

Statistica inferenziale

In un problema di statistica abbiamo un esperimento aleatorio al quale sappiamo associare uno spazio probabilizzabile, ma non sappiamo scegliere su di esso una misura di probabilità P : riusciamo solo a stabilire che P deve far parte di una certa famiglia \mathcal{P} , più o meno grande, di misure di probabilità. Il compito dello statistico è quello di raccogliere informazioni sull'esperimento aleatorio, favorendo la scelta della misura di probabilità più adeguata per descrivere l'esperimento.

Definizione Sia (Ω, \mathcal{A}) uno spazio probabilizzabile e sia \mathcal{P} una famiglia di misure di probabilità su (Ω, \mathcal{A}) . La terna $(\Omega, \mathcal{A}, \mathcal{P})$ è detta modello statistico. Esso si dice parametrico se si ha $\mathcal{P} = \{P_\theta, \theta \in D\}$ con $D \subseteq \mathbb{R}^k$, non parametrico altrimenti.

Noi ci limiteremo all'analisi di modelli statistici parametrici. In definitiva, compito preliminare per lo studio di un problema statistico è quello di associare all'esperimento aleatorio un modello statistico.

Esempi (1) (controllo di qualità). Una popolazione è composta da individui di tipo A e di tipo B. Non conoscendo il rapporto effettivo tra il n° di individui di tipo A e l'intera popolazione, si sceglie un campione di N individui: le osservazioni di questo

(419)

esperimento sono allora (come abbiamo già visto) rappresentate da v.a. X_1, \dots, X_N , a valori in $\{0, 1\}$, definite su un opportuno spazio probabilizzabile (Ω, \mathcal{A}) : ovviamente, $X_i = 1$ se l' i -esimo individuo del campione è di tipo A , $X_i = 0$ altrimenti. Non siamo in grado però di scegliere una misura di probabilità P su (Ω, \mathcal{A}) , pur sapendo che le v.a. X_i sono tutte indipendenti e bernoulliane con un certo parametro $\theta \in [0, 1]$. Dunque, ad uno un modello statistico (Ω, \mathcal{A}, P) con $P = \{P^\theta, \theta \in [0, 1]\}$, ove, per ogni valore del parametro sconosciuto θ , le v.a. X_i hanno legge di Bernoulli: $B(\theta)$ secondo la misura P^θ .

(2) (misura di una grandezza fisica). Per effettuare una misura di una certa grandezza con un determinato strumento, si esegue un certo numero N di misurazioni, il cui risultato, secondo il teorema limite centrale, sarà descritto da v.a. X_1, \dots, X_N indipendenti e gaussiane, pur non essendo note a priori né la speranza μ , né la varianza σ^2 comuni a tutte le X_i . Si potrà quindi scegliere, su un opportuno spazio probabilizzabile (Ω, \mathcal{A}) , una famiglia di misure di probabilità $\{P^{\mu, \sigma^2}; (\mu, \sigma^2) \in \mathbb{R} \times \mathbb{R}^+\}$, in modo tale che, per ogni (μ, σ^2) , le X_i abbiano legge normale $N(\mu, \sigma^2)$ secondo la misura di probabilità P^{μ, σ^2} .

Definizione Dato un modello statistico $(\Omega, \mathcal{A}, \mathcal{P})$, una 420
v.a. X , definita su (Ω, \mathcal{A}) , è detta una statistica -

Definizione Sia $(\Omega, \mathcal{A}, \{P^\theta: \theta \in D\})$ un modello statistico parametrico. Una successione $\{X_n\}$ di statistiche si dice indipendente se le X_n sono indipendenti rispetto a ciascuna delle probabilità P^θ . Una statistica X è integrabile se X è integrabile rispetto a ciascuna delle probabilità P^θ . Scriveremo in tal caso

$$E^\theta[X], \quad \text{Var}^\theta[X]$$

per indicare la speranza e la varianza di X rispetto alla probabilità P^θ .

