

Spesso ad una fissata popolazione vengono associate coppe di dati: età e peso, età e reddito, reddito e scelte politiche, ecc.) - In questo caso l'insieme dei dati è dato da coppe numeriche

$$(x_1, y_1), \dots, (x_N, y_N).$$

In genere, non ci sarà una legge funzionale precisa che legni le y_i alle x_i . Ma in certi casi potrà succedere che a valori $x_i < \bar{x}$ corrispondano valori $y_i < \bar{y}$, e che viceversa quando $x_i > \bar{x}$ si abbia anche $y_i > \bar{y}$. Oppure potrà succedere l'opposto, cioè che per $x_i < \bar{x}$ si abbia $y_i > \bar{y}$ e per $x_i > \bar{x}$ si abbia $y_i < \bar{y}$. O magari non succederà nulla di tutto ciò.

Una misura numerica del modo in cui le x_i si associano alle y_i è data dalla covarianza di X e Y ,

$$\text{Cov}(X, Y) = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y}).$$

Nei due casi prima formulati, il 1° corrisponde a una covarianza positiva, il 2° a una covarianza negativa.

Se $\text{Cov}(X, Y) = 0$ si dice che X e Y sono non correlate.

Elenchiamo le proprietà della covarianza:

$$\text{Cov}(X, Y) = \frac{1}{N} \sum_{i=1}^N x_i y_i - \bar{x} \bar{y} = E[XY] - E[X]E[Y]$$

$$\text{Cov}(X, X) = \sigma_X^2 = \text{Var}[X],$$

$$\text{Cov}(aX + b, cY + d) = ac \text{Cov}(X, Y);$$

Inoltre

$$\text{Var}[X+Y] = \text{Var}[X] + \text{Var}[Y] + 2 \text{Cov}(X, Y).$$

Il numero

$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$$

si chiama coefficiente di correlazione di X e Y.

Si ha $\rho(aX+b, cY+d) = \rho(X, Y)$ se $ac > 0$. Inoltre

$$-1 \leq \rho(X, Y) \leq 1:$$

infatti per la disuguaglianza di Cauchy-Schwarz

$$\begin{aligned}
|\text{Cov}(X, Y)| &= \frac{1}{N} \left| \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y}) \right| \leq \\
&\leq \frac{1}{N} \left[\sum_{i=1}^N (x_i - \bar{x})^2 \right]^{1/2} \left[\sum_{i=1}^N (y_i - \bar{y})^2 \right]^{1/2} = \\
&= \sigma_X \sigma_Y.
\end{aligned}$$

È facile provare, inoltre, che

$$|\rho(X, Y)| = 1 \iff \exists a \neq 0, b \in \mathbb{R}: Y = aX + b.$$

Infatti, se $Y = aX + b$, dalle proprietà della covarianza segue

$$\text{Cov}(X, Y) = a \cdot \text{Cov}(X, X) = a \text{Var}[X];$$

d'altronde

$$\text{Var}[Y] = a^2 \text{Var}[X],$$

così che

$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{a}{|a|} = \begin{cases} +1 & \text{se } a > 0 \\ -1 & \text{se } a < 0. \end{cases}$$

Viceversa, se $\rho=1$ allora

$$\text{Var} \left[\frac{X}{\sigma_X} - \frac{Y}{\sigma_Y} \right] = \frac{\text{Var}[X]}{\sigma_X^2} + \frac{\text{Var}[Y]}{\sigma_Y^2} - 2 \frac{\text{Cov}(X,Y)}{\sigma_X \sigma_Y} = 2(1-\rho) = 0,$$

per cui

$$\frac{x_i}{\sigma_X} = \frac{y_i}{\sigma_Y}, \quad i=1, \dots, N$$

ossia esiste $c \in \mathbb{R}$ tale che

$$\frac{x_i}{\sigma_X} - \frac{y_i}{\sigma_Y} = c, \quad i=1, \dots, N$$

Dunque tutte le coppie (x_i, y_i) appartengono alla retta

$$\frac{x}{\sigma_X} - \frac{y}{\sigma_Y} = c,$$

ossia

$$y = \frac{\sigma_Y}{\sigma_X} x - c \sigma_Y$$

vale a dire

$$y = aX + b \quad \text{con } a = \frac{\sigma_Y}{\sigma_X}, \quad b = -c\sigma_Y.$$

