

TORNIAMO AD ESAMINARE IL LANCIO DI UN DADO. AVEMMO DESCRITTO LA SITUAZIONE TRAMITE LO SPAZIO DI PROBABILITÀ (Ω, \mathcal{A}, P) , DOVE $\Omega = \{1, \dots, 6\}$, $\mathcal{A} = \{\text{FAMIGLIA DEI SOTTOINSIEMI DI } \Omega\}$, $P: \mathcal{A} \rightarrow [0, 1]$ DEFINITA DA $P(A) = \frac{|A|}{|\Omega|}$ O, EQUIVAMENTE ESSENDO Ω FINITO, $P(\{i\}) = \frac{1}{|\Omega|}$ (CIOÈ BASTA DEFINIRE P SUGLI EVENTI ELEMENTARI).

LA NOSTRA SCELTA DI PORRE $P(\{1\}) = \dots = P(\{6\}) = \frac{1}{6}$ È DETTATA DA CONSIDERAZIONI RELATIVE AL SINGHERIA DEL DADO.

SE IL DADO FOSSE SBILANCIATO POTREMMO INVECE DECIDERE DI PORRE $P(\{1\}) = \frac{1}{4}$ E $P(\{2\}) = \dots = P(\{6\}) = \frac{3}{20}$

PIÙ IN GENERALE, POSTO $\theta \in [1, \infty)$, POTREMMO PORRE

$$P_\theta(\{1\}) = \frac{1}{\theta} \quad \text{E} \quad P_\theta(\{2\}) = \dots = P_\theta(\{6\}) = \left(1 - \frac{1}{\theta}\right) \frac{1}{5}$$

QUINDI A SECONDA DELLA SCELTA θ ABBIAMO UNA SPECIFICA PROBABILITÀ $P_\theta: \mathcal{A} \rightarrow [0, 1]$ E UN RELATIVO SPAZIO DI PROBABILITÀ $(\Omega, \mathcal{A}, P_\theta)$.

A SECONDA DI COME È FATTO IL DADO, UNO DEI VALORI DEL PARAMETRO θ FORNIRÀ UNO SPAZIO CHE DESCRIVE NELLO LA SITUAZIONE. TALE VALORE VA TROVATO TRAMITE OSSERVAZIONI.

DEFINIZIONE CHIAMEREMO MODELLO STATISTICO (PARAMETRICO) UNA FAMIGLIA DI SPAZI DI PROBABILITÀ $(\Omega, \mathcal{A}, P_\theta)$ INDICIZZATI DA $\theta \in \Theta$, DOVE Θ È UNA OPPORTUNA FAMIGLIA DI PARAMETRI. SI NOTI CHE L'INSIEME DEGLI EVENTI ELEMENTARI Ω E LA σ -ALGEBRA \mathcal{A} SONO SEMPRE GLI STESSI, È LA PROBABILITÀ P_θ CHE VARIA CON IL PARAMETRO.

SUPPONIAMO ORA DI AVERE UNA VARIABILE ALEATORIA DISCRETA $X: \Omega \rightarrow \mathbb{R}$. SIA q LA SUA DENSITÀ. PER DEFINIZIONE

$$q(x) = P\{X=x\} \quad \forall x \in \mathbb{R}$$

QUINDI SE VARIANO P IN FUNZIONE DI UN PARAMETRO θ , ANCHE LA VARIABILE ANCHE LA DENSITÀ DI X :

$$q(x) = P_\theta\{X=x\}$$

SIMILI CONSIDERAZIONI VALGONO NEL CASO ASSOLUTAMENTE CONTINUO. LA DENSITÀ $f: \mathbb{R} \rightarrow \mathbb{R}$ È TALE CHE

$$\int_{-\infty}^x f(t) dt = P\{X \leq x\} \quad \forall x \in \mathbb{R}$$

QUINDI SE VARIANO P IN FUNZIONE DI θ , CI TROVEREMO UNA DENSITÀ $f = f(t, \theta)$ TALE CHE

$$\int_{-\infty}^x f(t, \theta) dt = P_\theta\{X \leq x\}$$

Molto spesso ci siano trovati ad operare con una v.a. X senza conoscere esplicitamente lo spazio di probabilità (Ω, \mathcal{A}, P) sottostante, ma conoscendo solo la densità di X .

Similmente, ci troviamo spesso nella situazione non di avere un modello statistico $\{(\Omega, \mathcal{A}, P_\theta) : \theta \in \Theta\}$ e una v.a. $X: \Omega \rightarrow \mathbb{R}$ con densità che varia con θ , ma direttamente una famiglia di v.a. (discrete o continue) con una densità che varia a seconda di un parametro θ :

$$\{q(x, \theta) : \theta \in \Theta\}, \quad \{f(x, \theta) : \theta \in \Theta\}$$

Pensaremo a questa famiglia di v.a. come ad un'unica v.a. definita su uno spazio di probabilità in cui la misura di probabilità P varia con θ .

