

Stima Puntuale e Test d'Ipotesi

Corso di Modelli Statistici e Statistical Learning

a.a. 2023/2024 - Primo Semestre

***Corso di Laurea Magistrale in
Ingegneria Informatica (percorso AI-ML)
DIIMES***

Prof. Filippo DOMMA

(filippo.domma@unical.it)

Dipartimento di Economia, Statistica e Finanza, Università della Calabria

Proprietà Asintotiche

Gli studi di tipo asintotico riguardano il comportamento degli stimatori al divergere della dimensione campionaria. In tale contesto, ci si aspetta che all'aumentare della dimensione campionaria la *performance* dello stimatore migliori in quanto aumenta l'informazione circa l'obiettivo di stima.

A tal fine, dato lo stimatore $T=T(\mathbf{X})$ della funzione $g(\theta)$, consideriamo la successione di stimatori costruita nel seguente modo

$$\begin{array}{ll} X_1 & \longrightarrow T_1 = T(X_1) \\ X_1, X_2 & \longrightarrow T_2 = T(X_1, X_2) \\ X_1, X_2, X_3 & \longrightarrow T_3 = T(X_1, X_2, X_3) \\ \dots & \dots \\ X_1, X_2, X_3, \dots, X_n & \longrightarrow T_n = T(X_1, X_2, X_3, \dots, X_n) \\ \dots & \dots \end{array}$$

Ad esempio, consideriamo lo stimatore media campionaria

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

della media di una popolazione $g(\theta)=E(X)$, la successione di stimatori media campionaria è data da:

X_1	\longrightarrow	$\bar{X}_1 = X_1$
X_1, X_2	\longrightarrow	$\bar{X}_2 = \frac{X_1 + X_2}{2}$
X_1, X_2, X_3	\longrightarrow	$\bar{X}_3 = \frac{X_1 + X_2 + X_3}{3}$
....	
$X_1, X_2, X_3, \dots, X_n$	\longrightarrow	$\bar{X}_n = \frac{X_1 + X_2 + X_3 + \dots + X_n}{n}$
....	

Stimatori asintoticamente non-distorti

Se uno stimatore è distorto per n fissato, possiamo chiederci se la sua distorsione si riduce all'aumentare della dimensione campionaria e, al limite, se si annulla.

Definizione. **Stimatore asintoticamente non-distorto**

Lo stimatore $T_n = T(X_1, X_2, X_3, \dots, X_n)$ di $g(\theta)$ è asintoticamente non-distorto se e solo se

$$\lim_{n \rightarrow \infty} E(T_n) = g(\theta) \quad \forall \theta \in \Theta$$

Esempio. Abbiamo visto che lo stimatore naturale della varianza di una popolazione, $S^{2'}$, è uno stimatore distorto per $V(X)$ perché $E(S^{2'}) = \left(1 - \frac{1}{n}\right) \times V(X)$.

D'altra parte, il $\lim_{n \rightarrow \infty} E(S^{2'}) = \lim_{n \rightarrow \infty} \left(1 - \frac{1}{n}\right) \times V(X) = V(X)$ possiamo, quindi, concludere che $S^{2'}$ è uno stimatore asintoticamente non-distorto.

La Consistenza

Una proprietà asintotica legata alla dispersione dello stimatore intorno alla funzione del parametro da stimare è la cosiddetta consistenza. Si richiede che all'aumentare della dimensione campionaria e, quindi, all'aumentare delle informazioni sulla quantità da stimare, lo stimatore debba fornire stime sempre più “vicine” alla funzione $g(\theta)$. In altri termini, all'aumentare di n ci si aspetta che la dispersione di T_n intorno a $g(\theta)$ diminuisca. Tale proprietà è formalizzata nella seguente

Definizione. **Consistenza forte**

Lo stimatore $T_n = T(X_1, X_2, X_3, \dots, X_n)$ di $g(\theta)$ è fortemente consistente se converge quasi certamente a $g(\theta)$, cioè se

$$P \left\{ \lim_{n \rightarrow \infty} T_n = g(\theta) \right\} = 1 \quad \forall \theta \in \Theta$$

Una condizione più debole della consistenza forte è la consistenza semplice espressa nella definizione riportata di seguito

Definizione. **Consistenza semplice (o debole)**

Lo stimatore $T_n = T(X_1, X_2, X_3, \dots, X_n)$ di $g(\theta)$ è debolmente consistente per $g(\theta)$ se, scelto un $\varepsilon > 0$ qualsiasi, si ha

$$\lim_{n \rightarrow \infty} P\{|T_n - g(\theta)| \leq \varepsilon\} = 1 \quad \forall \theta \in \Theta$$

