

Statistica inferenziale

Test di ipotesi

Bari, 7 Gennaio 2008

Coefficiente di correlazione

Abbiamo visto come misurare l'associazione tra variabili aleatorie discrete. L'analisi della correlazione mira a misurare l'associazione tra v.a. di tipo continuo, posto che tra esse vi sia una relazione di tipo lineare.

Ricordiamo che il coefficiente di correlazione ρ tra le v.a. X e Y é definito come:

$$\rho_{xy} = \frac{\mathbb{E}[(X - \mu_x)(Y - \mu_y)]}{\sqrt{\text{Var}(X)}\sqrt{\text{Var}(Y)}}$$

dove $\mu_x = \mathbb{E}(X)$ e $\mu_y = \mathbb{E}(Y)$. Inoltre $-1 \leq \rho_{xy} \leq 1$ e se $\rho_{xy} = 0$ allora le v.a. X e Y sono non correlate.

Coefficiente di correlazione

Consideriamo un campione accoppiato costituito da n osservazioni delle v.a. $X : N(\mu_x, \sigma_x)$ e $Y : N(\mu_y, \sigma_y)$:

$$(X_1, Y_1), \dots, (X_n, Y_n).$$

Uno stimatore di ρ é il coefficiente di correlazione di Pearson dato da:

$$r = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}}$$

Ipotesi nulla e alternativa

Vogliamo testare le seguenti ipotesi:

H_0 : $(\rho = 0)$ X e Y sono non correlate;

H_1 : $(\rho \neq 0)$ X e Y sono correlate.

Si può dimostrare che, sotto H_0 vera, la variabile:

$$T_0 = r \sqrt{\frac{n-2}{1-r^2}}$$

é una t di Student con $(n-2)$ gradi di libertà $t(n-2)$.

Regione critica e test

La regione critica D di livello α é data da:

$$D = \left\{ |T_0| > t_{1-\frac{\alpha}{2}}(n-2) \right\}.$$

Esecuzione del test:

- calcola il valore empirico t_0 di T_0 ;
- confronta t_0 con il quantile $t_{1-\frac{\alpha}{2}}(n-2)$;
- se $|t_0| > t_{1-\frac{\alpha}{2}}(n-2)$ allora rigetta H_0 e accetta H_1 , altrimenti accetta H_0 .

Esempio

Esaminare la relazione tra la percentuale di bambini vaccinati all'età di un anno contro DPT in un Paese ed il corrispondente tasso di mortalità in bambini al di sotto dei 5 anni in quel Paese.

Paese	% Vaccinati	Morti ogni 1000 nati vivi	Paese	% Vaccinati	Morti ogni 1000 nati vivi
<i>Bolivia</i>	77	118	<i>Giappone</i>	87	6
<i>Brasile</i>	69	65	<i>Grecia</i>	54	9
<i>Cambogia</i>	32	184	<i>India</i>	89	124
<i>Canada</i>	85	8	<i>Italia</i>	95	10
<i>Cina</i>	94	43	<i>Messico</i>	91	33
<i>Egitto</i>	89	55	<i>Polonia</i>	98	16
<i>Etiopia</i>	13	208	<i>RegnoUnito</i>	90	9
<i>Fed.Russa</i>	73	32	<i>Rep.Ceca</i>	99	12
<i>Finlandia</i>	95	7	<i>Senegal</i>	47	145
<i>Francia</i>	95	9	<i>Turchia</i>	76	87

Esempio (cont.)