Questa è la situazione tipica in cui sappiamo che il fenomeno è modellato da una v.a. di un dato tipo ma non ne conosciamo i parametri. Per esempio

- $\{E(\lambda) : \lambda \in (0, +\infty)\}$ è la famiglia delle v.a. esponenziali.

Il parametro λ gioca il ruolo di θ e $\Theta = (0, +\infty)$

- $\{N(\mu, \sigma^2) : \mu \in \mathbb{R} \text{ e } \sigma \in (0, +\infty)\}$ è la famiglia delle v.a.

Gaussiane. Qui $\theta = (\mu, \sigma) \in \Theta = \mathbb{R} \times (0, +\infty)$.

CON LE NOZIONI DI PROBABILITÀ STUDIATE FINORA ABBIAMO POTUTO RISPONDERE AD UNA SERIE (RELATIVAMENTE ESTESA) DI DOMANDE DEL TIPO

- QUAL È LA PROBABILITÀ DELL'EVENUTO...
- QUAL È LA PROBABILITÀ CHE LA V.A. ...
- QUAL È LA MEDIA DELLA V.A. ...

ECC. NEL RISPONDERE A QUESTE DOMANDE ABBIAMO GIOCATO CONOSCENDO UN PUNTO FERMO: LA PROBABILITÀ p .

IN STATISTICA LE DOMANDE SONO SIMILI, MA ABBIAMO p_θ INVECE DI p . PER STIMARE p_θ CI SERVIRANNO DI OSSERVAZIONI/DATI

x_1, \dots, x_m OTTENUTI TRAMITE ESPERIMENTI, E LE RISPOSTE ALLE DOMANDE SOPRA SARANNO DEL TIPO $\psi(\theta)$, CIÒE DIPENDENTI DA θ . LA COSTRUZIONE DEL MODELLO STATISTICO $(\Omega, \mathcal{F}, p_\theta)$ È DELICATA E NON BANALE. PER ESEMPIO, NEL COSTRUIRE QUELLO CON IL DADO, SIAMO SICURI CHE LA SCELTA DI p_θ CHE ABBIAMO FATTO SIA VALIDA?

LE OSSERVAZIONI x_1, \dots, x_m VIENGONO DI SOLITO PENSATE COME VALORI ASSUMTI DA UNA FAMIGLIA X_1, \dots, X_m DI V.A. LA CUI LEGGE CONGIUNTA DIPENDE DAL PARAMETRO θ

Un uso frequente è quello in cui X_1, \dots, X_n sono tra loro indipendenti e con la stessa distribuzione: è la formalizzazione matematica del caso in cui lo sperimentatore decide di ripetere n volte l'esperimento, in condizioni di indipendenza.

Diamo dunque la seguente definizione.

Definizione Sia $\{(\Omega, \mathcal{F}, P_\theta) : \theta \in \Theta\}$ un modello statistico. Chiameremo campione ogni successione $\{X_n\}_{n \in \mathbb{N}}$ di v.a. su Ω che per ogni probabilità P_θ , $\theta \in \Theta$, sono i.i.d. (indipendenti ed identicamente distribuite).

L' n -pla (X_1, \dots, X_n) è detta campione di ampiezza n .

Nel caso discreto [risp. continuo] ognuna delle X_n avrà densità $\{p(x, \theta) : \theta \in \Theta\}$ [risp. $\{f(x, \theta) : \theta \in \Theta\}$].

Lo sperimentatore userà le v.a. X_1, \dots, X_n per stimare il parametro incognito θ , o per dare direttamente una risposta alla domanda $\psi(\theta)$. Per questo userà una "elaborazione" $H = h(X_1, \dots, X_n)$ di queste v.a.

La situazione ideale sarebbe che valesse l'uguaglianza

$$H(\omega) = \psi(\theta) \quad \forall \omega \in \Omega \text{ e } \forall \theta \in \Theta$$

cioè $h(x_1, \dots, x_m) = \psi(\theta)$ per ogni m -pla di valori osservati (x_1, \dots, x_m) . Troppo ottimistico, è più sensato chiedere che valga in media:

$$\mathbb{E}_\theta[H] = \psi(\theta) \quad \forall \theta \in \Theta$$

Con il simbolo \mathbb{E}_θ stiamo indicando il valore medio rispetto alla probabilità P_θ .

