

Tema di Statistica Matematica

Primo appello

23 gennaio 2019

1) Sia (X_1, X_2, \dots, X_n) un campione casuale da una distribuzione Uniforme su $(0, \theta)$ e sia dato lo stimatore $\hat{\theta}_k = k \cdot \max(X_1, X_2, \dots, X_n)$ per $k > 1$. Trovare

- a) l'errore quadratico medio di $\hat{\theta}_k$;
- b) il valore di k per il quale $\hat{\theta}_k$ é non distorto per θ ;
- c) il valore di k che minimizza l'errore quadratico medio e commentare il risultato ottenuto.

Dimostrare, infine, che il valore $k = 1$ restituisce lo stimatore di massima verosimiglianza $\hat{\theta}_n$ di θ e confrontare quest'ultimo stimatore con quello ottenuto al punto c).

2) Sia (X_1, X_2, \dots, X_n) un campione casuale proveniente da una distribuzione la cui densit  sia data da

$$f_X(x; \theta) = \theta x^{\theta-1} \mathbb{I}_{(0,1)}(x), \quad \theta > 0.$$

Esiste una funzione di θ per la quale esiste uno stimatore non distorto la cui varianza raggiunge il limite inferiore di Rao-Cram r? Se   cos , individuarla; se no, dimostrare perch  no.

3) Sia (x_1, x_2, \dots, x_n) la determinazione di un campione casuale proveniente da una distribuzione continua con densit 

$$f_X(x; \theta) = \theta (x+1)^{-(\theta+1)} \mathbb{I}_{[0,\infty)}(x)$$

con $\theta > 0$ parametro ignoto. Indicata con F_X la funzione di distribuzione di X ,

- a) dimostrare che la variabile $V = -2 \ln(1 - F_X(X))$ ha distribuzione chi-quadrato con $\nu = 2$ gradi di libert ;
- b) usando il risultato ottenuto al precedente punto a), ricavare una statistica pivot per l'inferenza su θ e costruire un intervallo di confidenza per il parametro θ di livello (esatto) $(1 - \alpha) = 0.95$;
- c) stabilire se lo stimatore di massima verosimiglianza di θ   o meno non distorto per θ ;
(**Sugg.:** *disuguaglianza di Jensen*: $\mathbb{E}(g(X)) \geq g(\mathbb{E}(X))$)
- d) ricavare la regione critica pi  potente di livello $\alpha = 0.05$ per il problema della verifica del sistema di ipotesi $H_0 : \theta = \theta_0$ vs. $H_1 : \theta > \theta_0$.

Ex. 1)

$$\begin{aligned}
 a) \text{MSE}_{\theta}(\hat{\theta}_k) &= [B_{\theta}(\hat{\theta}_k)]^2 + \text{Var}_{\theta}(\hat{\theta}_k) \\
 &= \left(\frac{nk}{n+1} - 1\right)^2 \theta^2 + \frac{nk^2 \theta^2}{(n+2)(n+1)^2} \\
 &= \frac{\theta^2}{(n+1)^2} \left[(nk - (n+1))^2 + \frac{nk^2}{n+2} \right] \\
 &= \frac{\theta^2}{(n+1)^2(n+2)} \left[(nk - (n+1))^2 (n+2) + nk^2 \right] \\
 &= \frac{\theta^2}{(n+1)^2(n+2)} \left[(n^2 k^2 - 2(n+1)nk + (n+1)^2)(n+2) + nk^2 \right] \\
 &= \frac{\theta^2}{(n+1)^2(n+2)} \left[(n + (n+2) \cdot n^2) k^2 - 2n(n+1)(n+2)k + (n+1)^2(n+2) \right]
 \end{aligned}$$

$$b) \mathbb{E}_{\theta}(\hat{\theta}_k) = \theta \Rightarrow k = \frac{n+1}{n} \quad \text{dal momento che}$$