Definizione Sia $(\Omega, \mathcal{A}, \{P^\theta: \theta \in D\})$ un modello statistico parametrico, e sia $\{L(\theta), \theta \in D\}$ una famiglia di leggi di probabilità. Una sequenza finita di statistiche X_1, \dots, X_N definite su $(\Omega, \mathcal{A}, \{P^\theta: \theta \in D\})$ si dice campione statistico di taglia N , estratto da una popolazione di legge $L(\theta)$, se le statistiche sono indipendenti e tutte dotate di legge $L(\theta)$ secondo la probabilità P^θ .

Osserviamo che, fissata una famiglia di leggi di probabilità $\{L(\theta), \theta \in D\}$, si può sempre costruire un modello statistico e, su

di esso, un campione statistico di taglia N estratto da una popolazione di legge $L(\theta)$. Esso si chiama modello statistico campionario di taglia N . (421)

Nel seguito ci occuperemo di 3 problemi generali, che illustrano i principali metodi di inferenza della statistica.

1. il problema dello stima puntuale;
2. il problema dello stima insiemistica;
3. il problema dei test d'ipotesi.

Il problema 1 consiste nella scelta di uno stimate, cioè di una applicazione $T: \Omega \rightarrow D$, che rappresenta la strategia seguente: ci si impegna, qualunque sia la realizzazione w dell'esperimento, ad attribuire al parametro sconosciuto θ il valore $T(w)$. Il problema dello stima consistere dunque nella scegliere lo stimate T in modo da minimizzare l'errore, ossia minimizzando certe quantità legate a T , espresse per mezzo delle probabilità P_θ .

Il problema 2 consiste nella scelta di un'applicazione $S: \Omega \rightarrow \mathcal{P}(D)$, che rappresenta la strategia seguente: ci si impegna, qualunque sia la realizzazione w dell'esperimento, a stimare il valore di θ come appartenente all'insieme $S(w)$ (insieme di fiducia).

Il problema 3 consiste nella scelta di un test, che permetta di verificare o di confutare una fissata ipotesi sul parametro sconosciuto θ , che stabilisca che esso appartiene ad un dato sottoinsieme $D_0 \subset D$. Il test consiste nella costruzione di una partizione $\{B, B^c\}$ di Ω , e rappresenta la strategia seguente: ci si impegna, qualunque sia la realizzazione w dell'esperimento, a rifiutare l'ipotesi se $w \in B$, e ad accettarla se $w \in B^c$.

In definitiva, lo statistico agisce seguendo una regola di decisione che tende a minimizzare certe conseguenze le quali, ragionevolmente, appaiono nocive; ma la regola di decisione deve essere a priori, vale a dire che ci si impegna a seguirle prima di compiere l'esperimento e qualunque sia il suo risultato w . Dunque l'azione dello statistico è determinata da w , ma attraverso una regola definita a priori.

Teoria dello stimo

Introduciamo la nozione di stimatore. Sia $(\Omega, \mathcal{A}, \{P^\theta\}_{\theta \in D})$ un modello statistico e sia assegnato su di esso un campione statistico X_1, \dots, X_N di taglia N . Sia poi $\psi: D \rightarrow I$ una funzione del parametro θ , a valori nell'insieme $I \subseteq \mathbb{R}$.

Definizione Uno stimatore della funzione $\psi(\theta)$ è una statistica T della forma

$$T = t(X_1, \dots, X_N)$$

con t funzione da \mathbb{R}^N in I .

Intuitivamente, assegnare uno stimatore $T = t(X_1, \dots, X_N)$ di $\psi(\theta)$ significa fissare la seguente regola: se i dati raccolti dall'osservazione di un risultato ω sono

$$(x_1, \dots, x_N) = (X_1(\omega), \dots, X_N(\omega)),$$

si stimerà la quantità sconosciuta $\psi(\theta)$ con il numero $t(x_1, \dots, x_N)$, detto appunto lo stimo di $\psi(\theta)$.