Equivalentemente, si può scrivere

$$\lim_{n \rightarrow \infty} P\{g(\theta) - \varepsilon \leq T_n \leq g(\theta) + \varepsilon\} = 1 \quad \forall \theta \in \Theta \text{ e } \varepsilon > 0$$

In letteratura, alcuni autori fanno riferimento alla **consistenza in media quadratica**, affermando che lo stimatore $T_n(\mathbf{X})$ di $g(\theta)$ converge in media quadratica se

$$\lim_{n \rightarrow \infty} EQM(T_n) = 0 \quad \forall \theta \in \Theta$$

Metodi di Stima

Dato il modello parametrico

$$\mathcal{M} = \{\mathcal{P}, \mathcal{X}\}$$

Supponiamo che la famiglia di distribuzione \mathcal{P} sia parametrizzata da un vettore di parametri di dimensione r , cioè

$$\mathcal{P} = \{f(.; \boldsymbol{\theta}): \boldsymbol{\theta} \in \boldsymbol{\Theta} \subset \mathbb{R}^r\}$$

L'obiettivo è quello di descrivere dei metodi che consentono di costruire stimatori per gli elementi del vettore dei parametri $\boldsymbol{\theta}$.

In letteratura esistono diverse tecniche, in questa parte vedremo solo il Metodo dei Momenti e il Metodo della Massima Verosimiglianza. Il Metodo dei Minimi Quadrati verrà esposto nelle prossime lezioni.

Metodo dei Momenti

Abbiamo visto che il momento dall'origine di ordine k , indicato con $\mu_k = E[X^k] = \mu_k(\boldsymbol{\theta})$, in generale, è una funzione del vettore sconosciuto dei parametri della popolazione. Dato un campione casuale di dimensione n , indichiamo con M_j il momento j -esimo campionario, cioè $M_j = \frac{1}{n} \sum_{i=1}^n X_i^j$.

Uguagliando, ordinatamente, i primi r momenti campionari ai primi r momenti della popolazione, otteniamo un sistema di r -equazioni in r -incognite, cioè

$$\begin{cases} M_1 = \mu_1(\boldsymbol{\theta}) \\ M_2 = \mu_2(\boldsymbol{\theta}) \\ \dots\dots\dots \\ M_r = \mu_r(\boldsymbol{\theta}) \end{cases} \Rightarrow \begin{cases} \frac{1}{n} \sum_{i=1}^n X_i = \mu_1(\theta_1, \dots, \theta_r) \\ \frac{1}{n} \sum_{i=1}^n X_i^2 = \mu_2(\theta_1, \dots, \theta_r) \\ \dots\dots\dots \\ \frac{1}{n} \sum_{i=1}^n X_i^r = \mu_r(\theta_1, \dots, \theta_r) \end{cases}$$

Se $\hat{\boldsymbol{\theta}} = (\hat{\theta}_1, \dots, \hat{\theta}_r)^t$ è l'unica soluzione del sistema allora $\hat{\boldsymbol{\theta}} = (\hat{\theta}_1, \dots, \hat{\theta}_r)^t$ è lo stimatore di $\boldsymbol{\theta} = (\theta_1, \dots, \theta_r)$ ottenuto con il metodo dei momenti.

Esempio. Dato il c.c. $\mathbf{X} = (X_1, X_2, X_3, \dots, X_n)$ di dimensione n , estratto da una v.c. Normale di parametri μ e σ^2 , determinare lo stimatore dei momenti del vettore di parametri $\theta = (\mu, \sigma^2)$. In questo caso, abbiamo due parametri da stimare e, quindi, dobbiamo costruire un sistema di 2 equazioni in due incognite, imponendo che i primi due momenti della popolazione siano uguali ai primi due momenti campionari, cioè

$$\begin{cases} M_1 = \mu_1(\boldsymbol{\theta}) \\ M_2 = \mu_2(\boldsymbol{\theta}) \end{cases} \Rightarrow \begin{cases} \frac{1}{n} \sum_{i=1}^n X_i = \mu \\ \frac{1}{n} \sum_{i=1}^n X_i^2 = \mu^2 + \sigma^2 \end{cases} \Rightarrow \begin{cases} \hat{\mu} = \bar{X} \\ \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}^2 \end{cases}$$

Esempio. Sia $\mathbf{X} = (X_1, \dots, X_n)$ un campione casuale estratto da una popolazione con funzione di densità data da:

$$f(x; \theta) = \frac{3x^2}{\theta^3} \quad 0 \leq x \leq \theta$$

- a) determinare lo stimatore di θ con il metodo dei momenti;
- b) stabilire se lo stimatore ottenuto è non-distorto e consistente in media quadratica

Metodo della Massima Verosimiglianza

Tale metodo si basa sull'idea che fd (o fp) diverse diano luogo, in generale, a campioni diversi e, quindi, sia più plausibile che il campione osservato provenga da una fd (o fp) con determinati valori di θ piuttosto che da altre fd (o fp) per le quali il suo realizzarsi sarebbe meno plausibile.

