

D. Statistiche descrittive

Le statistiche descrittive si occupano della presentazione e della sintesi di ~~un~~ insiemi di dati (di solito ~~numerici~~).

Presentazione dei dati in forma di tabelle o grafici (diagrammi); ~~sp~~ spesso utile classificazione (raggruppamenti in classi) dei dati.

Sintesi dei dati attraverso ~~le~~ statistiche, cioè funzioni numeriche definite sull'insieme dei dati ("calcolabili dei dati").

Consideriamo prima dati univariati, cioè ~~insiemi~~ ^{successioni} finite

di numeri reali:

Sia $X = \{(x_i)_{i \in \{1, \dots, n\}}\} \subset \mathbb{R}$ l'insieme dei dati
 (il campione). \hookrightarrow ^{lunghezza} ~~campione~~ $|X| = n$ è la

numerosità del campione.

Statistiche elementari per la sintesi dei dati:



a) "centro dei dati":

media campionaria, mediana campionaria,
moda o valori modal

b) "dispersione dei dati":

varianza campionaria e deviazione standard campionaria

c) "distribuzione dei dati": percentili campionari

(casi particolari: quartili, mediane)

2) Statistiche del "centro dei dati":

Sia $(x_i)_{i \in \{1, \dots, n\}} \subset \mathbb{R}$ il campione.

Def.: La medie campionaria \bar{x} del campione è

dato da
$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad (\text{medie aritmetica della successione})$$

Esempio: $x_i = i, \quad i \in \{1, \dots, n\}$

$$\leadsto \sum_{i=1}^n x_i = \sum_{i=1}^n i = \frac{n(n+1)}{2}$$

$$\leadsto \bar{x} = \frac{n+1}{2}$$

In questo caso, la medie campionaria è il "valore centrale".

La medie campionaria è ~~molto~~ ^{fortemente} influenzata da valori estremi, anche poco frequenti: Ad esempio,

~~n=100~~ $n = 100, \quad x_i = i \text{ per } i \in \{1, \dots, 99\},$
 $x_{100} = 10^6$

$$\leadsto \bar{x} = \frac{1}{100} \left(\frac{100}{2} \cdot 99 + 10^6 \right) = 10049,5.$$

La medie campionaria si può riscrivere come medie pesate dei valori assunti dai dati:



Sia v_1, \dots, v_K un'enumerazione dell'insieme dei valori dei dati, cioè

$$\{v_j : j \in \{1, \dots, K\}\} = \{x_i : i \in \{1, \dots, n\}\} \text{ e } v_j \neq v_l \text{ per } j \neq l.$$

Sia f_j la frequenza assoluta del valore v_j :

$$f_j = \#\{i \in \{1, \dots, n\} : x_i = v_j\}.$$

La media campionaria $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ si può allora scrivere come la media pesata

$$\bar{x} = \sum_{j=1}^K \underbrace{\frac{f_j}{n}}_{\text{frequenza relativa del valore } v_j} \cdot v_j$$

media dei valori pesata (ponderata) dalle corrispondenti frequenze relative.

Def.: I valori che hanno frequenze massime, cioè i valori v_j con $j \in \operatorname{argmax} \{f_j : j \in \{1, \dots, K\}\}$ si dicono valori modali. Se ~~il~~ ^{esiste} ~~valore~~ un unico valore di frequenze massime, cioè $\operatorname{argmax} \{f_j : j \in \{1, \dots, K\}\} = \{j_x\}$ per un (unico) $j_x \in \{1, \dots, K\}$, allora v_{j_x} si dice moda campionaria.

Sia σ una permutazione degli indici $\{1, \dots, n\}$ (cioè σ una biiezione $\{1, \dots, n\} \rightarrow \{1, \dots, n\}$)

tale che $X_{\sigma(1)} \leq X_{\sigma(2)} \leq \dots \leq X_{\sigma(n)}$

(ovvero $X_{\sigma(i)} \leq X_{\sigma(i+1)}$ per ogni $i \in \{1, \dots, n-1\}$).

La successione $(X_{\sigma(i)})_{i \in \{1, \dots, n\}}$ è quindi un riordinamento del campione in ordine crescente.