Osservazione 1: lo stimatore T dipende dalla taglia N del campione: perciò, di solito, si costruisce una successione $\{T_n\}_{n \in \mathbb{N}}$ dove, per ogni n , la statistica T_n è uno stimatore di $\psi(\theta)$ su

tramite un campione di taglia n .

(424)

Osservazione 2 Uno stimatore, essendo una v.a., non assumerà mai il valore $\psi(\theta)$, ma naturalmente si spera che prenda valori non troppo distanti da esso. Poiché qualunque funzione delle osservazioni è uno stimatore, sarà necessario fissare qualche criterio per stabilire quali di esse sono "buoni" stimatori e quali no. Noi ci limiteremo a studiare buoni stimatori solo in casi molto semplici, legati alla media e alla varianza della legge del campione.

Definizione La statistica T è uno stimatore corretto (o non distorto) del parametro $\psi(\theta)$ se risulta

$$E^\theta [T] = \psi(\theta) \quad \forall \theta \in D.$$

In caso contrario, T si dice stimatore distorto.

Se lo stimatore si discosta da $\psi(\theta)$, la sostituzione del valore $\psi(\theta)$ con T comporta un errore, e quindi un costo, che si esprime mediante una funzione positiva $C(\theta, a)$, che misura la perdita proveniente dal sostituire $\psi(\theta)$ con il valore a . Dunque, se $T(w)$ si discosta da $\psi(\theta)$, il costo sarà la v.a.

$$w \mapsto C(\theta, T(w)).$$

(6.25)

Definizione Il rischio dello stimatore T è il suo costo medio, ossia la funzione

$$R_T(\theta) = E^\theta [C(\theta, T)], \quad \theta \in D.$$

Generalmente, si sceglie come costo la funzione quadratica

$$C(\theta, a) = \psi(\theta) |a - \theta|^2,$$

della appunto costo quadratico. Con questa scelta, il rischio di T è la funzione

$$R_T(\theta) = E^\theta [(\psi(\theta) - T)^2],$$

della rischio quadratico.

Si noti che se T è uno stimatore corretto, allora

$$\begin{aligned} R_T(\theta) &= E^\theta [(\psi(\theta) - T)^2] = \psi(\theta)^2 - 2\psi(\theta) E^\theta [T] + E^\theta [T^2] = \\ &= E^\theta [T^2] - E[\theta]^2 = \text{Var}^\theta [T]. \end{aligned}$$

Diciamo infine, se S, T sono stimatori di $\psi(\theta)$, che T è preferibile a S se risulta $R_T(\theta) \leq R_S(\theta)$; se \mathcal{T} è una famiglia di stimatori di $\psi(\theta)$, diremo che $T \in \mathcal{T}$ è uno stimatore ottimale (rispetto a \mathcal{T}) se T è preferibile a S per ogni $S \in \mathcal{T}$.

Passiamo ora dal generale al particolare, supponendo di avere un campione stocastico (X_1, \dots, X_N) di taglia N su un modello stocastico $(\Omega, \mathcal{A}, \mathbb{P})$. Supponiamo che il campione sia costituito da v.a. gaussiane con media μ e varianza σ^2 , entrambe finite ma sconosciute. Dunque

$$\mathbb{P} = \left\{ \begin{array}{l} \mu, \sigma \\ (\mu, \sigma) \in \mathbb{R} \times \mathbb{R}^+ \end{array} \right.$$

Definizione La media empirica, o media campionaria, è lo stimatore \bar{X} di μ già definito:

$$\bar{X} = \frac{X_1 + \dots + X_N}{N}$$

(dunque $\bar{X} = t(X_1, \dots, X_N)$ con $t(x) = \frac{1}{N} \sum_{i=1}^N x_i$).