$$\mathbb{E}_{\theta}(X_{(n)}) = \frac{n}{n+1} \theta \quad \text{dove}$$

$$X_{(n)} = \max(X_1, \dots, X_n)$$

c) Come si può vedere al punto a) l'espressione quadratica MEDIO di $\hat{\theta}_k$ visto come funzione di k è data da

$$\theta^2 \left[\frac{n + (n+2)n^2}{(n+1)^2(n+2)} k^2 - \frac{2n(n+1)(n+2)}{(n+1)^2(n+2)} k + \frac{(n+1)^2(n+2)}{(n+1)^2(n+2)} \right]$$

si ottiene

$$\frac{d}{dk} \text{MSE}_{\theta}(\hat{\theta}_k) = \theta^2 \left[\frac{2(n+(n+2) \cdot n^2)}{(n+1)^2(n+2)} k - \frac{2n(n+1)(n+2)}{(n+1)^2(n+2)} \right] = 0$$

e dunque, con qualche passaggio, si ottiene

$$k = \frac{n+2}{n+1}$$

Insomma

$$\frac{d^2}{dk^2} \text{MSE}_{\theta}(\hat{\theta}_k) = \frac{2(n+(n+2)n^2)}{(n+1)^2(n+2)} > 0$$

e dunque $k = \frac{n+2}{n+1}$ è punto di minimo.

Commento: dai calcoli appena fatti emerge che lo stimatore di minimo errore quadratico medio in questa famiglia di distribuzioni è distorto se confrontato con $k = \frac{n+1}{n}$ per lo stimatore (non distorto) di massima verosimiglianza al pto b) (||| opportunamente corretto)

d) $k=1$.

La funzione di verosimiglianza per la famiglia di distribuzioni in questione è data da

$$L(\theta; x) = \frac{1}{\theta^n} \prod_{i=1}^n \Pi_{(0, \theta)}(x_i)$$

Questa funzione risulta massima in corrispondenza del valore di θ più piccolo possibile ma non più piccolo del massimo campionario altrimenti uno degli indicatori sarà zero.

Perciò

$$\hat{\theta}_n = X_{(n)} = \max(X_1, X_2, \dots, X_n)$$

Commento: i) $\hat{\theta}_n$ non è stimatore non distorto di θ
poiché
$$\mathbb{E}_{\theta}(\hat{\theta}_n) \neq \theta$$

ii) $\hat{\theta}_n$ non è lo stimatore a minimo
errore quadratico medio.

Ex. 2)

In questo caso

$$L(\theta; \underline{x}) = \prod_{i=1}^n \theta x_i^{\theta-1} \mathbb{1}_{(0,1)}(x_i) = \theta^n \left(\prod_{i=1}^n x_i \right)^{\theta-1} \cdot \prod_{i=1}^n \mathbb{1}_{(0,1)}(x_i)$$

si ottiene

$$l(\theta; \underline{x}) = n \ln(\theta) + (\theta-1) \sum_{i=1}^n \ln(x_i) + \ln \prod_{i=1}^n \mathbb{1}_{(0,1)}(x_i)$$

La funzione score è data da

$$\frac{d}{d\theta} l(\theta; \underline{x}) = \frac{n}{\theta} + \sum_{i=1}^n \ln(x_i) \equiv S(\theta; \underline{x})$$

e, ricordando che $\mathbb{E}_{\theta}(S(\theta; \underline{x})) = 0$ si ha

$$\mathbb{E}_{\theta} \left(\frac{n}{\theta} + \sum_{i=1}^n \ln(x_i) \right) = \frac{n}{\theta} + \mathbb{E}_{\theta} \left(\sum_{i=1}^n \ln(x_i) \right) = 0$$

si ottiene

$$\mathbb{E}_{\theta} \left(-\frac{1}{n} \sum_{i=1}^n \ln(x_i) \right) = \frac{1}{\theta}$$

da cui segue che $T_n = -\frac{1}{n} \sum_{i=1}^n \ln(x_i)$ è uno stimatore
NON DISTORTO di $\eta(\theta) = \frac{1}{\theta}$.