Def.: La mediana campionaria del campione è data dal valore

$$\bar{m} = \begin{cases} X_{\sigma(\frac{n+1}{2})} & \text{se } n \text{ è dispari,} \\ \frac{1}{2} (X_{\sigma(\frac{n}{2})} + X_{\sigma(\frac{n+2}{2})}) & \text{se } n \text{ è pari.} \end{cases}$$

Esempio (cf. media campionaria):

dati già
in ordine crescente

1) $X_i = i$, $i \in \{1, \dots, n\}$, allora

$$\bar{m} = \frac{n+1}{2} \quad (\text{sia per } n \text{ pari che dispari}).$$

In questo caso, $\bar{m} = \bar{x}$, media e mediana campionaria coincidono!



$$2) \quad n = 100, \quad X_i = i, \quad i \in \{1, \dots, 99\},$$

$$X_{100} = 10^6$$

dati già
in ordine crescente

$$\leadsto \quad \bar{m} = \frac{1}{2} \left(\overset{50}{\cancel{49}} + \overset{51}{\cancel{49.5}} \right) = 50.5,$$

$$\text{mentre } \bar{x} = 10049.5 \quad !$$

Infatti, \bar{m} non cambierebbe se X_{100} assumesse

un qualsiasi altro valore $\geq \cancel{49.5} 51$.

La mediana divide i dati in ~~due~~ ^{tre} parti; ~~uguali~~

$$\{i \in \{1, \dots, n\} : X_i < \bar{m}\}, \quad \{i \in \{1, \dots, n\} : X_i = \bar{m}\}, \quad \{i \in \{1, \dots, n\} : X_i > \bar{m}\}$$

(possibilmente vuoto)

in modo che

$$\# \{i \in \{1, \dots, n\} : X_i < \bar{m}\} = \# \{i \in \{1, \dots, n\} : X_i > \bar{m}\}.$$

"metà dei dati sotto, metà sopra la mediana".

Attenzione: conta anche la frequenza dei valori.

b) Statistiche della "dispersione dei dati"

Sia $(x_i)_{i=1, \dots, n} \in \mathbb{R}$ il campione.

Vogliamo avere una misura ~~di~~ per la dispersione dei dati intorno alla media campionaria; si dà più peso a deviazioni grandi, meno a deviazioni piccole.

Def.: Per $n \geq 2$, la varianza campionaria del campione è data da

$$s^2 = \frac{1}{n-1} \left(\sum_{i=1}^n (x_i - \bar{x})^2 \right)$$

$$\text{La quantità } s = \sqrt{\frac{1}{n-1} \left(\sum_{i=1}^n (x_i - \bar{x})^2 \right)}$$

si dice deviazione standard campionaria.

Osservazioni:

- 1) La varianza campionaria si basa ~~sulla~~ sullo scarto quadratico tra media campionaria e valore dei dati; è "quasi" uguale alla media campionaria degli scarti quadratici (asintoticamente uguale, $n \rightarrow \infty$).
- 2) Il prefattore $\frac{1}{n-1}$ invece di $\frac{1}{n}$ serve per ottenere uno stimatore "corretto" o "non distorto" quando i dati sono quantità statistiche (cf. infre).



3) Sieno ~~anche~~ $a, b \in \mathbb{R}$. Poniamo

$$y_i = a \cdot x_i + b, \quad i \in \{1, \dots, n\}.$$

Sia s_y^2 la varianza campionaria di (y_i) ,
 s_x^2 " " " " " (x_i) .

Allora $s_y^2 = a^2 \cdot s_x^2$ e $s_y = |a| \cdot s_x$

Il campione (y_i) è una trasformazione lineare-affine di (x_i) . La varianza campionaria non dipende dalla traslazione b , e si trasforma con a^2 , mentre la deviazione standard si trasforma con il fattore $|a|$ ~~mantenere la stessa unità~~ ~~mantenere la stessa unità~~ (oppure a se $a > 0$). La deviazione standard ha quindi la stessa unità di misura dei dati.

4) Sviluppando i quadrati, la varianza campionaria si riscrive come

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n x_i^2 - \frac{n}{n-1} \bar{x}^2.$$

utile
per i
calcoli

Esempio: ➤ Foglio 1

c) Statistiche per la distribuzione dei dati

Come per la mediana campionaria, sia

$(x_{(i)})_{i \in N}$ un riordinamento dei dati in ordine crescente.

Def.: Il percentile K-esimo, con $K \in \{0, \dots, 100\}$,
del campione è dato dal valore

$$\bar{p}_K = \begin{cases} x_{(\lceil \frac{n \cdot K}{100} \rceil)} & \text{se } n \cdot \frac{K}{100} \notin \mathbb{N}_0, \\ \frac{1}{2} (x_{(\frac{n \cdot K}{100})} + x_{(\frac{n \cdot K}{100} + 1)}) & \text{se } n \cdot \frac{K}{100} \in \mathbb{N}_0. \end{cases}$$

[~~Def.~~ Per $x \in \mathbb{R}$ poniamo $\lceil x \rceil = \min \{l \in \mathbb{Z} : l \geq x\}$.]
 \uparrow arrotondamento per eccesso
 \swarrow numeri interi

Osservazioni:

1) Il percentile 50-esimo coincide con la mediana.

