

5. Stimatori di massima verosimiglianza

Sia X_1, \dots, X_n un campione statistico e sia $Y = \varphi(X_1, \dots, X_n)$ una sua statistica. Se Y ha lo scopo di stimare un parametro θ della distribuzione del campione, diciamo che Y è uno *stimatore del parametro* θ .

Supponiamo di conoscere la distribuzione del campione a meno di un parametro θ e supponiamo che tale distribuzione sia discreta o assolutamente continua e dunque dotata di densità (discreta o meno). Tale densità dipenderà dal parametro θ e la indico col simbolo $g(x|\theta)$. La densità congiunta si indica col simbolo $f(x_1, \dots, x_n|\theta)$ e sappiamo che, grazie all'indipendenza delle v.a. che costituiscono il campione, si ha

$$f(x_1, \dots, x_n|\theta) = g(x_1|\theta) \cdot \dots \cdot g(x_n|\theta) = \prod_{i=1}^n g(x_i|\theta).$$

Interpreto $f(x_1, \dots, x_n|\theta)$ come la *plausibilità* che la n -upla x_1, \dots, x_n si realizzi nel campione empirico quando il parametro incognito prende il valore θ . Consideriamo infatti i due casi, discreto e assolutamente continuo

- Se il campione ha distribuzione discreta, allora

$$\mathbb{P}(X_1 = x_1, \dots, X_n = x_n) = \prod_{i=1}^n \mathbb{P}(X_i = x_i) = \prod_{i=1}^n g(x_i|\theta).$$

- Se il campione ha distribuzione assolutamente continua, e se la densità g è continua, allora, per $\delta > 0$ e sufficientemente piccolo si ha

$$\begin{aligned} \mathbb{P}\left(|X_1 - x_1| < \frac{\delta}{2}, \dots, |X_n - x_n| < \frac{\delta}{2}\right) &= \prod_{i=1}^n \mathbb{P}\left(|X_i - x_i| < \frac{\delta}{2}\right) \\ &= \prod_{i=1}^n \mathbb{P}\left(X_i \in \left(x_i - \frac{\delta}{2}, x_i + \frac{\delta}{2}\right)\right) \simeq \prod_{i=1}^n (g(x_i|\theta) \delta)^n \\ &= \delta^n \prod_{i=1}^n g(x_i|\theta) = \delta^n f(x_1, \dots, x_n|\theta) \end{aligned}$$

La funzione f , vista come funzione di θ , si dice *funzione di verosimiglianza*.

Dunque: dato il campione empirico x_1, \dots, x_n , cerco $\hat{\theta} = \hat{\theta}(x_1, \dots, x_n)$ che massimizza la funzione $f(x_1, \dots, x_n|\theta)$. La statistica $\hat{\theta} = \hat{\theta}(X_1, \dots, X_n)$ si dirà *stimatore di massima verosimiglianza del parametro* θ .

Osservazione 5.0.1. Poiché la funzione $\ln: (0, +\infty) \rightarrow \mathbb{R}$ è strettamente monotona crescente, massimizzare $f(x_1, \dots, x_n|\theta) = \prod_{i=1}^n g(x_i|\theta)$ equivale a massimizzare la funzione

$$h(x_1, \dots, x_n|\theta) := \ln f(x_1, \dots, x_n|\theta) = \sum_{i=1}^n \ln g(x_i|\theta)$$

Si ha

$$\frac{\partial}{\partial \theta} h(x_1, \dots, x_n|\theta) = \frac{\partial}{\partial \theta} \sum_{i=1}^n \ln g(x_i|\theta) = \sum_{i=1}^n \frac{\partial}{\partial \theta} \ln g(x_i|\theta) = \sum_{i=1}^n \frac{1}{g(x_i|\theta)} \frac{\partial g(x_i|\theta)}{\partial \theta}$$

5.1 Distribuzione di Bernoulli

In questo caso la distribuzione dipende dal solo parametro $p = \mathbb{P}(X_i = 1) = \mathbb{E}[X_i]$. Sia dunque X_1, \dots, X_n un campione statistico di Bernoulli di parametro incognito $p \in [0, 1]$. Realizzo n prove di Bernoulli e ottengo il campione empirico x_1, \dots, x_n , $x_i \in \{0, 1\}$. Sia $k = k(x_1, \dots, x_n) := \sum_{i=1}^n x_i$. Abbiamo

$$\begin{aligned} f(x_1, \dots, x_n|p) &= \mathbb{P}(X_1 = x_1, \dots, X_n = x_n) = p^k(1-p)^{n-k}, \\ h(x_1, \dots, x_n|p) &= \ln(p^k(1-p)^{n-k}) = k \ln p + (n-k) \ln(1-p). \\ \frac{\partial h}{\partial p} &= \frac{k}{p} - \frac{n-k}{1-p} = \frac{k-np}{p(1-p)} \geq 0 \iff k-np \geq 0 \iff p \leq \frac{k}{n}. \end{aligned}$$

Poiché $k = \sum_{i=1}^n x_i$, lo stimatore di massima verosimiglianza per il parametro p è $\frac{\sum_{i=1}^n X_i}{n}$ cioè la media campionaria \bar{X} . Ricordiamo che la media campionaria è uno stimatore corretto di $\mathbb{E}[X_i] = p$.