In altri termini, una volta osservato il campione (x_1, \dots, x_n) , si cerca di determinare quale distribuzione ha generato con maggiore plausibilità il campione stesso. Così, fissate le osservazioni campionarie (x_1, \dots, x_n) si osserva come variano i valori di $f(x; \theta)$ al variare di θ nello spazio parametrico Θ . In tal modo, otteniamo una funzione definita sullo spazio parametrico Θ , cioè $L: \Theta \rightarrow \mathbb{R}^+$.

ESEMPIO: GRAFICO

Definizione. Dato il modello parametrico $\mathcal{M} = \{P, X\}$, sia (X_1, \dots, X_n) un campione casuale indipendente ed identicamente distribuito estratto da $f(\mathbf{x}; \theta)$, e sia $\mathbf{x} = (x_1, \dots, x_n)$ il campione osservato. La funzione di densità (o di probabilità) congiunta $f(\mathbf{x}; \theta)$ definita sullo spazio parametrico Θ , viene detta funzione di verosimiglianza del campione osservato \mathbf{x} ed indicata con:

$$L(\theta; \mathbf{x}) = f(\mathbf{x}; \theta)$$

Quindi, la funzione di verosimiglianza è una applicazione definita nello spazio parametrico che associa valori nell'insieme dei numeri reali positivi, cioè

$$L: \Theta \rightarrow \mathbb{R}^+$$

Possiamo dire che massimizzando $L(\theta; \mathbf{x})$ individuiamo quel valore di θ , diciamo $\hat{\theta} = \hat{\theta}(\mathbf{x})$ e, quindi, quella fd (o fp) che con maggior verosimiglianza ha generato le osservazioni campionarie (x_1, \dots, x_n) . Da qui la seguente:

Definizione. Dato il modello parametrico $\mathbf{M} = \{\mathcal{X}, \mathbf{P}\}$, sia $L(\theta; \mathbf{x})$ la f.v. del campione osservato $\mathbf{x} = (x_1, \dots, x_n)$. Il valore $\hat{\theta} = \hat{\theta}(\mathbf{x})$ tale per cui

$$L[\hat{\theta}(\mathbf{x}); \mathbf{x}] = \sup_{\theta \in \Theta} L[\theta; \mathbf{x}]$$

se esiste è detta stima di massima verosimiglianza. Tale valore valutato nel campione casuale $\mathbf{X}=(X_1, \dots, X_n)$ fornisce lo stimatore di massima verosimiglianza $\hat{\theta} = \hat{\theta}(\mathbf{X})$.

Sotto l'ipotesi di indipendenza, la fd (o fp) congiunta campionaria fattorizza nel prodotto delle marginali e, quindi, la funzione di verosimiglianza può essere scritta nel seguente modo

$$L(\theta; \mathbf{x}) = f(\mathbf{x}; \theta) = \prod_{i=1}^n f(x_i; \theta)$$

Spesso, per semplificare i calcoli, si preferisce far riferimento al logaritmo della funzione di verosimiglianza, denominata funzione di log-verosimiglianza

$$\ell(\theta; \mathbf{x}) = \ln L(\theta; \mathbf{x}) = \ln \prod_{i=1}^n f(x_i; \theta) = \sum_{i=1}^n \ln f(x_i; \theta)$$

poiché la trasformazione è monotona il punto che massimizza $L(\theta; \mathbf{x})$ massimizzerà anche $\ell(\theta; \mathbf{x})$; infatti, si ha:

$$\frac{\partial \ell(\theta; \mathbf{x})}{\partial \theta} = \frac{1}{L(\theta; \mathbf{x})} \frac{\partial L(\theta; \mathbf{x})}{\partial \theta} = 0 \iff \frac{\partial L(\theta; \mathbf{x})}{\partial \theta} = 0.$$