Si dicono primo, secondo, terzo e ~~quarto~~ ^{quarto} quartile
i percentili 25-esimo, 50-esimo, 75-esimo e 100-esimo.

2) La proporzione degli indici $i \in \{1, \dots, n\}$ tali che $x_i \leq \bar{p}_K$
è maggiore di o uguale a $\frac{K}{100}$:

$$\frac{\#\{i \in \{1, \dots, n\} : x_i \leq \bar{p}_K\}}{n} \geq \frac{K}{100}. \quad \text{Più precisamente,}$$

$$\#\{i \in \{1, \dots, n\} : x_i \leq \bar{p}_K\} \geq \lceil n \cdot \frac{K}{100} \rceil, \quad \#\{i \in \{1, \dots, n\} : x_i > \bar{p}_K\} \leq n - \lceil n \cdot \frac{K}{100} \rceil.$$

La varianza campionaria (equivalentemente, la deviazione standard) permette di stimare la proporzione dei dati che sono vicini alla (o lontani dalla) media campionaria, grazie alla disuguaglianza di Chebyshev. L'unità di misura spesso più conveniente è la deviazione standard

Disuguaglianza di Chebyshev (versione campionaria):

Siano \bar{X} la media ~~e~~ ^{sic} s la deviazione standard campionaria del campione $(x_i)_{i \in \{1, \dots, n\}} \subset \mathbb{R}$.

Se $s > 0$, allora per ogni $\alpha > 0$:

$$a) \quad \frac{\#\{i \in \{1, \dots, n\} : |x_i - \bar{X}| < \alpha \cdot s\}}{n} \geq 1 - \frac{(n-1)}{n\alpha^2} > 1 - \frac{1}{\alpha^2},$$

$$b) \quad \frac{\#\{i \in \{1, \dots, n\} : |x_i - \bar{X}| \leq \alpha \cdot s\}}{n} > 1 - \frac{(n-1)}{n\alpha^2}.$$

unilaterale

$$c) \quad \frac{\#\{i \in \{1, \dots, n\} : x_i - \bar{X} \geq \alpha \cdot s\}}{n} < \frac{1}{1 + \alpha^2}.$$

Se $s = 0$, allora $x_i = \bar{X}$ per ogni $i \in \{1, \dots, n\}$.

Punti a) e b) interessanti per $\alpha \geq 1$ (almeno $\alpha > \sqrt{\frac{n-1}{n}}$).



Osservazioni:

- 1) La disuguaglianza di Chebyshev (nelle ~~varie~~ ^{varie} versioni) fornisce ~~una~~ stime per il caso peggiore, ma è valida per un qualsiasi campione di dati.
- 2) Le stime di sopra possono essere migliorate quando si hanno più informazioni sulla distribuzione dei dati (in particolare, per dati "approssimativamente normali").
- 3) Scegliendo $\alpha \in \{2, 3, 5\}$ nella 2) otteniamo:

$$\frac{\#\{i \in \{1, \dots, n\} : |x_i - \bar{x}| < 2s\}}{n} > \frac{3}{4} = 75\%$$

$$\frac{\#\{i \in \{1, \dots, n\} : |x_i - \bar{x}| < 3s\}}{n} > \frac{8}{9} \approx 89\%$$

$$\frac{\#\{i \in \{1, \dots, n\} : |x_i - \bar{x}| < 5s\}}{n} > \frac{24}{25} = 96\%$$

Deviazioni di più di ~~quattro~~ ^{cinque} volte (2 deviazione standard) sono quindi ~~molto~~ ^{abbastanza} rare, persino nel caso peggiore.

Spesso i dati non sono univariati,
ma bi- o multivariati, cioè elementi di \mathbb{R}^d
($d=2$ nel caso bivariato), oppure con una struttura
ancora più generale (2 valori in uno spazio metrico).

In questo caso, ~~tra~~ è di interesse quantificare
una possibile dipendenza tra le componenti (marginali)
dei dati. Una statistica fondamentale per
misurare la dipendenza tra due marginali di
un campione di dati multivariati è la
correlazione campionaria, a sua volta definita
in termini delle deviazioni standard delle componenti
e della covarianza campionaria.

