

STATISTICA

venerdì 18 agosto 2023 09:39

Statistica descrittiva

analisi dei dati X_1, \dots, X_n e della loro distribuzione, senza l'interpretazione di un modello probabilistico
(non ci interessa la distribuzione di prob. associata all'esperimento)

Indici statistici

media campionaria $\rightarrow \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

mediana campionaria \rightarrow ordine in senso crescente, n dispari è la data centrale, n pari è la media dei 2 dati centrali

varianza campionaria $\rightarrow S^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$ indice di dispersione

coefficiente di correlazione campionario $\rightarrow r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\left(\sum_{i=1}^n (x_i - \bar{x})^2\right)^{\frac{1}{2}} \left(\sum_{i=1}^n (y_i - \bar{y})^2\right)^{\frac{1}{2}}} \in [-1, 1]$ indica quanto i dati sono allineati
 $|r| = 0 \Rightarrow$ poco allineati
 $|r| \approx 1 \Rightarrow$ dati molto vicini a una retta.

} Indici di centralità

retta di regressione campionario $\rightarrow y = \alpha^* x + \beta^*$ con $\alpha^* = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$, $\beta^* = \bar{y} - \alpha^* \bar{x}$
è la retta che meglio approssima i dati

Statistica inferenziale

A partire dai dati x_1, \dots, x_n ottenuto da un esperimento ripetuto ricavare info sulla distrib. di prob.

Statistica inferenziale parametrica

Supp che la distrib. dell'esp. sia nota a meno di uno o più parametri. (Siamo interessati a determinare uno o più parametri incogniti)

Modello Statistico parametrico

È una terna $(S, J, (Q_\theta)_{\theta \in \Theta})$

- (S, J) : spazio misurabile
- Θ : insieme dei parametri $\neq \emptyset$
- Q_θ prob su (S, J) $\forall \theta \in \Theta$
rappresenta la distribuzione di un carattere di un esperimento aleatorio

esempio

Lancio di moneta

$$\left. \begin{array}{l} \theta = p = \text{prob di testa} \\ \Theta = [0,1] \end{array} \right\} \theta \in \Theta$$

$$(S, \mathcal{J}) = ([0,1], \mathcal{P}([0,1]))$$

$$Q_\theta = B(\theta) \text{ Bernoulli di parametro } \theta$$

Campione

Sia $(S, \mathcal{J}, (Q_\theta)_{\theta \in \Theta})$ modello statistico.

Si ottiene campione ogni successione $(X_n)_{n \geq 1}$ di v.a. reali definite su S

$\exists (P_\theta)_{\theta \in \Theta}$ fam. di prob, t.c. $\forall \theta \in \Theta$ le X_i sono i.i.d. sotto P_θ con legge Q_θ

(con distribuzione che può dipendere da θ)

(X_1, X_2, \dots, X_n) è detta campione di taglia n e legge Q_θ .

Tipicamente

Se Q_θ è la distribuzione di un esperimento aleatorio X_1, \dots, X_n rappresentano gli esiti di n ripetizioni

es

n lanci di moneta

$$X_i = \begin{cases} 1 & \text{se esce T all'i-esimo lancio} \\ 0 & \text{altrimenti} \end{cases}$$

$$\Omega = \{0,1\}^n \quad \mathcal{J} = \mathcal{P}(\Omega) \quad P_\theta = Q_\theta^{\otimes n} \quad \text{con } P_\theta(x_1, \dots, x_n) = q_\theta(x_1) \dots q_\theta(x_n)$$

Nella statistica inferenziale osservo i risultati x_1, \dots, x_n di un campione X_1, \dots, X_n di legge Q_θ e voglio studiare il parametro incognito θ a partire da questi risultati

Dato una fam. di prob. $(Q_\theta)_{\theta \in \Theta}$ su $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ tutte discrete o tutte ass. cont.

\exists un campione (X_1, \dots, X_n) di legge Q_θ :

$$\Omega = \mathbb{R}^n, \quad \mathcal{J} = \mathcal{B}(\mathbb{R}^n)$$

$$X_i(w_1, \dots, w_n) = w_i$$

P_θ caso discreto $P_\theta(x_1, \dots, x_n) = q_\theta(x_1) \dots q_\theta(x_n)$

caso ass. cont. $f_\theta(x_1, \dots, x_n) = g_\theta(x_1) \dots g_\theta(x_n)$

caso generale $Q_\theta = P_\theta^{\otimes n}$

Statistica campionaria

Sia $(X_n)_{n \geq 1}$ un campione.