Similmente, se $\rho=-1$ si trova

$$\text{Var} \left[\frac{X}{\sigma_X} + \frac{Y}{\sigma_Y} \right] = 2(1+\rho) = 0,$$

e analogamente si arriva a

$$Y = aX + b \quad \text{con } a = \frac{\sigma_Y}{\sigma_X}, \quad b = c\sigma_Y.$$

Il coefficiente di correlazione esprime una misura della dipendenza lineare fra X e Y . L'insieme delle coppie (x_i, y_i) è il grafico di dispersione: se l'insieme di questi punti è

molto addensato intorno a una retta, allora $\rho(X, Y)$ sarà vicino a $+1$ o a -1 , a seconda che il coefficiente angolare di tale retta sia positivo o negativo. Se invece l'insieme dei punti (x_i, y_i) è abbastanza "rotondo", allora $\rho(X, Y)$ sarà vicino a 0 . 415

Esempio L'insieme dei dati

$$(-2, 4), (-1, 1), (0, 0), (1, 1), (2, 4)$$

soddisfa la relazione $Y = X^2$, e si ha

$$E[X] = 0, \quad E[Y] = 2, \quad \text{Var}[X] = 2, \quad \text{Var}[Y] = \frac{14}{5}, \quad \text{Cov}(X, Y) = 0,$$

e dunque $\rho(X, Y) = 0$.

In molte applicazioni, le variabili X, Y possono essere legate da una relazione lineare $Y = aX + b$, ma, a causa di errori di misura, non è possibile determinare a, b ; in altri casi, tale legame non è precisamente lineare. Tuttavia, se si ritiene che il legame statistico fra X e Y sia approssimabile con un'opportuna funzione lineare $Y = aX + b$, si deve cercare un metodo per selezionare la retta "migliore" fra le infinite possibili.

La logica è la seguente: se i punti (x_i, y_i) appartenessero

tutti alle stesse retta $y = ax + b$, avremmo

$$(y_i - ax_i - b)^2 = 0, \quad i = 1, \dots, N.$$

Se una tale retta non esiste, ci si accontenta di determinare la retta che rende minimo la funzione

$$f(a, b) = \sum_{i=1}^N (y_i - ax_i - b)^2.$$

Tale retta si chiama retta di regressione di Y su X e il metodo è quello dei minimi quadrati. Per trovare i coefficienti di a e b basterà annullare il gradiente di f:

$$\nabla f(a, b) = 0 \iff \begin{cases} -2 \sum_{i=1}^N x_i (y_i - ax_i - b) = 0 \\ -2 \sum_{i=1}^N (y_i - ax_i - b) = 0, \end{cases}$$

ossia

$$\begin{cases} \sum_{i=1}^N x_i y_i - a \sum_{i=1}^N x_i^2 - b \sum_{i=1}^N x_i = 0 \\ \sum_{i=1}^N y_i - a \sum_{i=1}^N x_i - Nb = 0. \end{cases}$$

Ciò equivale a

$$\begin{cases} N E[XY] - Na E[X^2] - Nb E[X] = 0 \\ N E[Y] - Na E[X] - Nb = 0, \end{cases}$$

da cui, semplificando N e ricavando b dalla 2^a equazione, si ottiene

$$\begin{cases} a = \frac{\text{Cov}(X, Y)}{\sigma_x^2} = \rho(X, Y) \frac{\sigma_y}{\sigma_x} \\ b = E[Y] - E[X] \rho(X, Y) \frac{\sigma_y}{\sigma_x}. \end{cases}$$

La retta di regressione è dunque

$$y - E[Y] = \rho \frac{\sigma_y}{\sigma_x} (x - E[X]),$$

ovvero, in forma più simmetrica,

$$\frac{y - E[Y]}{\sigma_y} = \rho \frac{x - E[X]}{\sigma_x}.$$

In modo completamente analogo, la retta di regressione di X su Y, che ha la forma $x = ay + b$, è data da

$$x - E[X] = \rho \frac{\sigma_x}{\sigma_y} (y - E[Y])$$

ovvero

$$\frac{x - E[X]}{\sigma_x} = \rho \frac{y - E[Y]}{\sigma_y}.$$

Si noti che il punto $(E[X], E[Y])$ appartiene sempre alla retta di regressione.

Si noti infine che

$$\rho(X, Y) = \pm 1 \iff \text{Var} \left[\frac{y - E[Y]}{\sigma_y} - \rho \frac{x - E[X]}{\sigma_x} \right] = 0,$$