

CORRELATION COEFFICIENT

Utile per analizzare la correlazione fra due variabili numeriche.

$$r_{A,B} = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{n \cdot \sigma_A \cdot \sigma_B} = \frac{\sum_{i=1}^n (a_i b_i) - n \bar{A} \bar{B}}{n \cdot \sigma_A \cdot \sigma_B}$$

Dove:

- a_i e b_i : valori attuali in A e B
- \bar{A} e \bar{B} : media dei valori in A e B
- σ_A e σ_B : deviazione standard in A e B
- n : numero di tuple nel DataSet

L'intuizione della formula è la seguente:

Se assumiamo che A e B hanno lo stesso trend, per esempio entrambi crescono in un determinato punto, quello che succede nella formula è che il prodotto al numeratore mantiene il segno positivo. Se per esempio assumiamo che i valori siano ordinati in modo crescente, se partiamo da a_0 e b_0 , essi sono necessariamente $a_0 < \bar{A}$ e $b_0 < \bar{B}$ perché il trend è crescente. Sotto questa luce abbiamo:

$$\begin{aligned}(a_i - \bar{A}) &\Rightarrow \text{NEGATIVO} \\ (b_i - \bar{B}) &\Rightarrow \text{NEGATIVO} \\ (a_i - \bar{A})(b_i - \bar{B}) &\Rightarrow \text{POSITIVO}\end{aligned}$$

Ma questo prodotto, sotto queste ipotesi è sempre positivo.

Le cose cambiano il trend, con A che cresce e B decresce. Ora otteniamo che il prodotto è sempre negativo.

$$\begin{array}{llll} (a_1 - \bar{A}) & \Rightarrow \text{NEGATIVO} & \dots & (a_n - \bar{A}) \Rightarrow \text{POSITIVO} \\ (b_1 - \bar{B}) & \Rightarrow \text{POSITIVO} & \dots & (b_n - \bar{B}) \Rightarrow \text{NEGATIVO} \\ (a_1 - \bar{A})(b_1 - \bar{B}) & \Rightarrow \text{NEGATIVO} & \dots & (a_n - \bar{A})(b_n - \bar{B}) \Rightarrow \text{NEGATIVO} \end{array}$$

SS $\sum_{A,B} > 0$ ALLORA A e B **POSITIVAMENTE CORRELATE**

SS $\sum_{A,B} < 0$ ALLORA A e B **NEGATIVAMENTE CORRELATE**

SS $\sum_{A,B} \cong 0$ ALLORA A e B **INDIPENDENTI**

Quest'ultimo caso è motivato dal fatto che, essendo A e B indipendenti, il prodotto a volte è negativo altre è positivo (Confrontiamo valori randomici), quindi nella somma a volte sottraggo, altre sommo.

ATTENZIONE

Questo metodo rileva solo dipendenze lineari.

VALORI: $-1 < \sum_{A,B} < 1$

Standardizzazione

Per valutare la correlazione, i valori vengono spesso standardizzati:

$$a'_k = \frac{(a_k - \bar{A})}{\sigma_A}, \quad b'_k = \frac{(b_k - \bar{B})}{\sigma_B}$$

COVARIANZA

$$\text{Cov}(A, B) = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{n}$$

Da cui possiamo dire che $\tau_{A,B} = \frac{\text{Cov}(A,B)}{\sigma_A \sigma_B}$

Valgono le stesse considerazioni fatte per $\tau_{A,B}$, tranne per $\text{Cov}(A,B) \neq 0$.

Per $\text{Cov}(A,B) \neq 0$, A e B sono indipendenti. Il contrario non è sempre vero, ovvero che:

$$A, B : \text{indipendenti} \Rightarrow \text{Cov}(A,B) \neq 0$$