Diamo dunque le seguenti definizioni:

Definizioni Sia $\{(\Omega, \mathcal{A}, P_\theta) : \theta \in \Theta\}$ un modello statistico e sia $\{X_m\}_{m \in \mathbb{N}}$ un campione su di esso.

Una successione di variabili aleatorie $\{H_m\}_{m \in \mathbb{N}}$ della forma

$$H_m = h_m(X_1, \dots, X_m)$$

dove le $h_m: \mathbb{R}^m \rightarrow \mathbb{R}$ sono funzioni regolari, viene detta statistica campionaria (o stimatori).

Per regolari intendo che trasformano v.a. in v.a.

La singola variabile H_m è detta statistica campionaria basata su un campione di taglia m .

Una statistica campionaria $\{H_m\}_{m \in \mathbb{N}}$ si dice stimatori corretto per una funzione $\psi: \Theta \rightarrow \mathbb{R}$ se le H_m hanno valore medio finito rispetto ogni P_θ e

$$E_{\theta}[H_n] = \psi(\theta) \quad \forall \theta \in \Theta \quad \text{e} \quad \forall n \in \mathbb{N}$$

FISSATA LA STATISTICA $H_n = h_n(X_1, \dots, X_n)$, QUANDO PRODUCIAMO UN DETERMINATO CAMPIONE (x_1, \dots, x_n) , IL VALORE $h_n(x_1, \dots, x_n)$ È DETTO STIMA PER $\psi(\theta)$.

NOTA NELLE DEFINIZIONI USIAMO SUCCESSIONI $\{X_n\}_{n \in \mathbb{N}}$ INVECE CHE n -PLE SOLO PER INDICARE CHE I CAMPIONI POSSONO ESSERE ARBITRARIAMENTE ESTESI.

Ricapitolano

- Lo stimatore H_n è una variabile aleatoria
- La stina $h_n(x_1, \dots, x_n)$ è un numero reale
- Le v.a. X_1, \dots, X_n sono i possibili output degli esperimenti. Precisamente, X_k rappresenta il risultato del k -esimo esperimento, o il k -esimo elemento del campione.

Se lo stato è $\omega \in \Omega$, allora l'osservazione restituirà $x_1 = X_1(\omega), \dots, x_n = X_n(\omega)$. E equivalentemente, se osserviamo x_1, \dots, x_n significa che si è verificato l'evento

$$\{X_k = x_k, k=1, \dots, n\}$$

- Se a posteriori decidiamo di calcolare $h_n(x_1, \dots, x_n)$ per stimare $\psi(\theta)$, il modello di questa stina è a priori appunto la statistica campionaria $H_n = h_n(X_1, \dots, X_n)$. Tale statistica ha la proprietà che in media restituisce $\psi(\theta)$.

ESEMPIO CONSIDERIAMO UN MODELLO STATISTICO IN CUI
 $q(x, \lambda)$ È LA DENSITÀ DI TIPO POISSON $\mathcal{P}(\lambda)$. STIAMO PREN-
DENDO COME PARAMETRO $\theta = \lambda \in \Theta = (0, +\infty)$.

PONIAMO $h_m(x_1, \dots, x_m) = \frac{1}{m} \sum_{k=1}^m x_k$, CIOÈ PRENDIAMO COME h_m
LA MEDIA CAMPIONARIA SU m DATI. ABBIAMO DUNQUE CHE
LA STATISTICA CAMPIONARIA $H_m = h_m(X_1, \dots, X_m)$ COINCIDE CON
LA MEDIA CAMPIONARIA \bar{X}_m .

PRENDENDO COME ψ LA FUNZIONE IDENTITÀ, CIOÈ $\psi(\lambda) = \lambda$,

ABBIAMO

$$E_{\lambda}[H_m] = \frac{1}{m} \sum_{k=1}^m E_{\lambda}[X_k] = \lambda = \psi(\lambda)$$

QUINDI $\{H_m\}$ È UNO STIMATORE CORRETTO PER ψ , CIOÈ
PER IL PARAMETRO θ STESSO.

ESEMPIO CONSIDERIAMO UN MODELLO STATISTICO IN CUI
 $f(x, \theta)$ È LA DENSITÀ DI TIPO GAUSSIANO $\mathcal{N}(\mu, \sigma^2)$. STIAMO
PRENDENDO COME PARAMETRO $\theta = (\mu, \sigma^2) \in \Theta = \mathbb{R} \times (0, +\infty)$.