5.2 Distribuzione di Poisson

La distribuzione di Poisson è concentrata sugli interi nonnegativi e dipende da un solo parametro:

$$g(x|\lambda) = e^{-\lambda} \frac{\lambda^x}{x!}, \quad x = 0, 1, \dots$$

Dunque

$$\begin{aligned} f(x_1, \dots, x_n|\lambda) &= \prod_{i=1}^n \left(e^{-\lambda} \frac{\lambda^{x_i}}{x_i!} \right), \\ h(x_1, \dots, x_n|\lambda) &= \ln f(x_1, \dots, x_n|\lambda) = \sum_{i=1}^n \ln \left(e^{-\lambda} \frac{\lambda^{x_i}}{x_i!} \right) \\ &= \sum_{i=1}^n (-\lambda + x_i \ln(\lambda) - \ln(x_i!)) = -n\lambda + n\bar{x} \ln(\lambda) - \sum_{i=1}^n \ln(x_i!) \end{aligned}$$

Da cui

$$\frac{\partial}{\partial \lambda} h(x_1, \dots, x_n|\lambda) = n \left(-1 + \frac{\bar{x}}{\lambda} \right) \geq 0 \iff \lambda \leq \bar{x}.$$

Quindi anche in questo caso lo stimatore di massima verosimiglianza per il parametro λ è la media campionaria \bar{X} (che è uno stimatore corretto).

5.3 Distribuzione gaussiana

In questo caso la densità dipende da due parametri, $\mu \in \mathbb{R}$ e $\sigma > 0$:

$$\begin{aligned} f(x_1, \dots, x_n | \mu, \sigma) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{-(x_i - \mu)^2}{2\sigma^2}\right) \\ &= (2\pi)^{-\frac{n}{2}} (\sigma)^{-n} \exp\left(\frac{-\sum_{i=1}^n (x_i - \mu)^2}{2\sigma^2}\right) \end{aligned}$$

cosicché

$$\begin{aligned} h(x_1, \dots, x_n | \mu, \sigma) &= \ln f(x_1, \dots, x_n | \mu, \sigma) = -\frac{n}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \\ &= -\frac{n}{2} \ln(2\pi) - n \ln(\sigma) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2. \end{aligned}$$

Si ha quindi

$$\begin{aligned} \frac{\partial}{\partial \mu} h(x_1, \dots, x_n | \mu, \sigma) &= \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = n(\bar{x} - \mu), \\ \frac{\partial}{\partial \sigma} h(x_1, \dots, x_n | \mu, \sigma) &= \frac{-n}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^n (x_i - \mu)^2 = \frac{1}{\sigma^3} \left(-n\sigma^2 + \sum_{i=1}^n (x_i - \mu)^2 \right). \end{aligned}$$

Dunque le due derivate parziali si annullano contemporaneamente se e solo se

$$\mu = \bar{x}, \quad \sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{n-1}{n} \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{n-1}{n} s^2.$$

Dunque la media campionaria \bar{X} è uno stimatore di massima verosimiglianza (ed è uno stimatore corretto) per il valore atteso μ mentre $\frac{n-1}{n} S^2$ è uno stimatore di massima verosimiglianza per la varianza σ^2 .

5.4 Distribuzione uniforme su un intervallo

Se (a, b) è l'intervallo, allora la densità del campione è

$$g(x|a, b) = \begin{cases} \frac{1}{b-a} & x \in [a, b], \\ 0 & \text{altrimenti} \end{cases}$$

da cui

$$f(x_1, \dots, x_n | a, b) = \begin{cases} \frac{1}{(b-a)^n} & x_i \in [a, b] \quad \forall i = 1, \dots, n, \\ 0 & \text{altrimenti.} \end{cases}$$

Devo massimizzare $\frac{1}{(b-a)^n}$ con il vincolo $a \leq x_i \leq b$ per ogni $i = 1, \dots, n$. Devo dunque minimizzare la lunghezza dell'intervallo $b - a$ con il vincolo $a \leq x_i \leq b$ per ogni $i = 1, \dots, n$. È dunque

$$a = \min \{x_1, \dots, x_n\}, \quad b = \max \{x_1, \dots, x_n\}.$$

Dunque

$$\min \{X_1, \dots, X_n\}, \quad \max \{X_1, \dots, X_n\}$$

sono stimatori di massima verosimiglianza rispettivamente per l'estremo inferiore e per l'estremo superiore dell'intervallo.

5.5 Efficienza degli stimatori

Sia X_1, \dots, X_n un campione statistico. Supponiamo di conoscere la *forma* della distribuzione ma che essa dipenda da parametro incognito θ (scalare o vettoriale).