Nota: è immediato osservare che $T_n = -\frac{1}{n} \sum_{i=1}^n \ln(x_i)$

è stimatore di massima verosimiglianza di $\eta(\theta)$

si ottiene

$$\hat{\theta}_n = -\frac{n}{\sum_{i=1}^n \ln(x_i)}$$

è stimatore di massima verosimiglianza di θ .

On pseudo

$$T_n = -\frac{1}{n} \sum_{i=1}^n \ln(X_i)$$

stimatore NON DISTORSO di $g(\theta) = \frac{1}{\theta}$ e FUNZIONE di STATISTICA
SUFFICIENTE MINIMALE di θ , è anche STIMATORE UMVU di $g(\theta) = \frac{1}{\theta}$
in virtù delle CONOSCENZE della TEORIA di RAO-BLACKWELL.

Verifichiamo ora se T_n è anche STIMATORE EFFICIENTE
per $g(\theta) = \frac{1}{\theta}$.

Per far ciò, calcoliamo la VARIANZA di T_n e confrontiamola
con il LIMITE INFERIORE di RAO-CRAMER per la VARIANZA di un
qualsiasi stimatore NON DISTORSO di $g(\theta)$.

Cominciamo con come

$$Y_i = -\ln(X_i) \Rightarrow X_i = e^{-Y_i}, \quad \left| \frac{d}{dY_i} X_i \right| = e^{-Y_i}$$

sicché

$$\begin{aligned} f_{Y_i}(y; \theta) &= \theta e^{-y(\theta-1)} e^{-y} \mathbb{1}_{[0, +\infty)}(y) \\ &= \theta e^{-\theta y} \mathbb{1}_{[0, +\infty)}(y) \end{aligned}$$

si ha

$$Y_i \sim \text{Exp}(\theta) \quad \text{sicché} \quad E_{\theta}(Y_i) = \frac{1}{\theta} \quad \text{e} \quad \text{Var}_{\theta}(Y_i) = \frac{1}{\theta^2}$$

Dunque per la PROPRIETÀ di RIPRODUCIBILITÀ

$$T_n = \frac{1}{n} \sum_{i=1}^n Y_i$$

$$\text{sicché} \quad \text{Var}_{\theta}(T_n) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}_{\theta}(Y_i) = \frac{1}{n^2} \frac{n}{\theta^2} = \frac{1}{n\theta^2}$$

Ora, tenuto conto del fatto che il LIMITE INFERIORE per la varianza di uno STIMATORE NON DISTORTO di una funzione $g(\theta)$ del parametro θ è dato da

$$\text{Var}_{\theta}(T_n) \geq \frac{[g'(\theta)]^2}{I_n(\theta)}$$

ed essendo

$$[g'(\theta)]^2 = \left[-\frac{1}{\theta^2}\right]^2 = \frac{1}{\theta^4}$$

$$I_n(\theta) = -E_{\theta} \left[\frac{d}{d\theta} S(\theta; \underline{X}) \right] = -E \left[-\frac{n}{\theta^2} \right] = \frac{n}{\theta^2}$$

si ha

$$\text{Var}_{\theta}(T_n) \geq \frac{\theta^2}{n} \cdot \frac{1}{\theta^4} = \frac{1}{n\theta^2}$$

ricchi

$$\text{Var}_{\theta}(T_n) = \frac{1}{n\theta^2}$$

Cm $T_n = -\frac{1}{n} \sum_{i=1}^n \ln(X_i)$, da cui la conclusione che

T_n è STIMATORE EFFICIENTE di $g(\theta) = \frac{1}{\theta}$.