DATO UN CAMPIONE $\{X_m\}_{m \in \mathbb{N}}$ i.i.d. CON $X_m \sim \mathcal{N}(\mu, \sigma^2)$,

CONSIDERIAMO LE STATISTICHE CAMPIONARIE

$$\bar{X}_m = \frac{1}{m} \sum_{k=1}^m X_k \quad \text{E} \quad S_m^2 = \frac{1}{m-1} \sum_{k=1}^m (X_k - \bar{X}_m)^2$$

RICORDO CHE S_m^2 È LA V.A. VARIANZA CAMPIONARIA DELLE V.A.

X_1, \dots, X_n . In questo caso

$$h_n(x_1, \dots, x_n) = \frac{1}{n-1} \sum_{k=1}^n \left(x_k - \frac{1}{n} \sum_{j=1}^n x_j \right)^2$$

cioè prendiamo come h_n quella che è chiamata

VARIANZA CAMPIONARIA SU n DATI (NOTARE IL FATTORE

$\frac{1}{n-1}$ CHE DIFFERISCE DAL FATTORE $\frac{1}{n}$ USATO NELLA DEFINIZIONE DI VARIANZA). CONSIDERIAMO ORA LE DUE FUNZIONI

$$\psi_1(\mu, \sigma^2) = \mu \quad \text{e} \quad \psi_2(\mu, \sigma^2) = \sigma^2$$

cioè LE PROIEZIONI SULLA PRIMA E SULLA SECONDA COORDINATA.

Abbiamo $E_{\theta}[\bar{X}_n] = \mu = \psi_1(\theta)$ e quindi \bar{X}_n

è uno stimatore corretto per ψ_1 . Inoltre

$$(n-1)E_{\theta}[S_n^2] = \sum_{k=1}^n E_{\theta}[(X_k - \bar{X}_n)^2]$$

USANO
RIPETUTAMENTE
LA LINEARITÀ
DEL VALORE
MEDIO

$$\begin{aligned} &= \sum_{k=1}^n \left(E_{\theta}[X_k^2] + E_{\theta}[\bar{X}_n^2] - 2E[X_k \bar{X}_n] \right) \\ &= \left(\sum_{k=1}^n E_{\theta}[X_k^2] \right) + n E_{\theta}[\bar{X}_n^2] - 2E\left[\sum_{k=1}^n X_k \bar{X}_n\right] \\ &= \left(\sum_{k=1}^n E_{\theta}[X_k^2] \right) - n E_{\theta}[\bar{X}_n^2] \end{aligned}$$

$$= n(\mu^2 + \sigma^2) - n\left(\mu^2 + \frac{\sigma^2}{n}\right) = (n-1)\sigma^2$$

* RICORDANDO CHE $E[X^2] = \text{VAR } X + E[X]^2$ E CHE $\text{VAR } \bar{X}_n = \frac{\sigma^2}{n}$

Dunque S_n^2 è uno stimatore corretto per ψ_2 .

OSSERVAZIONE NELL'ESEMPIO IL FATTO DI OPERARE CON LA V.A. GAUSSIANA NON È STATO SFRUTTATO: ABBIAMO SOLO USATO IL FATTO CHE LE X_m HANNO MEDIA μ E VARIANZA σ^2 . NE SEGUE IL SEGUENTE RISULTATO GENERALE.

LEMMA DATO UN MODELLO STATISTICO $\{(\Omega, \mathcal{F}, P_\theta) : \theta \in \Theta\}$ E UN CAMPIONE $\{X_m\}_{m \in \mathbb{N}}$ CON X_m AVENTE DISTRIBUZIONE SECONDO FINITO RISPETTO OGNI P_θ , DEFINIAMO $\mu, \sigma^2 : \Theta \rightarrow \mathbb{R}$ TRAMITE

$$\mu(\theta) = E_\theta[X_m] \quad \text{e} \quad \sigma^2(\theta) = \text{Var}_\theta X_m \quad \forall \theta \in \Theta$$

NOTARE CHE LE DEFINIZIONI SONO BEN POSTE PERCHÉ IL LATO DESTRO NON DIPENDE DA m : LE X_m HANNO TUTTE LA STESSA DISTRIBUZIONE RISPETTO OGNI DATA P_θ .

ALLORA GLI STATISTORI \bar{X}_m E S_m^2 SONO CORRETTI PER $\mu = \mu(\theta)$ E $\sigma^2 = \sigma^2(\theta)$ RISPETTIVAMENTE.