Sia $(X^{(l)})_{l \in \{1, \dots, n\}} \subset \mathbb{R}^d$ un campione di
numerosità n di dati d -variati. Notazione: $X^{(l)} = (X_1^{(l)}, \dots, X_d^{(l)})$

Def.: La covarianza campionaria tra componenti
 i -esima e j -esima del campione ($i, j \in \{1, \dots, d\}$)

$$\text{è data da } \text{cov}_{i,j} = \frac{1}{n-1} \left(\sum_{l=1}^n (X_i^{(l)} - \bar{X}_i)(X_j^{(l)} - \bar{X}_j) \right),$$

dove \bar{X}_i, \bar{X}_j sono le relative medie campionarie marginali

Osservazione:

La matrice $d \times d$ data da

$$\text{cov} = (\text{cov}_{ij})_{i,j \in \{1, \dots, d\}} \quad \text{si dice}$$

matrice di covarianza. Essa è

simmetrica (cioè $\text{cov}_{ij} = \text{cov}_{ji} \quad \forall i, j$).

semidefinita positiva (cioè $x^T \text{cov} x \geq 0 \quad \forall x \in \mathbb{R}^d$).

Le entrate sulla diagonale sono non-negative, infatti

$\forall i \in \{1, \dots, d\}$: cov_{ii} è la varianza della componente i -esima.

Def.: La correlazione campionaria tra le componenti i e j del campione

è data da

$$\text{corr}_{i,j} = \frac{\text{cov}_{i,j}}{s_i \cdot s_j},$$

dove s_i, s_j sono le deviazioni standard campionarie delle componenti i e j .

Osservazioni:

1) Per definizione,

$$\text{corr}_{i,j} = \frac{\sum_{l=1}^d (x_i^{(l)} - \bar{x}_i) \cdot (x_j^{(l)} - \bar{x}_j)}{\sqrt{\left(\sum_{l=1}^d (x_i^{(l)} - \bar{x}_i)^2\right) \cdot \left(\sum_{l=1}^d (x_j^{(l)} - \bar{x}_j)^2\right)}}$$

Nota: Stesso risultato se nelle definizioni di varianza e covarianza si usa il prefattore

$$\frac{1}{n} \text{ al posto di } \frac{1}{n-1}.$$

2) $\text{corr}_{i,j} \in [-1, 1]$. Inoltre,

<div style="border-left: 1px solid black; padding-left: 5px; margin-left: 5px;"> <u>corr</u> <u>misura per</u> <u>il grado</u> <u>di dipendenza</u> <u>lineare-affine</u> </div>	}	$\text{corr}_{i,j} = 1$ se e solo se $\exists a, b \in \mathbb{R}$ con <u>$b \geq 0$</u> :	$x_j^{(l)} = a + b \cdot x_i^{(l)} \quad \forall l \in \{1, \dots, n\}.$
		$\text{corr}_{i,j} = -1$ se e solo se $\exists a, b \in \mathbb{R}$ con <u>$b \leq 0$</u> :	$x_j^{(l)} = a + b \cdot x_i^{(l)} \quad \forall l \in \{1, \dots, n\}.$

Domanda: A cosa serve la probabilità
nel contesto della statistica?

Spesso: Dati di interesse riguardano una "popolazione"
(di individui, oggetti, "entità" generali) molto più
numerosa della numerosità del campione
a disposizione.

[Esempio: ci interessano le età degli ~~individui~~ "italiani" ~~di una~~]

Per poter trarre conclusioni dai dati del campione
sulle ~~distin~~ corrispondenti grandezze per la popolazione
bisogna scegliere un campione rappresentativo.

Metodo più affidabile: scegliere "a caso" un
sottoinsieme "abbastanza numeroso" della popolazione

~> dati del campione diventano risultato
di un "esperimento elettorale"

~> serve un modello matematico per
descrivere esperimenti elettorali.

Altri motivi per introdurre un modello matematico per "il caso":

- 1) Dati "perturbati", in particolare errori di osservazione / misurazione:

$$Y = X + \text{"rumore"}$$

\uparrow valore "vero" \uparrow valore rilevato

- 2) Impossibilità di prevedere grandezze di interesse con certezza; per mancanza di informazioni o per la natura del fenomeno
(2d esempio, previsioni meteo, giochi d'azzardo, finanza)

Alcuni modelli semplici che incontreremo:

- 1) Lancio di monete, dadi come ingrediente principale dell'esperimento (\nearrow Kids University)
- 2) dinamiche stocastiche: passeggiate stocastiche, catene di Markov