Ex. 3)

a) Si ha

$$\begin{aligned}P(V \leq v) &= P(-2 \ln(1 - F_X(x)) \leq v) \\&= P(F_X(x) \leq 1 - e^{-\frac{1}{2}v}) \\&= 1 - e^{-\frac{1}{2}v}\end{aligned}$$

essendo, come è noto, $F_X(x) \sim U(0,1)$. Quindi V ha una distribuzione ESPONENZIALE di parametro $\theta = 2$ ovvero

$$V \sim \text{Exp}(2) \equiv \mathcal{G}(1, 2) \equiv \chi^2_2$$

ricordando che ESPONENZIALE e CHI-QUADRATO sono casi PARTICOLARI di GAMMA.

b) In base ai risultati di cui al punto a) si ha che

$$-2 \sum_{i=1}^n \ln(1 - F_X(x_i)) \sim \mathcal{G}(n, 2) \equiv \chi^2_{2n}$$

D'altra parte

$$F_X(x) = \int_0^x \theta (t+1)^{-(\theta+1)} dt = \int_1^{x+1} w^{-(\theta+1)} dw = 1 - (x+1)^{-\theta}$$

Quindi

$$\begin{aligned}-2 \sum_{i=1}^n \ln(1 - F(x_i)) &= -2 \sum_{i=1}^n \ln(x_i+1)^{-\theta} \\&= 2\theta \sum_{i=1}^n \ln(x_i+1) \sim \chi^2_{2n}\end{aligned}$$

è la STATISTICA PIVOT CORRETTA per θ .

Portanto l'INTERVALLO DI CONFIDENZA ESATTO per θ di livello 0.95 è dato da

$$IC_{\theta}(0.95) = \left[\frac{\chi^2_{2n; 0.025}}{2 \sum_{i=1}^n \ln(x_i+1)}, \frac{\chi^2_{2n; 0.975}}{2 \sum_{i=1}^n \ln(x_i+1)} \right]$$

dove $\chi^2_{2n; \alpha}$ indica il QUANTILE di ordine α della DISTRIBUZIONE CHI-QUADRO con $2n$ gradi di libertà.

c) la funzione di VEROSIMILITUDINE è data da

$$L(\theta; \underline{x}) = \theta^n \prod_{i=1}^n (x_i+1)^{-(\theta+1)} \prod_{i=1}^n \mathbb{1}_{\mathbb{R}^+}(x_i)$$

ricchi

$$l(\theta; \underline{x}) = \ln L(\theta; \underline{x}) = n \ln(\theta) - (\theta+1) \sum_{i=1}^n \ln(x_i+1) + \ln \left(\prod_{i=1}^n \mathbb{1}_{\mathbb{R}^+}(x_i) \right),$$

e dunque

$$\hat{\theta}_n = \frac{n}{\sum_{i=1}^n \ln(x_i+1)}$$

è la soluzione di MASSIMA VEROSIMILITUDINE di θ $\left[\frac{d^2}{d\theta^2} l(\theta; \underline{x}) < 0 \right]$

Indicando con $S(\theta; \underline{x})$ la funzione score otteni

$$S(\theta; \underline{x}) = \frac{d}{d\theta} l(\theta; \underline{x}) = \frac{n}{\theta} - \sum_{i=1}^n \ln(x_i+1)$$

e ricordando che

$$\mathbb{E}_{\theta}(S(\theta; \underline{X})) = 0$$

sì ha che

$$\mathbb{E}_{\theta} \left[\sum_{i=1}^n \ln(x_i+1) \right] = \frac{n}{\theta}$$

Utilizzando la DIVERGENZA di JENSEN

$$\mathbb{E}_\theta(\hat{\theta}_n) = \mathbb{E}_\theta \left[\frac{n}{\sum_{i=1}^n \ln(X_i + 1)} \right] > \frac{n}{\mathbb{E}_\theta \left[\sum_{i=1}^n \ln(X_i + 1) \right]} = \theta$$

e pertanto la STIMATORE di massima VEROSIMILITUDINE risulta essere DISTORTO per θ .