Supponiamo che $T = t(X_1, \dots, X_n)$ sia uno stimatore del parametro (scalare) θ . T potrebbe non essere uno stimatore corretto. Definisco *bias dello stimatore θ* la quantità

$$b_\theta(T) := \mathbb{E}[T] - \theta.$$

Valuto l'efficacia di T considerando l'errore quadratico medio

$$\begin{aligned} r(T, \theta) &:= \mathbb{E}[(T - \theta)^2] = \mathbb{E}[(T - \mathbb{E}[T] + b_\theta(T))^2] \\ &= \mathbb{E}[(T - \mathbb{E}[T])^2 + 2b_\theta(T)(T - \mathbb{E}[T]) + b_\theta^2(T)] = \text{Var}[T] + b_\theta^2(T) \end{aligned}$$

Supponiamo che X_1, \dots, X_n sia un campione statistico con distribuzione uniforme sull'intervallo $[0, \theta]$. Sappiamo che

$$\mathbb{E}[X_i] = \frac{\theta}{2}, \quad \text{Var}[X_i] = \frac{\theta^2}{12}.$$

Come stimatore del parametro θ considero allora la statistica $T = 2\bar{X} = \frac{2}{n} \sum_{i=1}^n X_i$. Abbiamo

$$\mathbb{E}[T] = \theta, \quad b_\theta(T) = 0, \quad r(T, \theta) = \text{Var}[T] = 4 \text{Var}[\bar{X}] = \frac{\theta^2}{3n}.$$

Considero ora lo stimatore di massima verosimiglianza

$$\hat{T} := \max(X_1, \dots, X_n).$$

Per calcolarne valore atteso e varianza ne calcoliamo innanzitutto la legge:

$$F_{\hat{T}}(t) = \mathbb{P}(\hat{T} \leq t) = \mathbb{P}(X_1 \leq t, \dots, X_n \leq t) = \prod_{i=1}^n \mathbb{P}(X_i \leq t) = \begin{cases} 0 & t < 0, \\ \left(\frac{t}{\theta}\right)^n & 0 \leq t < \theta, \\ 1 & t \geq \theta. \end{cases}$$

Dunque $\mathbb{P}_{\hat{T}} = g(t)dt$ con

$$g(t) = \begin{cases} n \frac{t^{n-1}}{\theta^n}, & t \in [0, \theta], \\ 0 & \text{altrimenti} \end{cases}$$

Abbiamo dunque

$$\begin{aligned}\mathbb{E}[\widehat{T}] &= \int_{\mathbb{R}} t g(t) dt = \int_0^\theta \frac{n}{\theta^n} t^n dt = \frac{n\theta}{n+1}, \\ \mathbb{E}[\widehat{T}^2] &= \int_{\mathbb{R}} t^2 g(t) dt = \int_0^\theta \frac{n}{\theta^n} t^{n+1} dt = \frac{n\theta^2}{n+2}, \\ \text{Var}[\widehat{T}] &= \frac{n\theta^2}{n+2} - \left(\frac{n\theta}{n+1}\right)^2 = \frac{n\theta^2}{(n+2)(n+1)^2}.\end{aligned}$$

Quindi

$$\begin{aligned}b_\theta(\widehat{T}) &:= \mathbb{E}[\widehat{T}] - \theta = \frac{n\theta}{n+1} - \theta = \frac{-\theta}{n+1}, \\ r(\widehat{T}, \theta) &= \mathbb{E}[(\widehat{T} - \theta)^2] = \text{Var}[\widehat{T}] + b_\theta^2(\widehat{T}) = \frac{n\theta^2}{(n+2)(n+1)^2} + \frac{\theta^2}{(n+1)^2} \\ &= \frac{2\theta^2}{(n+1)(n+2)}\end{aligned}$$

Dunque, per n grande $r(\widehat{T}, \theta) \ll r(T, \theta)$ anche se $\mathbb{E}[\widehat{T}] = \frac{n\theta}{n+1} \neq \theta = \mathbb{E}[T]$ e dunque T è uno stimatore corretto mentre \widehat{T} non lo è.

Cerchiamo ora il *miglior stimatore* nella famiglia $T_c := c\widehat{T}$, $c \in \mathbb{R}$ dove per *miglior stimatore* intendiamo quello che minimizza l'errore quadratico medio $r(T_c, \theta)$. Abbiamo

$$\mathbb{E}[T_c] = c\mathbb{E}[\widehat{T}] = \frac{cn\theta}{n+1}, \quad \text{Var}[T_c] = c^2\text{Var}[\widehat{T}] = \frac{c^2 n\theta^2}{(n+2)(n+1)^2}.$$