Analogamente a quanto detto in precedenza, se la fd (o fp) è funzione di un vettore di parametri θ , con $\theta \in \Theta \subseteq \mathfrak{R}^k$, allora la stima di massima verosimiglianza $\hat{\theta}(\mathbf{x}) = (\hat{\theta}_1(\mathbf{x}), \dots, \hat{\theta}_k(\mathbf{x}))^t$ di $\theta = (\theta_1, \dots, \theta_k)^t$ è quel vettore, se esiste, tale per cui

$$L[\hat{\theta}(\mathbf{x}); \mathbf{x}] = \sup_{\theta \in \Theta} L[\theta; \mathbf{x}].$$

Sostituendo il campione osservato, \mathbf{x} , con il campione casuale, \mathbf{X} , si ottiene lo stimatore di m.v. $\hat{\theta}(\mathbf{X}) = (\hat{\theta}_1(\mathbf{X}), \dots, \hat{\theta}_k(\mathbf{X}))^t$ del vettore dei parametri $\theta = (\theta_1, \dots, \theta_k)^t$.

Nei casi in cui $\ell(\boldsymbol{\theta}; \mathbf{x})$ è differenziabile rispetto a $\boldsymbol{\theta}$ allora, utilizzando i metodi dell'analisi matematica, costruiamo il seguente sistema

$$\begin{cases} \frac{\partial \ell(\boldsymbol{\theta}; \mathbf{x})}{\partial \theta_1} = 0 \\ \dots\dots\dots \\ \frac{\partial \ell(\boldsymbol{\theta}; \mathbf{x})}{\partial \theta_k} = 0 \end{cases}$$

formato da k equazioni di verosimiglianza (le derivate parziali della log-verosimiglianza rispetto alle k componenti del vettore $\boldsymbol{\theta}$) in k incognite (gli elementi del vettore incognito $\boldsymbol{\theta}$). La soluzione di detto sistema fornisce la stima di massima verosimiglianza $\hat{\boldsymbol{\theta}} = (\hat{\theta}_1, \dots, \hat{\theta}_k)$ del vettore di parametri incogniti $\boldsymbol{\theta} = (\theta_1, \dots, \theta_k)$, una volta verificato che la matrice delle derivate seconde

Esempio. Bernoulli

Esempio. Normale.

Proprietà stimatori di massima verosimiglianza

Teorema (*Principle of invariance*).

Sia $g(\cdot)$ una funzione definita nello spazio parametrico con immagine in $\Omega \subset \mathfrak{R}$.
Se $\hat{\theta} = \hat{\theta}(\mathbf{x})$ è la stima di massima verosimiglianza di $\theta \in \Theta \subset \mathfrak{R}$, allora $g(\hat{\theta})$ è la corrispondente stima di massima verosimiglianza di $g(\theta)$.

In estrema sintesi, gli stimatori di massima verosimiglianza:

- non necessariamente sono non-distorti;
- sono consistenti
- sono asintoticamente normali e pienamente efficienti

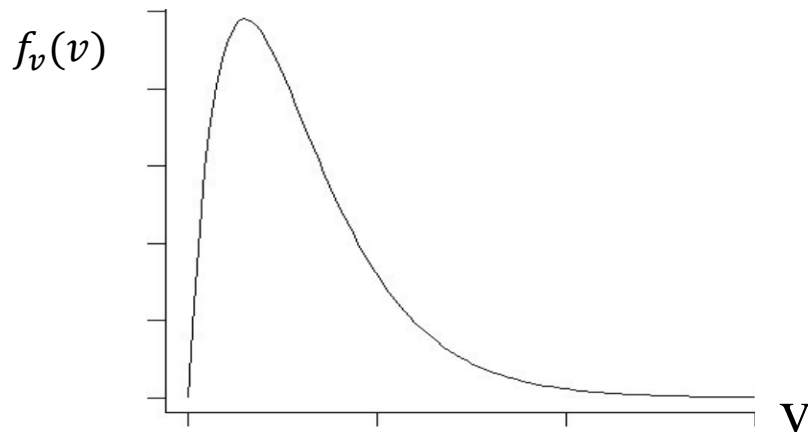
Alcune utili variabili casuali ottenute da trasformazioni
della variabile casuale Normale.