Una successione di v.a. $(Y_n)_{n \geq 1}$ della forma $Y_n = g_n(X_1, \dots, X_n)$

dove $g_n: \mathbb{R}^n \rightarrow \mathbb{R}$ è una f. misurabile si chiama statistica campionaria.

Stimatore

è una statistica che viene chiamata stimatore per via del ruolo speciale che ricopre.

Infatti uno stimatore è una statistica che non dipende dirett. da θ .

Lo scopo di uno stimatore è stimare $R(\theta)$ con $R: \Theta \rightarrow \mathbb{R}$ f. data

Criteri per vedere la bontà di uno stimatore

Stimatore corretto

Uno stimatore U si dice stimatore corretto di $R(\theta)$ se

- $\forall \theta \in \Theta$ U è p^θ -integrabile equiv. $U_n \in L^1(S, \mathcal{J}, (Q_\theta)_{\theta \in \Theta})$
- $E^\theta[U] = R(\theta) \quad \forall \theta \in \Theta$
↳ E^θ è il v. atteso rispetto alla prob. Q_θ

slogan "la media di U su tutti i campioni possibili è $R(\theta)$ "

Stimatore asintoticamente corretto

Una successione di stimatori $U_n = g_n(X_1, \dots, X_n)$ di $R(\theta)$

è asintoticamente corretto non distorta se $\forall n \in \mathbb{N}^+, \forall \theta \in \Theta$,

- U_n è p_θ -integrabile
- $\lim_n E^\theta[U_n] = R(\theta) \quad \forall \theta \in \Theta$

Mi serve indebolire la def di correttezza [così non $\forall n$ ma solo nel $\lim_{n \rightarrow \infty}$]

perché non è sempre possibile determinare stimatori con le

proprietà desiderate

Stimatore consistente

Una successione di stimatori $U_n = g_n(X_1, \dots, X_n)$ di $R(\theta)$ è

consistente se:

$$\lim_n P^\theta\{|U_n - R(\theta)| > \varepsilon\} = 0 \quad \forall \varepsilon > 0 \quad \forall \theta \in \Theta$$

cioè $(U_n)_n$ converge in p^θ -prob a $R(\theta)$

rischio quadratico

Dato U stimatore di $R(\theta)$

$$R_{\theta}(U) = E^{\theta} [(U - R(\theta))^2] = \text{rischio quadratico di } U$$

preferibilita'

Dati U e V stimatori di $R(\theta)$. U e' preferibile a V se:

$$R_{\theta}(U) \leq R_{\theta}(V) \quad \forall \theta \in \Theta \quad \text{cioe' } U \text{ e' meno incerto di } V$$

Oss.

$$\text{Se } U \text{ e' corretto } R_{\theta}(U) = \text{Var}^{\theta}(U)$$

Stimatori di massima verosomiglianza

notazione

Q_{θ} discreta con den. discreta $m_{\theta} = p_{\theta}$
 Q_{θ} ass. con " " " $m_{\theta} = f_{\theta}$

funzione di verosomiglianza

$$L: \Theta \times \mathbb{R}^n \rightarrow \mathbb{R}$$

$$L(\theta, x_1, \dots, x_n) = L_{\theta}(x_1, \dots, x_n) = m_{\theta}(x_1) \dots m_{\theta}(x_n)$$

caso discreto

$$L_{\theta}(x_1, \dots, x_n) = P_{\theta}(X_1 = x_1) \dots P_{\theta}(X_n = x_n) = P_{\theta}(X_1 = x_1, \dots, X_n = x_n)$$

caso ass. cont.

L_{θ} e' la densita' congiunta di un campione i.i.d. X_1, \dots, X_n di legge Q_{θ}

stimatore di massima verosomiglianza

U e' stimatore di massima verosomiglianza (MLE) di θ se:

$$L(U(x_1, \dots, x_n), x_1, \dots, x_n) = \sup_{\theta \in \Theta} L_{\theta}(x_1, \dots, x_n)$$

perche' $U = g(x_1, \dots, x_n)$ trovare U di max ver $\Leftrightarrow g$ t.c. $L_{g(x_1, \dots, x_n)}(x_1, \dots, x_n) = \sup_{\theta \in \Theta} L_{\theta}(x_1, \dots, x_n)$

Attenzione si usa se Θ interv. aperto di \mathbb{R} / a misura o \mathbb{R} stesso.