Quindi

$$\begin{aligned}b_\theta(T_c) &= \mathbb{E}[T_c] - \theta = \frac{cn\theta}{n+1} - \theta = \frac{\theta(nc - n - 1)}{n+1}, \\ r(T_c, \theta) &= \text{Var}[T_c] + b_\theta^2(T_c) = \frac{c^2 n\theta^2}{(n+2)(n+1)^2} + \frac{\theta^2(nc - n - 1)^2}{(n+1)^2}.\end{aligned}$$

Abbiamo allora

$$\begin{aligned}\frac{\partial}{\partial \theta} r(T_c, \theta) &= \frac{2cn\theta^2}{(n+2)(n+1)^2} + \frac{2n\theta^2(nc - n - 1)}{(n+1)^2} \\ &= \frac{2n\theta^2}{(n+1)^2} \left[\frac{c}{n+2} + (nc - n - 1) \right] = \frac{2n\theta^2}{(n+2)(n+1)^2} [c(n+1)^2 - (n+2)(n+1)] \\ &= \frac{2n\theta^2}{(n+2)(n+1)} [c(n+1) - (n+2)] \stackrel{\geq 0}{\leq 0} \iff c \stackrel{\geq}{\leq} \frac{n+2}{n+1}\end{aligned}$$

Quindi $c \in \mathbb{R} \mapsto r(T_c, \theta)$ ammette minimo assoluto nel punto $c := \frac{n+2}{n+1}$ ed il minimo vale

$$r\left(T_{\frac{n+2}{n+1}}, \theta\right) = \text{Var}[T_b] + b_\theta^2(T_c) = \frac{\left(\frac{n+2}{n+1}\right)^2 n\theta^2}{(n+2)(n+1)^2} + \frac{\theta^2\left(\frac{n+2}{n+1} - n - 1\right)^2}{(n+1)^2} = \frac{\theta^2}{(n+1)^2}$$

6. Intervalli di confidenza

La media campionaria e la varianza campionaria ci offrono una stima dei parametri valore atteso e varianza del campione statistico in esame. Abbiamo però bisogno di sapere *quanto ci si possa fidare di questa stima* ovvero quale sia la probabilità che il *vero* valore del parametro incognito non sia *troppo distante* dalla stima trovata.

Diamo perciò la seguente definizione:

Definizione 6.0.1 (Intervallo di confidenza). Sia X_1, \dots, X_n un campione statistico e sia θ un parametro (ignoto) che caratterizza la distribuzione del campione.

Siano $L_i = l_i(X_1, \dots, X_n)$ e $L_s = l_s(X_1, \dots, X_n)$ due statistiche del campione e sia $\alpha \in (0, 1)$. Dico che l'intervallo (L_i, L_s) è un *intervallo di confidenza* (o di fiducia) di livello $1 - \alpha$ se $\mathbb{P}(\theta \in (L_i, L_s)) \geq 1 - \alpha$, ovvero che (L_i, L_s) è un intervallo di confidenza (o di fiducia) di errore α se $\mathbb{P}(\theta \notin (L_i, L_s)) \leq \alpha$.

Dico che la semiretta $(L_i, +\infty)$ è un *intervallo di confidenza unilaterale superiore* di livello $1 - \alpha$ se $\mathbb{P}(\theta > L_i) \geq 1 - \alpha$

Dico che la semiretta $(-\infty, L_s)$ è un *intervallo di confidenza unilaterale inferiore* di livello $1 - \alpha$ se $\mathbb{P}(\theta < L_s) \geq 1 - \alpha$

Osservazione 6.0.1. 1. La scelta dei nomi delle due statistiche non è casuale: L_i sta per limitazione inferiore mentre L_s sta per limitazione superiore.

2. Di solito si è interessati a *piccoli* valori di α , più precisamente a $\alpha \in (10^{-2}, 10^{-1})$.

3. La disuguaglianza di Chebychev ci ha fornito un intervallo di confidenza per il valore atteso μ del campione nel caso in cui la varianza σ^2 sia nota

$$\mathbb{P}(|\bar{X} - \mu| \geq t) \leq \frac{\sigma^2}{t^2} \quad \forall t > 0$$

ovvero

$$\mathbb{P}(|\bar{X} - \mu| < t) \geq 1 - \frac{\sigma^2}{t^2} \quad \forall t > 0$$

cioè

$$\mathbb{P}(\bar{X} - t < \mu < \bar{X} + t) \geq 1 - \frac{\sigma^2}{t^2} \quad \forall t > 0.$$

Fissato $\alpha \in (0, 1)$ scelgo $t = \frac{\sigma}{\sqrt{\alpha}}$. La disuguaglianza di Chebychev si legge allora

$$\mathbb{P}\left(\bar{X} - \frac{\sigma}{\sqrt{\alpha}} < \mu < \bar{X} + \frac{\sigma}{\sqrt{\alpha}}\right) \geq 1 - \alpha \quad \forall \alpha \in (0, 1).$$

Dunque l'intervallo $\left(\bar{X} - \frac{\sigma}{\sqrt{\alpha}}, \bar{X} + \frac{\sigma}{\sqrt{\alpha}}\right)$ è un intervallo di confidenza di livello $1 - \alpha$ per il valore atteso μ del campione.