Variabile Casuale chi-quadro χ^2

Definizione. Date k variabili casuali Normali e *indipendenti*, $X_i \sim N(\mu_i, \sigma_i^2)$ per $i=1, \dots, k$, la variabile casuale così definita

$$V = \sum_{i=1}^k \left(\frac{X_i - \mu_i}{\sigma_i} \right)^2$$

si distribuisce secondo una chi-quadrato di parametro k . Viene indicata con $\chi^2(k)$.



$$E(V) = k$$

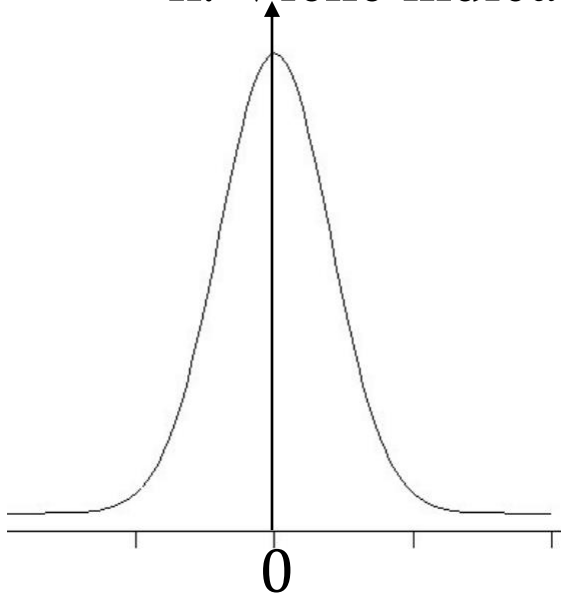
$$\text{Var}(V) = 2k$$

Variabile Casuale t – Student

Definizione. Se $X \sim N(\mu, \sigma^2)$ e $V \sim \chi^2(k)$ sono *indipendenti*, allora il rapporto

$$T = \frac{\frac{X - \mu}{\sigma}}{\sqrt{\frac{V}{k}}} = \frac{Z}{\sqrt{\frac{V}{k}}}$$

si distribuisce secondo una t-Student di parametro k . Viene indicata con $t(k)$.



$$E(T) = 0$$

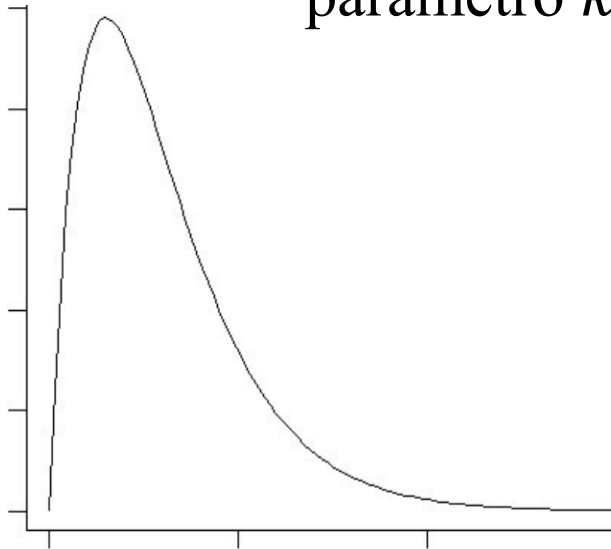
$$\text{Var}(T) = \frac{k}{k-2}$$

Variabile Casuale **F** di Fisher

Definizione. Se $V_1 \sim \chi^2(k_1)$ e $V_2 \sim \chi^2(k_2)$ sono *indipendenti*, allora il rapporto

$$F = \frac{\frac{V_1}{k_1}}{\frac{V_2}{k_2}}$$

si distribuisce secondo una F di Fisher di parametri parametro k_1 e k_2 . Viene indicata con $F(k_1, k_2)$.



$$E(F) = \frac{n}{n-2} \quad \text{per } n > 2$$

$$V(F) = \frac{n^2(2m + 2n - 4)}{m(n - 2)^2(n - 4)} \quad \text{per } n > 4$$

Gradi di libertà

I gradi di libertà sono il numero di componenti del campione casuale che possono essere scelti liberamente dato dalla dimensione campionaria meno il numero di vincoli sullo spazio campionario costituite dalle stime preliminari che bisogna effettuare per calcolare lo stimatore ‘*corrente*’.

Dato un c.c. \mathbf{X} di dimensione n , lo stimatore ‘*corrente*’
 $T_c = h(\mathbf{X}, T_1, \dots, T_r)$ ha $n-r$ gradi di libertà.

Esempio. Supponiamo che la media campionaria su un campione di dimensione 5 risulta essere pari a 20. Se vogliamo calcolare la devianza campionaria (numeratore della varianza) $\sum_{i=1}^5 (x_i - \bar{x})^2$ che rappresenta lo stimatore corrente, allora la media campionaria rappresenta un vincolo sullo spazio campionario e possiamo scegliere liberamente solo 4 componenti del campione.