n} x_i^2$.

Disegnare, su un sistema di assi cartesiani, la retta relativa al modello prescelto in (2.2).

Le covarianze si ottengono tenendo conto della definizione di correlazione:

$$\rho(X,Y) = \frac{\text{COV}(X,Y)}{\sqrt{V(X)V(Y)}} \quad \Rightarrow \quad \text{COV}(X,Y) = \rho(X,Y)\sqrt{V(X)V(Y)}$$

Le varianze si ottengono dalla relazione: $V(X) = M(X^2)-M(X)^2$.

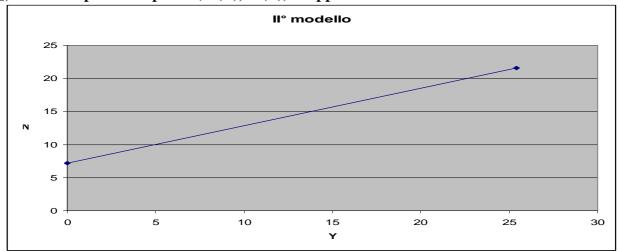
 $I^{\circ} \ modello: \ \ Z = \alpha_1 + \beta_1 \ X + \epsilon$

$$a_1 = M(Z)-b_1M(X)$$
 = 3.915
 $b_1 = COV(Z,X)/VAR(X)$ = 0.648

II° **modello:** $\mathbf{Z} = \alpha_2 + \beta_2 \mathbf{Y} + \boldsymbol{\epsilon}$

$$a_2 = M(Z)-b_2M(Y)$$
 = 7.352
 $b_2 = COV(Z,Y)/VAR(Y)$ = 0.561

Esiste un'uica retta passante per due punti: il primo punto è rappresentato dall'intercetta α_2 , il secondo punto è il punto (M(Y), M(Z)). Rappresentazione II° modello:



2.4) Per quale punto passa sicuramente la retta di regressione?

Passa sicuramente per i punti $(0, a_2)$ e (M(Y), M(Z)). Il primo per definizione di intercetta di una retta, il secondo perchè la stima di Z in corrispondenza di Y = M(Y) è:

$$Z' = a_2 + b_2 M(Y) = M(Z) - b_2 M(Y) + b_2 M(Y) = M(Z).$$

2.5) Si ottenga una misura della bontà del modello scelto.

I° modello:
$$R^2 = \rho^2 (Z, X) = 0.5625$$

II° modello: $R^2 = \rho^2 (Z, Y) = 0.6561$

2.6) Si ottega il coeffciente di correlazione parziale tra voto in statistica descrittiva e voto in lgebra lineare al netto dell'età.

La correlazione parziale tra Z e Y al netto di X è la correlazione tra Z' e Y', dove Z' e Y' sono i residui delle regressioni di Z su X e Y su X. In questo modo si evidenzia la correlazione dei dati "depurata" dall'effetto di X. Il calcolo dei residui è il seguente:

$$y_i' = y_i - \left(\overline{y} - \frac{\sigma_{XY}}{\sigma_X^2}\overline{x} + \frac{\sigma_{XY}}{\sigma_X^2}x_i\right) = (y_i - \overline{y}) - \frac{\sigma_{XY}}{\sigma_X^2}(x_i - \overline{x})$$

Analogamente:

$$z_i' = (z_i - \overline{z}) - \frac{\sigma_{XZ}}{\sigma_X^2} (x_i - \overline{x})$$

segue che:

$$\begin{split} V(Y') &= \sigma_{Y'}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \overline{y})^2 + \frac{\sigma_{XY}^2}{\sigma_X^4} \frac{1}{n} \sum_{i=1}^n (x_i - \overline{x})^2 - 2 \frac{\sigma_{XY}}{\sigma_X^2} \sum_{i=1}^n (x_i - \overline{x})(y_i - \overline{y}) \\ &= \sigma_Y^2 + \frac{\sigma_{XY}^2}{\sigma_X^4} \sigma_X^2 - 2 \frac{\sigma_{XY}}{\sigma_X^2} \sigma_{XY} = \sigma_Y^2 - \frac{\sigma_{XY}^2}{\sigma_X^2} = \sigma_Y^2 \left(1 - \frac{\sigma_{XY}^2}{\sigma_X^2 \sigma_Y^2} \right) = \sigma_Y^2 (1 - \rho_{XY}^2) \end{split}$$

Analogamente:

$$\sigma^2_{Z'} = \sigma^2_{Z} (1 - \rho^2_{XZ})$$

$$\sigma_{Y'Z'} = \sigma_{YZ} - \sigma_{XY} \sigma_{XZ} / \sigma_{X}^{2}$$

Infine:

$$\rho_{YZ|X} = \frac{\rho_{YZ} - \rho_{XY}\rho_{XZ}}{\sqrt{(1 - \rho_{XY}^2)(1 - \rho_{XZ}^2)}} = 0.79275$$