Nota: DIVERGENZA di JENSEN

$$\mathbb{E}_\theta(\varphi(X)) \geq \varphi(\mathbb{E}_\theta(X))$$

dove φ è una funzione vale se solo se $\varphi(X)$ è funzione LINEARE di X .

d) È facile verificare che la famiglia parametrica considerata costituisce una famiglia ESPONENZIALE (regolare) a $k=1$ parametri con STATISTICA NATURALE (SUFFICIENTE, MINIMALE e COMPLETA)

$$T_n(X_1, \dots, X_n) = \sum_{i=1}^n \ln(X_i + 1)$$

Pertanto, il test ottimo per il problema di verifica di ipotesi considerato ha regione critica C_α basata su $T_n(X_1, \dots, X_n)$ o opportuna trasformazione.

Sappiamo che sotto H_0

$$2\theta_0 T_n(\underline{X}) \sim \chi_{2n}^2$$

e inoltre che

$$T_n(\underline{X}) = \frac{n}{\hat{\theta}_n}$$

e grandi valori della statistica di verifica sono, evidentemente, contrari all'ipotesi nulla.

Ne segue che la regione critica più potente cercata è

$$C_\alpha = \{ \underline{x} \in \mathcal{X} : 2 \theta_0 T_n(\underline{x}) < \chi^2_{2n; \alpha} \}$$

con α livello di significatività del test, e, nel nostro caso fissato in $\alpha = 0.05$.

Es. 4)

a) Si può osservare che

$$\tau = P_\theta(X \leq 1) = P_\theta(X=0) + P_\theta(X=1) = \theta + \theta(1-\theta) = 2\theta - \theta^2$$

e perciò $\tau = g(\theta)$ con $g'(\theta) = 2 - 2\theta > 0$ per ogni valore di $\theta \in (0, 1)$. Quindi τ è funzione monotona crescente di θ .

D'altra parte, la funzione di verosimiglianza è

$$L(\theta; \underline{x}) = \prod_{i=1}^n \theta(1-\theta)^{x_i} = \theta^n (1-\theta)^{\sum_{i=1}^n x_i}$$

e quindi

$$\ell(\theta; \underline{x}) = \ln L(\theta; \underline{x}) = n \ln(\theta) + \ln(1-\theta) \cdot \sum_{i=1}^n x_i$$

Di conseguenza

$$S(\theta; \underline{x}) = \frac{d}{d\theta} \ell(\theta; \underline{x}) = \frac{n}{\theta} - \frac{1}{1-\theta} \sum_{i=1}^n x_i$$

da cui, uguagliando a zero la precedente e risolvendo rispetto θ si ha che la funzione di TV di θ è data da

$$\hat{\theta}_n = \frac{n}{n + \sum_{i=1}^n x_i} = \frac{1}{1 + \bar{x}_n}$$

e inoltre

$$\begin{aligned} I_n(\theta) &= -E_\theta \left[\frac{d^2}{d\theta^2} \ell(\theta; \underline{x}) \right] = -E_\theta \left[-\frac{n}{\theta^2} - \frac{1}{(1-\theta)^2} \sum_{i=1}^n x_i \right] \\ &= \frac{n}{\theta^2} + \frac{1}{(1-\theta)^2} \cdot n \cdot \frac{1-\theta}{\theta} = \frac{n(1-\theta) + n\theta}{\theta^2(1-\theta)} = \frac{n}{\theta^2(1-\theta)} \end{aligned}$$

dato che $E(X) = \frac{1-\theta}{\theta}$.

Ricordando che $\hat{\theta}_n$ è lo stimatore di MV di θ , per le sue note proprietà si ha

$$(\hat{\theta}_n - \theta) \underset{av}{\sim} N\left(0, \frac{1}{I_n(\theta)}\right) \equiv N\left(0, \frac{\theta^2(1-\theta)}{n}\right)$$

sicché

$$\left[\hat{\theta}_n - 1.96 \sqrt{\frac{\theta^2(1-\theta)}{n}}, \quad \hat{\theta}_n + 1.96 \sqrt{\frac{\theta^2(1-\theta)}{n}} \right]$$

costituisce un INTERVALLO di CONFIDENZA APPROSSIMATO per θ di livello 0,95.