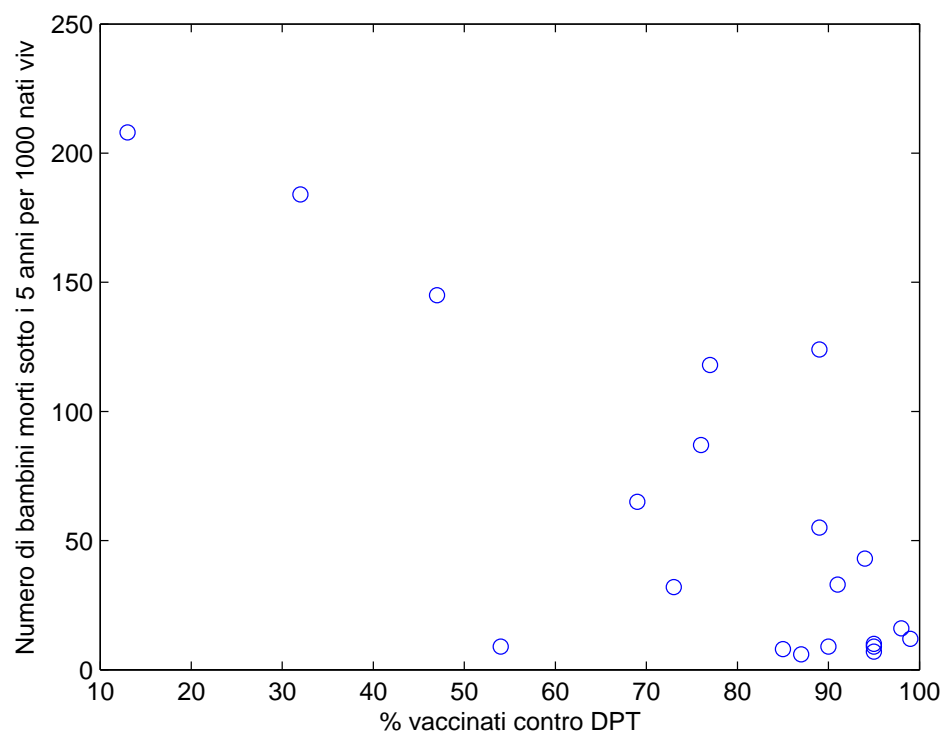


Figure 1: Diagramma a punti dei dati in tabella.

Esempio (cont.)

Calcoliamo le medie $\bar{x} = 77.4$, $\bar{y} = 59.0$ e le varianze $s_x^2 = 559.5$, $s_y^2 = 4078.8$. Inoltre $\sum_{i=1}^{20} (x_i - 77.4)(y_i - 59.0) = -22706$. Pertanto:

$$r = \frac{-22706}{\sqrt{(n-1)s_x^2}\sqrt{(n-1)s_y^2}} = \frac{-22706}{(19)(23.6541)(63.8658)} = -0.7911.$$

Esempio (cont.)

Per $\alpha = 0.05$, confrontiamo il valore di

$$t_0 = r \sqrt{\frac{n-2}{1-r^2}} = -0.7911 \sqrt{\frac{20-2}{1-(-0.7911)^2}} = -5.4870$$

con il quantile $t_{0.975}(20-2) = 2.1009$. Poiché $|t_0| > 2.1009$ allora a livello α rigettiamo l'ipotesi nulla e concludiamo che esiste una forte correlazione tra la percentuale di bambini vaccinati ed il tasso di mortalità al di sotto di 5 anni.

Regressione

Vogliamo affrontare il problema di predire o stimare il valore di una v.a. Y a partire dalla conoscenza del valore x assunto da una v.a. X .

Per esempio vogliamo predire il peso P di una persona conoscendo la sua altezza h .

Oppure, vogliamo analizzare il tipo di relazione esistente tra due v.a.

In generale le v.a. X e Y sono continue. Quando Y é discreta allora si parla di classificazione.

Funzione di regressione

Supponiamo di conoscere la funzione densità di probabilità congiunta $p(x, y)$ delle due v.a. X e Y . Allora si definisce "funzione di regressione" f la funzione:

$$f(x) = \mathbb{E}(Y|x) = \int_{-\infty}^{+\infty} yp(y|x)dy.$$

Nel caso in cui le v.a. X e Y sono Normali congiuntamente, allora la funzione di regressione é una retta.

Regressione lineare

Consideriamo un campione accoppiato $(x_1, y_1), \dots, (x_n, y_n)$. Determiniamo i coefficienti w e b della retta $y = wx + b$ che meglio approssima i dati. A tale scopo consideriamo l'errore quadratico medio:

$$\epsilon(w, b) = \frac{1}{n} \sum_{i=1}^n (y_i - (wx_i + b))^2$$

Determiniamo i coefficienti w^* e b^* che minimizzano l'errore quadratico medio. Questa fase viene detta "training".