Questo stimatore è corretto: infatti, per ogni $\theta = (\mu, \sigma)$ si ha

$$\begin{aligned} E^\theta[\bar{X}] &= E^{\mu, \sigma} \left[\frac{X_1 + \dots + X_N}{N} \right] = \frac{E^{\mu, \sigma}[X_1] + \dots + E^{\mu, \sigma}[X_N]}{N} = \\ &= \frac{N\mu}{N} = \mu. \end{aligned}$$

Inoltre, si osserva che (essendo in particolare le X_j indipendenti)

$$\text{Var}^\theta[\bar{X}] = \frac{\text{Var}^\theta[X_1] + \dots + \text{Var}^\theta[X_N]}{N^2} = \frac{N\sigma^2}{N^2} = \frac{\sigma^2}{N}.$$

Dunque la media empirica è uno stimatore corretto della media

del campione, ma la sua varianza è ridotta, rispetto (L27)
a quella del campione, di un fattore N .

Inoltre, per il teorema limite centrale, la statistica $Z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{N}}}$
è approssimativamente legge $\mathcal{N}(0,1)$ rispetto a ciascuna $P_{\mu, \sigma}$.

Dunque la legge di \bar{X} è approssimativamente $\mathcal{N}(\mu, \frac{\sigma^2}{N})$:

$$P(\bar{X} \leq c) = P(Z \leq \frac{\sqrt{N}(c-\mu)}{\sigma}) = \Phi\left(\frac{\sqrt{N}}{\sigma}(c-\mu)\right), \quad c \in \mathbb{R}.$$

Definizione La varianza empirica, o varianza campionaria,
è lo stimatore S^2 di σ^2 così definito:

$$S^2 = \frac{1}{N-1} \sum_{i=1}^N (X_i - \bar{X})^2,$$

mentre la statistica $S = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (X_i - \bar{X})^2}$ è detta deviazione standard empirica (o campionaria).

Perché dividere per $N-1$ e non per N ? Perché, così facendo,
 S^2 è uno stimatore corretto (e altrimenti no). Ciò si vede,
un po' faticosamente, così: partiamo dalla relazione

$$N\bar{X} = \sum_{i=1}^N X_i.$$

Allora

$$\sum_{i=1}^N (X_i - \bar{X})^2 = \sum_{i=1}^N X_i^2 - 2\bar{X} \sum_{i=1}^N X_i + N\bar{X}^2 = \sum_{i=1}^N X_i^2 - N\bar{X}^2.$$

Ne segue

$$S^2 = \frac{1}{N-1} \sum_{i=1}^N (X_i - \bar{X})^2 = \frac{1}{N-1} \left(\sum_{i=1}^N X_i^2 - N\bar{X}^2 \right),$$

(428)

da cui

$$(N-1)S^2 = \sum_{i=1}^N X_i^2 - N\bar{X}^2.$$

Prendiamo la speranza di entrambi i membri e ricordiamo che $E^{M,\sigma}[Y^2] = \text{Var}^{M,\sigma}[Y] + E^{M,\sigma}[Y]^2$ per ogni v.a. integrabile Y . Quindi

$$\begin{aligned} (N-1) E^{M,\sigma}[S^2] &= E^{M,\sigma} \left[\sum_{i=1}^N X_i^2 \right] - N E^{M,\sigma}[\bar{X}^2] = \\ &= N E^{M,\sigma}[X_1^2] - N E^{M,\sigma}[\bar{X}^2] = \\ &= N \text{Var}^{M,\sigma}[X_1] + N E^{M,\sigma}[X_1]^2 - \\ &\quad - N \cdot \text{Var}^{M,\sigma}[\bar{X}] - N E^{M,\sigma}[\bar{X}]^2 = \\ &= N \sigma^2 + N \mu^2 - N \frac{\sigma^2}{N} - N \mu^2 = (N-1) \sigma^2, \end{aligned}$$

da cui, finalmente, $E^{M,\sigma}[S^2] = \sigma^2$.

Il calcolo di $\text{Var}^{M,\sigma}[S^2]$ è assai intricato; lo vedremo dopo.

Qual è la legge della varianza empirica? C'è un teorema a questo proposito.