Ora in virtù del METODO DELTA

$$g(\hat{\theta}_n) \underset{av}{\sim} N(g(\theta), (g'(\theta))^2 I_n^{-1}(\theta))$$

con $g(\theta) = 2\theta - \theta^2$ e $\underbrace{g'(\theta) I_n^{-1}(\theta)}_{\parallel \text{Var}_{\theta}(g(\hat{\theta}_n))} = 4 \cdot \frac{\theta^2(1-\theta)^3}{n},$

da cui

$$\frac{g(\hat{\theta}_n) - g(\theta)}{\sqrt{\text{Var}_{\theta}(g(\hat{\theta}_n))}} \underset{av}{\sim} N(0, 1)$$

sicché

$$IC_{g(\theta)}^{(1-\alpha)} = \left[g(\hat{\theta}_n) - z_{1-\alpha/2} \sqrt{\text{Var}_{\theta}^{\hat{}}(g(\hat{\theta}_n))}, \quad g(\hat{\theta}_n) + z_{1-\alpha/2} \sqrt{\text{Var}_{\theta}^{\hat{}}(g(\hat{\theta}_n))} \right]$$

dove

$$g(\hat{\theta}_n) = 2\hat{\theta}_n - \hat{\theta}_n^2 \quad \text{e} \quad \text{Var}_{\theta}^{\hat{}}(\hat{\theta}_n) = 4 \cdot \frac{\hat{\theta}_n^2(1-\hat{\theta}_n)^3}{n}$$

b) Sotto il modello parametrico considerato, lo stimatore di massima verosimiglianza

$$\hat{\theta}_n = \frac{1}{1 + \bar{X}_n}$$

è stimatore consistente per $\theta = P_\theta(X=0)$.

Se la vera legge di X fosse un membro della classe delle DISTRIBUZIONI di Poisson, avremo

$$P_\theta(X=0) = e^{-\theta}$$

D'altra parte, per la legge dei Grandi Numeri,

$$\frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{P} E_\theta(X) = \theta$$

Allora, essendo $\hat{\theta}_n$ funzione continua di \bar{X}_n , al crescere di n , si avrebbe

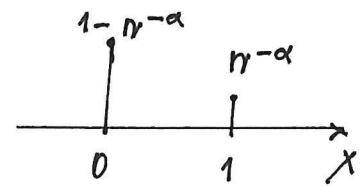
$$\hat{\theta}_n = \frac{1}{1 + \bar{X}_n} \xrightarrow{P} \frac{1}{1 + \theta} \neq e^{-\theta}$$

e quindi $\hat{\theta}_n$ NON RISPONDEREbbe al criterio di consistenza per $P_\theta(X=0)$ sotto modello di Poisson.

Ex 5)

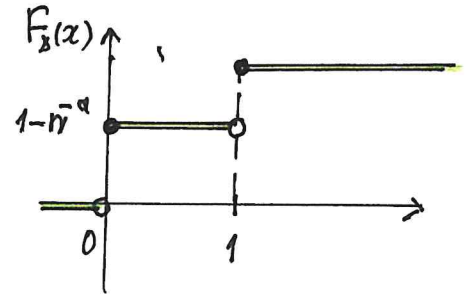
Sia $\{X_n\}_{n \in \mathbb{N}}$ una succ. di v.c. INDEPENDENT

$b(1, n^{-\alpha})$, con $\alpha > 0$.



Ora

$$F_{X_n}(x) = \begin{cases} 0 & x < 0 \\ 1 - \frac{1}{n^\alpha} & 0 \leq x < 1 \\ 1 & x \geq 1 \end{cases}$$



ricchi al crescere di n , qualunque sia $\alpha > 0$, X_n
Converge in distribuzione a zero.