

Concetti chiave e regole

I valori di sintesi

Per sintetizzare i dati di una distribuzione statistica si usano alcuni indici:

- la **media aritmetica** che è il rapporto tra la somma di tutti i dati e il loro numero:

– media semplice $M = \frac{\sum x_i}{n}$

– media ponderata $M = \frac{\sum x_i \cdot f_i}{\sum f_i}$

- la **moda**, che è il termine cui corrisponde la massima frequenza
- la **mediana**, che, una volta disposti i dati in ordine crescente o decrescente, è il termine che occupa il posto centrale della distribuzione.

La variabilità

Per studiare la dispersione dei dati attorno al valore medio si calcolano:

- lo **scarto quadratico medio**: media quadratica degli scarti dalla media aritmetica

– per dati semplici $\sigma = \sqrt{\frac{\sum (x_i - M)^2}{n}}$

– per dati ponderati $\sigma = \sqrt{\frac{\sum (x_i - M)^2 f_i}{\sum f_i}}$

- la **varianza**: quadrato dello scarto quadratico medio;
 - un'altra possibile formula per il calcolo della varianza è $\sigma^2 = \frac{\sum x_i^2}{n} - M^2$
- il **coefficiente di variazione**: rapporto fra σ e M : $CV = \frac{\sigma}{M}$

I rapporti statistici

Si chiama rapporto statistico il rapporto fra due dati di cui almeno uno sia di tipo statistico. In particolare si evidenziano:

- il rapporto di **composizione**: rapporto fra la frequenza assoluta e la totalità delle osservazioni (coincide con la frequenza relativa)
- rapporto di **coesistenza**: rapporto fra le frequenze di due fenomeni diversi riferiti allo stesso luogo o tempo
- rapporto di **derivazione**: rapporto fra le intensità di due fenomeni di cui il primo dipende dal secondo
- rapporto di **durata**: rapporto fra la consistenza media e il movimento medio
- rapporto di **ripetizione**: reciproco del rapporto di durata, indica quante volte una popolazione si rinnova nell'unità di tempo
- rapporto di **densità**: rapporto fra la frequenza di un fenomeno e la dimensione del campo su cui è stata fatta la rilevazione.

Altri rapporti statistici sono i **numeri indice** che esprimono il rapporto tra i dati di una serie statistica e uno di essi preso come base di riferimento.

La dipendenza statistica

Nell'analisi di un fenomeno statistico si cerca spesso di indagare sulla possibile dipendenza di due caratteri uno dall'altro.

La **teoria della correlazione** studia tale dipendenza nel caso in cui X e Y sono delle variabili statistiche, mettendo in rilievo, in particolare, se vi è dipendenza di tipo lineare.

L'indice che dà informazioni sulla dipendenza lineare è quello di **Bravais-Pearson** che è così definito:

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

essendo σ_X e σ_Y le deviazioni standard di X e di Y e $\text{cov}(X, Y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n}$

Se capita che:

- $\rho = \pm 1$ le due variabili sono perfettamente correlate e quindi, essendoci perfetta dipendenza lineare, i dati si distribuiscono su una retta
- $0 < \rho < 1$ le variabili sono correlate positivamente e più ρ si avvicina a 1, più la dipendenza lineare è forte, cioè i dati tendono a distribuirsi lungo una retta di coefficiente angolare positivo
- $-1 < \rho < 0$ le variabili sono correlate negativamente e più ρ si avvicina a -1 , più i dati tendono a distribuirsi lungo una retta di coefficiente angolare negativo
- $\rho = 0$ le variabili non sono correlate, cioè non vi è dipendenza lineare, anche se non si possono escludere altri tipi di dipendenza.

L'interpolazione statistica

Una funzione di **interpolazione statistica** è una curva che passa fra i punti (x_i, y_i) corrispondenti ai dati rilevati. Per scegliere la funzione che meglio approssima i dati, si applica il metodo dei minimi quadrati che consiste nello scegliere la funzione $f(x)$ che rende minima la quantità

$S = \sum_{i=1}^n [f(x_i) - y_i]^2$ che esprime la somma dei quadrati delle distanze dei punti rilevati dalla funzione $f(x)$.

La funzione interpolante più semplice da determinare è la funzione lineare che in molte situazioni riesce ad approssimare sufficientemente bene i dati. Essa ha equazione $y = mx + q$ dove i coefficienti m e q sono espressi dalle formule

$$m = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2} \quad q = \frac{\sum x_i^2 \sum y_i - \sum x_i \sum x_i y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

Un metodo alternativo per determinare la retta di interpolazione statistica è quello del baricentro dove la retta ha equazione

$$y - \bar{y} = m(x - \bar{x}) \quad \text{essendo:} \quad \bar{x} = \frac{\sum x_i}{n} \quad \bar{y} = \frac{\sum y_i}{n} \quad m = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

Se i dati rilevati riguardano una serie storica, la retta di interpolazione può descrivere il comportamento tendenziale o **trend** della variabile Y per valori di X che vanno oltre quelli rilevati.

La retta di regressione

Il coefficiente di Bravais-Pearson indica se fra due variabili vi è la tendenza ad una dipendenza di tipo lineare, ma non dà modo di trovare l'espressione di questa dipendenza.

La teoria della **regressione** consente invece di esprimere la dipendenza della variabile X da Y , o viceversa, determinando rispettivamente la retta di regressione di X su Y o di Y su X ; queste due rette si determinano con il metodo del baricentro oppure dei minimi quadrati ed hanno equazioni:

- regressione di Y su X : $y = ax + b$
- regressione di X su Y : $x = cy + d$

dove i coefficienti a , b , c , d si trovano con le formule relative all'interpolazione lineare, con l'avvertenza di scambiare i valori x_i con i valori y_i per la determinazione di c e d .