IN UN DATO MODELLO STATISTICO PARAMETRICO, VIENGONO UTILIZZATI DIFFERENTI STATISTORI CORRETTI PER UNA ASSEGNATA FUNZIONE $\psi = \psi(\theta)$. ABBIAMO VISTO COME LA VARIANZA SIA UN INDICE DI QUANTO UNA V.A. SIA CONCENTRATA INTORNO

AL SUO VALORE MEDIO. QUINDI, SE UNO STIMATORE H_m È CORRETTO PER ψ , CIOÈ $E_\theta[H] = \psi(\theta)$, LA $VAR_\theta H_m$ MISURERÀ QUANTO I VALORI DI H_m SONO CONCENTRATI INTORNO ψ .

DEFINIZIONE CHIAMEREMO RISCHIO QUADRATICO MEDIO DELLO STIMATORE H_m LA FUNZIONE

$$R_{H_m}(\theta) = E_\theta[(H_m - \psi(\theta))^2]$$

SE H_m È CORRETTO OVVIAMENTE $R_{H_m}(\theta) = VAR_\theta H_m$

LA CORRETTIEZZA DI UN ESTIMATORE È DATA SU UN NUMERO FINITO n DI OSSERVAZIONI. PUÒ ESSERE UTILE PERÒ ANCHE CONOSCERNE IL COMPORTAMENTO ASINTOTICO PER $n \rightarrow +\infty$.

DAL PUNTO DI VISTA PRATICO, POTREBBE SUGGERIRCI DI AUMENTARE IL NUMERO DI OSSERVAZIONI PER AVERE STIME MIGLIORI.

DEFINIZIONE UNA SUCCESSIONE $\{H_n\}_{n \in \mathbb{N}}$ DI STIMATORI DI $\psi = \psi(\theta)$ È DETTA CONSISTENTE SE CONVERGIE IN PROBABILITÀ RISPETTO OGNI P_θ A $\psi(\theta)$, CIOÈ

$$\lim_{n \rightarrow +\infty} P_\theta\{|H_n - \psi(\theta)| > \varepsilon\} = 0 \quad \forall \varepsilon > 0 \quad \forall \theta \in \Theta$$

DAL PUNTO DI VISTA INTUITIVO, $\{H_n\}_{n \in \mathbb{N}}$ È CONSISTENTE PER $\psi(\theta)$ SE, PER n GRANDE, H_n È UNA FUNZIONE DELLE OSSERVAZIONI CHE ASSUME VALORI PROSSIMI A $\psi(\theta)$ CON GRANDE PROBABILITÀ RISPETTO P_θ .

PER LA LEGGE DEI GRANDI NUMERI, LA MEDIA CAMPIONARIA RISULTA ESSERE UNO STIMATORE CONSISTENTE DEL VALORE MEDIO. SE INOLTRE $E_\theta[H_n^2]$ È FINITO $\forall n \in \mathbb{N}$ E $\forall \theta \in \Theta$, ALLORA ANCHE LA VARIANZA CAMPIONARIA RISULTA CONSISTENTE.

Ricapitolano II

LA MEDIA CAMPIONARIA \bar{X}_m E LA VARIANZA CAMPIONARIA S_m^2 SONO STIMATORI CORRETTI E CONSISTENTI PER LA MEDIA μ E LA VARIANZA σ^2 RISPETTIVAMENTE. MENTRE LA MEDIA CAMPIONARIA ERA LARGAMENTE ATTESA, PER LA VARIANZA CAMPIONARIA IL RISULTATO ERA PENO ASPETTATO. NOTARE CHE ORA È MOTIVATA LA DEFINIZIONE

$$S_m^2 = \frac{1}{m-1} \sum_{k=1}^m (X_k - \bar{X}_m)^2 \text{ INVECE DI } \frac{1}{m} \sum_{k=1}^m (X_k - \bar{X}_m)^2$$

ALTRIMENTI NON AVREMMO UNO STIMATORE CORRETTO!

NOTARE CHE S_m^2 PERMETTE DI STIMARE σ^2 SENZA CONOSCERE μ . A TAL PROPOSITO SI RIVEDA L'ESEMPIO SUI KIWI, RILEGGENDOLO IN QUEST'OTTICA.

DI CONTRO, SE μ È NOTA, ALLORA CONE STIMATORE POSSIAMO USARE

$$H_m = \frac{1}{m} \sum_{k=1}^m (X_k - \mu)^2$$

INFATTI

$$\begin{aligned} m E_{\theta} [H_m] &= \sum_{k=1}^m E_{\theta} [(X_k - \mu)^2] = \sum_{k=1}^m \left(E_{\theta} [X_k^2] + \mu^2 - 2\mu E_{\theta} [X_k] \right) \\ &= \sum_{k=1}^m \left(E_{\theta} [X_k^2] - \mu^2 \right) = \sum_{k=1}^m \sigma^2 = m \sigma^2 \end{aligned}$$