6.1 Stima per intervalli del valore atteso di campioni gaussiani

6.1.1 Campione gaussiano di cui è nota la varianza

Intervallo bilaterale

Sia X_1, \dots, X_n un campione gaussiano di valore atteso μ incognita e varianza σ^2 nota.

Sia Z una v.a. gaussiana standard e sia $\alpha \in (0, 1)$. Calcolo $\mathbb{P}\left(|Z| \leq z_{1-\frac{\alpha}{2}}\right)$:

$$\begin{aligned} \mathbb{P}\left(|Z| \leq z_{1-\frac{\alpha}{2}}\right) &= \mathbb{P}\left(-z_{1-\frac{\alpha}{2}} \leq Z \leq z_{1-\frac{\alpha}{2}}\right) \\ &= \mathbb{P}\left(Z \leq z_{1-\frac{\alpha}{2}}\right) - \mathbb{P}\left(Z \leq -z_{1-\frac{\alpha}{2}}\right) = \mathbb{P}\left(Z \leq z_{1-\frac{\alpha}{2}}\right) - \mathbb{P}\left(Z \leq z_{\frac{\alpha}{2}}\right) \\ &= \Phi\left(z_{1-\frac{\alpha}{2}}\right) - \Phi\left(z_{\frac{\alpha}{2}}\right) = 1 - \frac{\alpha}{2} - \frac{\alpha}{2} = 1 - \alpha. \end{aligned} \quad (6.1)$$

Sappiamo che $\mathbb{P}_{\bar{X}} = N\left(\mu, \frac{\sigma^2}{n}\right)$ e che dunque $\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$ ha distribuzione $N(0, 1)$. Applichiamo quindi la disuguaglianza (6.1) a $\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$. Si ha:

$$\begin{aligned} 1 - \alpha &= \mathbb{P}\left(\left|\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}\right| \leq z_{1-\frac{\alpha}{2}}\right) = \mathbb{P}\left(-z_{1-\frac{\alpha}{2}} \leq \frac{\mu - \bar{X}}{\frac{\sigma}{\sqrt{n}}} \leq z_{1-\frac{\alpha}{2}}\right) \\ &= \mathbb{P}\left(\frac{-\sigma z_{1-\frac{\alpha}{2}}}{\sqrt{n}} \leq \mu - \bar{X} \leq \frac{\sigma z_{1-\frac{\alpha}{2}}}{\sqrt{n}}\right) \\ &= \mathbb{P}\left(\bar{X} - \frac{\sigma z_{1-\frac{\alpha}{2}}}{\sqrt{n}} \leq \mu \leq \bar{X} + \frac{\sigma z_{1-\frac{\alpha}{2}}}{\sqrt{n}}\right) \end{aligned}$$

L'intervallo

$$\left(\bar{X} - \frac{\sigma z_{1-\frac{\alpha}{2}}}{\sqrt{n}}, \bar{X} + \frac{\sigma z_{1-\frac{\alpha}{2}}}{\sqrt{n}}\right)$$

è dunque un intervallo di confidenza di livello $1 - \alpha$ per il valore atteso μ del campione.

Osservazione 6.1.1 (Dimensionamento del campione). Fissato il livello di confidenza $1 - \alpha$, supponiamo di voler controllare l'ampiezza dell'intervallo di confidenza $L_s - L_i$. Nel caso in esame l'ampiezza dell'intervallo di confidenza è $\frac{2\sigma z_{1-\frac{\alpha}{2}}}{\sqrt{n}}$. Se fissiamo una limitazione superiore 2δ per l'ampiezza di tale intervallo, deve dunque essere

$$\frac{2\sigma z_{1-\frac{\alpha}{2}}}{\sqrt{n}} \leq 2\delta$$

ovvero

$$n \geq \left(\frac{\sigma z_{1-\frac{\alpha}{2}}}{\delta}\right)^2.$$

Intervallo unilaterale superiore

Sia Z una v.a. tale che $\mathbb{P}_Z = N(0, 1)$. Sappiamo che

$$\mathbb{P}(Z \leq t) = 1 - \alpha \quad \text{se e solo se} \quad t = z_{1-\alpha}.$$

Abbiamo dunque

$$1 - \alpha = \mathbb{P} \left(\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \leq z_{1-\alpha} \right) = \mathbb{P} \left(\bar{X} - \mu \leq \frac{\sigma z_{1-\alpha}}{\sqrt{n}} \right) = \mathbb{P} \left(\mu \geq \bar{X} - \frac{\sigma z_{1-\alpha}}{\sqrt{n}} \right).$$

Quindi la semiretta

$$\left(\bar{X} - \frac{\sigma z_{1-\alpha}}{\sqrt{n}}, +\infty \right)$$

è un intervallo di confidenza unilaterale superiore di livello $1 - \alpha$.