Teorema (di Cochran) Se \bar{X} e S^2 sono la media empirica e la varianza empirica di un campione statistico di taglia N e legge

normali $\mathcal{N}(\mu, \sigma^2)$, allora:

429

(i) \bar{X} e S^2 sono statistiche indipendenti;

(ii) \bar{X} ha legge normale $\mathcal{N}(\mu, \frac{\sigma^2}{N})$;

(iii) $W = \frac{N-1}{\sigma^2} S^2$ ha legge $\chi^2(N-1)$,

(iv) $T = \frac{\bar{X} - \mu}{S} \sqrt{N}$ ha legge $t(N-1)$.

dim. Omessa - \square .

Dunque possiamo fare stime probabilistiche sulle speranze e sulle varianze di campioni statistici con legge gaussiana.

Ma c'è di più: dal teorema di Cochran si sa che W ha legge $\chi^2(N-1)$, e dunque $\text{Var}[W] = 2(N-1)$; come abbiamo verificato qualche lezione fa. Dunque

$$\text{Var}[S^2] = \text{Var}\left[\frac{\sigma^2}{N-1} W\right] = \frac{\sigma^4}{(N-1)^2} \text{Var}[W] = \frac{2\sigma^4}{N-1},$$

e dunque, come accade per la media empirica, la varianza empirica ha varianza che tende a 0 per $N \rightarrow \infty$, ossia la

stima di σ^2 mediante S^2 è sempre più precisa all'aumentare della taglia del campione.

Esempio Per misurare una grandezza fisica si eseguono N misurazioni indipendenti ottenendo N risultati X_1, \dots, X_N . Poiché i dati ottenuti formano un campione di taglia N , per N grande la media e la varianza

empiriche sono stimatori corretti abbastanza precisi della media del campione e della varianza del campione. Quindi una stima della grandezza misurata è semplicemente

$$\bar{x} = \frac{1}{N} (X_1 + \dots + X_N)$$

e una stima della varianza del campione è

$$s^2 = \frac{1}{N-1} [(X_1 - \bar{x})^2 + \dots + (X_N - \bar{x})^2].$$

Esempio Il tempo di vita di un tipo di Compedine è, in media, 500 ore (h), con deviazione standard 80h. Comprate 16 Compedine, e assumendo che questo campione abbia legge normale, qual è la probabilità che la media empirica sia maggiore di 525h?

Siano X_j , $j=1, \dots, 16$, v.a. che rappresentano le durate delle j -sime Compedine. Le X_j hanno legge $N(500, 80^2)$.
Dunque \bar{X} ha legge $N(500, \frac{80^2}{16}) = N(500, 400)$.

Perché

$$P(\bar{X} > 525) = P\left(\frac{\bar{X} - 500}{20} > \frac{25}{20}\right) = 1 - \Phi\left(\frac{1.25}{1}\right) \approx 0.11.$$

Esempio Il tempo impiegato da un microprocessore per svolgere certi processi è descritto da una v.a. normale con media 30 ns (nanosecondi) e deviazione standard di 3 ns.

Se si osserva l'esecuzione di 16 processi, qual è la probabilità che la varianza empirica S^2 sia maggiore di 15 ns?

Se $X_j, 1 \leq j \leq 16$ è la v.a. che rappresenta la durata del j -esimo processo, X_j ha legge $N(30, 9)$. Per il teorema di Cochran, S^2 ha legge $\chi^2(15)$. Si ha

$$P(S^2 > 15) = P\left(\frac{n-1}{\sigma^2} S^2 > \frac{15 \cdot 15}{9}\right) = P\left(\frac{n-1}{\sigma^2} S^2 > 25\right) = 1 - P\left(\frac{n-1}{\sigma^2} S^2 \leq 25\right).$$

Dalle tabelle della legge $\chi^2(15)$, il valore 25 corrisponde circa al quantile $\chi^2_{0.95}(15)$. Perciò

$$P(S^2 > 15) = 1 - 0.95 = 0.05.$$