Calcolo dei coefficienti

A tale scopo poniamo uguale a zero le derivate parziali di ϵ rispetto a w e b :

$$\begin{aligned}\frac{\partial \epsilon}{\partial w} = 0 &\Rightarrow \\ -\frac{2}{n} \sum_{i=1}^n (y_i - (wx_i + b))x_i = 0 &\Rightarrow \sum_{i=1}^n (y_i x_i - wx_i^2 - bx_i) = 0 \Rightarrow \\ \sum_{i=1}^n y_i x_i = w \sum_{i=1}^n x_i^2 + b \sum_{i=1}^n x_i. &\quad (1)\end{aligned}$$

Calcolo dei coefficienti (cont.)

Analogamente:

$$\frac{\partial \epsilon}{\partial b} = 0 \Rightarrow$$

$$-\frac{2}{n} \sum_{i=1}^n (y_i - (wx_i + b)) = 0 \Rightarrow \sum_{i=1}^n (y_i - wx_i - b) = 0 \Rightarrow$$

$$\sum_{i=1}^n y_i = w \sum_{i=1}^n x_i + nb. \quad (2)$$

Le equazioni (1) e (2) costituiscono un sistema lineare di 2 equazioni nelle incognite w e b .

Calcolo dei coefficienti (cont.)

Da (2) segue che $b = \bar{y} - w\bar{x}$. Sostituendo in (1) si ha:

$$\sum_{i=1}^n y_i x_i = w \sum_{i=1}^n x_i^2 + (\bar{y} - w\bar{x}) \sum_{i=1}^n x_i \Rightarrow$$

$$\sum_{i=1}^n y_i x_i = w \sum_{i=1}^n x_i^2 + \bar{y} \sum_{i=1}^n x_i - w\bar{x} \sum_{i=1}^n x_i \Rightarrow$$

$$w \left(\sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \right) = \sum_{i=1}^n x_i y_i - \frac{1}{n} \sum_{i=1}^n y_i \sum_{i=1}^n x_i$$

Calcolo dei coefficienti (cont.)

Da cui:

$$w^* = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n y_i \sum_{i=1}^n x_i}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2} \quad (3)$$

Sostituendo w^* in $b = \bar{y} - w\bar{x}$ si ha:

$$b^* = \frac{\sum_{i=1}^n x_i^2 \sum_{i=1}^n y_i - \sum_{i=1}^n x_i \sum_{i=1}^n x_i y_i}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2}. \quad (4)$$

Errore di training

Una volta calcolati w^* e b^* che minimizzano l'errore quadratico medio, possiamo calcolare l'errore di training:

$$\epsilon(w^*, b^*) = \frac{1}{n} \sum_{i=1}^n (y_i - (w^* x_i + b^*))^2$$

Piú piccolo é ϵ , meglio i dati sono descritti da una retta.

Errore di predizione

Supponiamo di voler conoscere il valore y associato ad un "nuovo" x , ossia ad un dato non presente nel nostro campione. Stimiamo y attraverso $\hat{y} = w^*x + b^*$. Qual é l'errore medio che si commette quando stimiamo y con \hat{y} ?

Si definisce errore di predizione o errore di generalizzazione la quantità:

$$\mathbb{E}[(y - \hat{y})^2] = \int (y - (w^*x + b^*))^2 p(x, y) dx dy$$

Poiché la densità $p(x, y)$ é incognita, allora l'errore di predizione non si può calcolare.

Leave-One-Out Error (Luntz and Brailovsky 1969)

Una stima dell'errore di generalizzazione che non fa uso della conoscenza della densità di probabilità dei dati é basata sulla seguente procedura: per ogni $i = 1, 2, \dots, n$:

- escludi l'esempio (x_i, y_i) dal campione;
- calcola w_i^* e b_i^* sui dati rimanenti $(x_1, y_1), \dots, (x_{i-1}, y_{i-1}), (x_{i+1}, y_{i+1}), \dots, (x_n, y_n)$;
- calcola l'errore $\epsilon_i = (y_i - (w_i^* x_i + b_i^*))^2$.

Allora $LOOE = \frac{1}{n} \sum_{i=1}^n \epsilon_i$ é una stima non distorta dell'errore di generalizzazione.