Intervallo unilaterale inferiore

Sia Z una v.a. tale che $\mathbb{P}_Z = N(0, 1)$. Sappiamo che

$$\mathbb{P}(Z \geq t) = 1 - \alpha \quad \text{se e solo se} \quad \mathbb{P}(Z \leq t) = \alpha \quad \text{se e solo se} \quad t = z_\alpha.$$

Abbiamo dunque

$$1 - \alpha = \mathbb{P} \left(\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \geq z_\alpha \right) = \mathbb{P} \left(\bar{X} - \mu \geq \frac{\sigma z_\alpha}{\sqrt{n}} \right) = \mathbb{P} \left(\mu \leq \bar{X} - \frac{\sigma z_\alpha}{\sqrt{n}} \right).$$

Quindi la semiretta

$$\left(-\infty, \bar{X} - \frac{\sigma z_\alpha}{\sqrt{n}} \right) = \left(-\infty, \bar{X} + \frac{\sigma z_{1-\alpha}}{\sqrt{n}} \right)$$

è un intervallo di confidenza unilaterale inferiore di livello $1 - \alpha$.

6.1.2 Campione gaussiano di cui non è nota la varianza

Intervallo bilaterale

Sia X_1, \dots, X_n un campione gaussiano di valore atteso μ varianza σ^2 , entrambe incognite.

Sappiamo che la v.a. $T := \frac{(\bar{X} - \mu)\sqrt{n}}{S}$ segue la distribuzione t di Student con $n - 1$ gradi di libertà:

$$\mathbb{P}_T = t(n - 1).$$

Sia $t_{n-1, 1-\frac{\alpha}{2}}$ il relativo quantile di livello $1 - \frac{\alpha}{2}$:

$$\mathbb{P} \left(T \leq t_{n-1, 1-\frac{\alpha}{2}} \right) = 1 - \frac{\alpha}{2}.$$

Calcolo $\mathbb{P} \left(|T| \leq t_{n-1, 1-\frac{\alpha}{2}} \right)$:

$$\begin{aligned} \mathbb{P} \left(|T| \leq t_{n-1, 1-\frac{\alpha}{2}} \right) &= \mathbb{P} \left(-t_{n-1, 1-\frac{\alpha}{2}} \leq T \leq t_{n-1, 1-\frac{\alpha}{2}} \right) \\ &= \mathbb{P} \left(T \leq t_{n-1, 1-\frac{\alpha}{2}} \right) - \mathbb{P} \left(T \leq -t_{n-1, 1-\frac{\alpha}{2}} \right) \\ &= \mathbb{P} \left(T \leq t_{n-1, 1-\frac{\alpha}{2}} \right) - \mathbb{P} \left(T \leq t_{n-1, \frac{\alpha}{2}} \right) = 1 - \frac{\alpha}{2} - \frac{\alpha}{2} = 1 - \alpha. \end{aligned}$$

Abbiamo dunque

$$\begin{aligned} 1 - \alpha &= \mathbb{P} \left(|T| \leq t_{n-1, 1-\frac{\alpha}{2}} \right) = \mathbb{P} \left(\frac{|\bar{X} - \mu| \sqrt{n}}{S} \leq t_{n-1, 1-\frac{\alpha}{2}} \right) \\ &= \mathbb{P} \left(|\bar{X} - \mu| \leq \frac{S}{\sqrt{n}} t_{n-1, 1-\frac{\alpha}{2}} \right) \\ &= \mathbb{P} \left(\frac{-S}{\sqrt{n}} t_{n-1, 1-\frac{\alpha}{2}} \leq \mu - \bar{X} \leq \frac{S}{\sqrt{n}} t_{n-1, 1-\frac{\alpha}{2}} \right) \\ &= \mathbb{P} \left(\bar{X} - \frac{S}{\sqrt{n}} t_{n-1, 1-\frac{\alpha}{2}} \leq \mu \leq \bar{X} + \frac{S}{\sqrt{n}} t_{n-1, 1-\frac{\alpha}{2}} \right) \end{aligned}$$

L'intervallo

$$\left(\bar{X} - \frac{S t_{n-1, 1-\frac{\alpha}{2}}}{\sqrt{n}}, \bar{X} + \frac{S t_{n-1, 1-\frac{\alpha}{2}}}{\sqrt{n}} \right)$$

è dunque un intervallo di confidenza di livello $1 - \alpha$ per il valore atteso μ del campione.

Intervallo unilaterale superiore

Sappiamo che

$$\mathbb{P}(T \leq t) = 1 - \alpha \quad \text{se e solo se} \quad t = t_{n-1, 1-\alpha}.$$

Abbiamo dunque

$$\begin{aligned} 1 - \alpha &= \mathbb{P} \left(\frac{(\bar{X} - \mu) \sqrt{n}}{S} \leq t_{n-1, 1-\alpha} \right) = \mathbb{P} \left(\bar{X} - \mu \leq \frac{S t_{n-1, 1-\alpha}}{\sqrt{n}} \right) \\ &= \mathbb{P} \left(\mu \geq \bar{X} - \frac{S t_{n-1, 1-\alpha}}{\sqrt{n}} \right). \end{aligned}$$

Quindi la semiretta

$$\left(\bar{X} - \frac{S t_{n-1, 1-\alpha}}{\sqrt{n}}, +\infty \right)$$

è un intervallo di confidenza unilaterale superiore di livello $1 - \alpha$.

Intervallo unilaterale inferiore

Sappiamo che

$$\mathbb{P}(T \geq t) = 1 - \alpha \quad \text{se e solo se} \quad \mathbb{P}(T \leq t) = \alpha \quad \text{se e solo se} \quad t = t_{n-1, \alpha}.$$

Abbiamo dunque

$$1 - \alpha = \mathbb{P} \left(\frac{(\bar{X} - \mu) \sqrt{n}}{S} \geq t_{n-1, \alpha} \right) = \mathbb{P} \left(\bar{X} - \mu \geq \frac{S t_{n-1, \alpha}}{\sqrt{n}} \right) = \mathbb{P} \left(\mu \leq \bar{X} - \frac{S t_{n-1, \alpha}}{\sqrt{n}} \right).$$

Quindi la semiretta

$$\left(-\infty, \bar{X} - \frac{S t_{n-1, \alpha}}{\sqrt{n}} \right) = \left(-\infty, \bar{X} + \frac{S t_{n-1, 1-\alpha}}{\sqrt{n}} \right)$$

è un intervallo di confidenza unilaterale inferiore di livello $1 - \alpha$.

6.2 Stima per intervalli della varianza di campioni gaussiani

Intervallo bilaterale

Sia X_1, \dots, X_n un campione gaussiano di valore atteso μ (incognita o nota) e varianza σ^2 incognita.

Sappiamo che la v.a. $V := (n-1)\frac{S^2}{\sigma^2}$ segue la distribuzione χ^2 a $n-1$ gradi di libertà. Per ogni $\alpha \in (0, 1)$ indico con $\chi_{n-1, \alpha}^2$ il quantile di livello α della v.a. V :

$$F_V(\chi_{n-1, \alpha}^2) = \alpha \quad \forall \alpha \in (0, 1).$$

Osservazione 6.2.1. $\chi_{n-1, \alpha}^2 > 0$ per ogni $\alpha \in (0, 1)$.

Calcolo $\mathbb{P}\left(\chi_{n-1, \frac{\alpha}{2}}^2 < V < \chi_{n-1, 1-\frac{\alpha}{2}}^2\right)$:

$$\begin{aligned} \mathbb{P}\left(\chi_{n-1, \frac{\alpha}{2}}^2 < V < \chi_{n-1, 1-\frac{\alpha}{2}}^2\right) &= \mathbb{P}\left(V < \chi_{n-1, 1-\frac{\alpha}{2}}^2\right) - \\ &\quad - \mathbb{P}\left(V < \chi_{n-1, \frac{\alpha}{2}}^2\right) = 1 - \frac{\alpha}{2} - \frac{\alpha}{2} = 1 - \alpha. \end{aligned}$$

Dunque

$$\begin{aligned} 1 - \alpha &= \mathbb{P}\left(\chi_{n-1, \frac{\alpha}{2}}^2 < (n-1)\frac{S^2}{\sigma^2} < \chi_{n-1, 1-\frac{\alpha}{2}}^2\right) \\ &= \mathbb{P}\left(\frac{1}{\chi_{n-1, 1-\frac{\alpha}{2}}^2} < \frac{\sigma^2}{(n-1)S^2} < \frac{1}{\chi_{n-1, \frac{\alpha}{2}}^2}\right) = \mathbb{P}\left(\frac{(n-1)S^2}{\chi_{n-1, 1-\frac{\alpha}{2}}^2} < \sigma^2 < \frac{(n-1)S^2}{\chi_{n-1, \frac{\alpha}{2}}^2}\right) \end{aligned}$$

Quindi l'intervallo

$$\left(\frac{(n-1)S^2}{\chi_{n-1, 1-\frac{\alpha}{2}}^2}, \frac{(n-1)S^2}{\chi_{n-1, \frac{\alpha}{2}}^2}\right)$$

è un intervallo di confidenza di livello $1 - \alpha$ per la varianza σ^2 del campione.

Intervallo unilaterale superiore

Sappiamo che

$$\mathbb{P}(V \leq t) = 1 - \alpha \quad \text{se e solo se} \quad t = \chi_{n-1, 1-\alpha}^2.$$

Dunque

$$1 - \alpha = \mathbb{P}\left((n-1)\frac{S^2}{\sigma^2} < \chi_{n-1, 1-\alpha}^2\right) = \mathbb{P}\left(\sigma^2 > (n-1)\frac{S^2}{\chi_{n-1, 1-\alpha}^2}\right).$$

Quindi la semiretta

$$\left(\frac{(n-1)S^2}{\chi_{n-1, 1-\alpha}^2}, +\infty\right)$$

è un intervallo di confidenza di livello $1 - \alpha$ per la varianza σ^2 del campione.

Intervallo unilaterale inferiore

Sappiamo che

$$\mathbb{P}(V \geq t) = 1 - \alpha \quad \text{se e solo se} \quad \mathbb{P}(V \leq t) = \alpha \quad \text{se e solo se} \quad t = \chi_{n-1, \alpha}^2.$$

Dunque

$$1 - \alpha = \mathbb{P}\left((n-1)\frac{S^2}{\sigma^2} > \chi_{n-1, \alpha}^2\right) = \mathbb{P}\left(\sigma^2 \leq (n-1)\frac{S^2}{\chi_{n-1, \alpha}^2}\right).$$

Quindi l'intervallo

$$\left(0, \frac{(n-1)S^2}{\chi_{n-1, \alpha}^2}\right)$$

è un intervallo di confidenza di livello $1 - \alpha$ per la varianza σ^2 del campione.

Esempio 6.2.1. Calcoliamo gli intervalli di confidenza per il carattere Totpor dei dati tratti da [2], nell'ipotesi che si tratti della realizzazione di v.a. normali.

```
> setwd("~/Documents/didattica/2017-18_analisi_reale/alcuni_appunti/esempio_statistica")
>
> library(readr)
>
> table2 <- read_delim("~/Documents/didattica/2017-18_analisi_reale/alcuni_appunti/
table2.csv", "\t", escape_double = FALSE, trim_ws = TRUE)
Parsed with column specification:
cols(
  Code = col_character(),
  Totpor = col_double(),
  PRA = col_double(),
  PV = col_double(),
  Densi = col_double(),
  TenStr = col_double(),
  CO2SBW = col_double(),
  FirTemp = col_integer()
)
>
> ## definisco la funzione che calcola l'intervallo bilaterale con varianza nota
>
> bilat.norm = function(x, sigma, conf) { n = length(x); xbar=mean(x);
+ alpha = 1 - conf;
+ zstar = qnorm(1-alpha/2);
+ SE = sigma/sqrt(n);
+ xbar + c(-zstar*SE, zstar*SE)}
>
> # definisco la funzione che calcola l'intervallo bilaterale con varianza ignota
>
> bilat.stud = function(x, conf) { n = length(x);
+ m = n-1;
+ xbar=mean(x);
+ alpha = 1 - conf;
+ zstar = qt(1-alpha/2, m, lower.tail=TRUE);
```

```

+ SE = sd(x)/sqrt(n);
+ xbar + c(-zstar*SE,zstar*SE)
+ }
>
> # definisco la funzione che calcola l'intervallo bilaterale per la varianza
>
> bilat.chi = function(x,conf) {
+   n = length(x);
+   m = n-1;
+   alpha = 1 - conf;
+   zsup = qchisq(alpha/2, m, lower.tail=TRUE);
+   zinf = qchisq(1 - alpha/2, m, lower.tail=TRUE);
+   SE = sd(x)*sd(x)*m;
+   c(SE/zinf,SE/zsup)
+ }
>
>
> numSummary(table2[,c("Totpor", "PRA", "PV", "Densi", "TenStr", "CO2SBW", "FirTemp")],
+ statistics=c("mean", "sd", "quantiles"), quantiles=c(0,.25,.5,.75,1))
      mean      sd      0%      25%      50%      75%      100%  n NA
Totpor  40.1193548  7.0371760  26.850  36.0550  40.900  44.4200  54.640 31  0
PRA      0.6732581  0.4760389   0.158   0.4220   0.622   0.7305   2.657 31  0
PV       55.3290323 28.5498417  10.200  30.4500  59.400  80.7000  88.600 31  0
Densi    1.6929032  0.1701214   1.340   1.5600   1.680   1.8150   2.020 31  0
TenStr   0.6092258  0.3143682   0.143   0.4065   0.527   0.7165   1.405 31  0
CO2SBW   0.5816667  0.5259152   0.050   0.2900   0.390   0.4950   1.960 30  1
FirTemp 764.8387097 52.9698636 730.000 740.0000 740.000 750.0000 960.000 31  0
>
> bilat.norm(table2$Totpor, 7.04, .9)
[1] 38.03957 42.19914
> bilat.norm(table2$Totpor, 7.04, .95)
[1] 37.64113 42.59758
>
> bilat.stud(table2$Totpor, .9)
[1] 37.97416 42.26455
> bilat.stud(table2$Totpor, .95)
[1] 37.53810 42.70061
>
> bilat.chi(table2$Totpor, .9)
[1] 33.94002 80.33757
> bilat.chi(table2$Totpor, .95)
[1] 31.62366 88.48047
>

```