Per poter calcolare lo stimatore MLE si ricorre alla funzione log-verosomiglianza ottenuta attraverso l'applicazione del logaritmo naturale, quindi risulta: $l(x, \theta) = \log L(x, \theta)$

Dato che la funzione logaritmica e' una trasformazione monotona crescente, con il passaggio alla log-verosomiglianza non si perdono le caratteristiche della funzione $L(x, \theta)$ in termini di crescita e decrescenza e soprattutto si ottiene una forma analitica piu' semplice da trattare.

Non si usa per $U(\theta, \theta)$

es. log-verosomiglianza

$$Q_{\theta} = B(\theta) \quad \theta \in (0, 1)$$

$$m_{\theta}(x) = \begin{cases} \lambda - \theta & \text{se } x=0 \\ \theta & \text{se } x=1 \end{cases} = \theta^x (1-\theta)^{1-x} \quad x \in \{0, 1\}$$

$$\omega_\theta = \mathcal{D}(\theta) \quad \theta \in (0,1)$$

$$m_\theta(x) = \begin{cases} \theta^{1-x} & \text{se } x=0 \\ \theta & \text{se } x=1 \end{cases} = \theta^x (1-\theta)^{1-x} \quad x \in \{0,1\}$$

$$L_\theta(x_1, \dots, x_n) = \prod_{i=1}^n m_\theta(x_i) = \prod_{i=1}^n \theta^{x_i} (1-\theta)^{1-x_i} = \theta^{\sum_{i=1}^n x_i} (1-\theta)^{n - \sum_{i=1}^n x_i} = \theta^{n\bar{x}} (1-\theta)^{n(1-\bar{x})}$$

$$\text{con } \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

Cerco l'MLE con log-likelihood.

$$\log L_\theta(x_1, \dots, x_n) = \log(\theta^{n\bar{x}} (1-\theta)^{n(1-\bar{x})}) = \log(\theta^{n\bar{x}}) + \log((1-\theta)^{n(1-\bar{x})}) = n\bar{x} \log(\theta) + n(1-\bar{x}) \log(1-\theta)$$

Cerco il max della f.

$$\frac{d}{d\theta} \log L_\theta(x_1, \dots, x_n) = \frac{n\bar{x}}{\theta} + \frac{n(1-\bar{x})}{1-\theta} \quad n(\bar{x} - \bar{x}\theta + \theta)$$

$$\frac{d}{d\theta} = 0 \quad (\Rightarrow) \quad \frac{n\bar{x}(1-\theta) + n\theta(1-\bar{x})}{\theta(1-\theta)} = \frac{n\bar{x} - n\bar{x}\theta + n\theta - n\theta\bar{x}}{\theta(1-\theta)} = 0 \quad (\Rightarrow) \quad \bar{x} = \theta$$

Inoltre $\log L_\theta(x_1, \dots, x_n) \rightarrow -\infty$ per $\theta = 0, 1 \Rightarrow \exists!$ max in $\bar{x} = \theta$

$\Rightarrow \bar{x}$ è l'MLE cercato

Mo dello esponenziale

$(\mathbb{R}, \mathcal{B}(\mathbb{R}), (\mathcal{Q}_\theta)_{\theta \in \Theta}) \quad \Theta \subseteq \mathbb{R}$ intervallo t.c.

discreto $\exists T: \mathbb{N} \rightarrow \mathbb{R}, g: \mathbb{N} \rightarrow \mathbb{R}$ t.c. $\forall \theta \in \Theta \quad \mathcal{Q}_\theta$ ha den. discr. su \mathbb{N}

$$p_\theta(k) = c_\theta g(k) e^{\theta T(k)} \quad k \in \mathbb{N} \quad c_\theta > 0 \text{ cost.}$$

ass. cont. $\exists T: \mathbb{R} \rightarrow \mathbb{R}, g: \mathbb{R} \rightarrow \mathbb{R}$ boreliana t.c. $\forall \theta \in \Theta \quad \mathcal{Q}_\theta$ ha den.

$$f_\theta(x) = c_\theta g(x) e^{\theta T(x)} \quad x \in \mathbb{R} \quad c_\theta > 0 \text{ cost.}$$

Teorema consistenza dello stim. di MLE

$(\mathbb{R}, \mathcal{B}(\mathbb{R}), (\mathcal{Q}_\theta)_{\theta \in \Theta})$ modello statistico t.c.:

1. $\forall \theta_1 \neq \theta_2 \quad \mathcal{Q}_{\theta_1} \neq \mathcal{Q}_{\theta_2}$ (D libro)
2. $\Theta \subseteq \mathbb{R}$ interv. aperto (A libro)
3. $(\mathcal{Q}_\theta)_{\theta \in \Theta}$ è mo dello esp. ass. cont. $\rightarrow f_\theta(x) = c_\theta g(x) e^{\theta T(x)}$ (B+C) libro
4. $x \mapsto g(x) e^{\theta T(x)}$ è integrabile $\forall \theta \in \Theta$ (E libro) $\Leftrightarrow \int_{\mathbb{R}} c_\theta (g(x) e^{\theta T(x)})^2 dx < +\infty$
 $f(x, t) = g(x) e^{tT(x)}$ caso cont.
5. Sia (X_1, \dots, X_n) campione iid. di legge \mathcal{Q}_θ , e supp che \exists Un MLE a valori in Θ

\Rightarrow Allora $\hat{\theta}_n$ è unico e consistente

Lemma

Y_n v. z. reali $Y_n \xrightarrow{P} c, \varphi: \mathbb{R} \rightarrow \mathbb{R}$ cont. $\Rightarrow \varphi(Y_n) \xrightarrow{P} \varphi(c)$

Intervallo di fiducia

$(\mathcal{S}, \mathcal{J}, (\mathcal{Q}_\theta)_{\theta \in \Theta})$ modello statistico
 (X_1, \dots, X_n) campione iid di legge \mathcal{Q}_θ

regione di fiducia

Dato $\alpha \in (0,1)$, una regione di fiducia a livello $1-\alpha$ per θ è un insieme aleatorio $D(\omega) \subseteq \Theta, \omega \in \Omega$
 (cioè una mappa $\Omega \rightarrow \mathcal{P}(\Theta)$
 $\omega \mapsto D(\omega)$)

Dato $\alpha \in (0,1)$, una regione di fiducia a livello $1-\alpha$ per θ è un insieme aleatorio $D(\omega) \subseteq \Theta$, $\omega \in \Omega$
(cioè una mappa $\Omega \rightarrow \mathcal{P}(\Theta)$)
 $\omega \mapsto D(\omega)$)

tale che $P_{\theta}(\theta \in D) \geq 1-\alpha \quad \forall \theta \in \Theta$
 $= \{ \omega \in \Omega \mid \theta \in D(\omega) \} \subseteq \Omega$
 $\{ \theta \in \Theta \} \subseteq \mathcal{P}(\Theta)$

dato un campione, posso determ D te, con prob alta, $\theta \in D$

Oss: $\{ \theta \in D \}$ (l'aleatorietà è in D non in θ che è deterministico)

$P(\theta \in D) \geq 1-\alpha$ va inteso come in almeno una frazione $1-\alpha$ di tutti i campioni, θ cade in D

A questo intervallo di confidenza si associa quindi un valore di probabilità cumulativa che caratterizza, indirettamente in termini di probabilità, la sua ampiezza rispetto ai valori massimi assumibili dalla variabile aleatoria.

Cioè il valore di probabilità cumulativa indica la probabilità che l'evento casuale descritto dalla variabile aleatoria cada all'interno di suddetto intervallo di confidenza, graficamente pari all'area sottesa dalla curva di distribuzione di probabilità della variabile aleatoria nell'intervallo considerato.

È bene non confondere l'intervallo di confidenza con la probabilità.

Perciò l'espressione "vi è un livello di confidenza del 95% che μ sia nell'intervallo", non indica la probabilità che μ cada nell'intervallo, in quanto μ non è una variabile aleatoria nella logica frequentista (μ è invece interpretata come una costante non nota); bensì, indica che nel 95% dei casi in cui questa tecnica viene adottata, questa produce un intervallo che contiene il [valore vero](#) di μ

- La **stima puntuale** è lo specifico valore assunto da una statistica, calcolata in corrispondenza dei dati campionari e che viene utilizzata per stimare il vero valore non noto di un parametro di una popolazione
- Uno **stimatore per intervallo** è un intervallo costruito attorno allo stimatore puntuale, in modo tale che sia nota e fissata la probabilità che il parametro appartenga all'intervallo stesso
- Tale probabilità è detta **livello di confidenza** ed è in generale indicato con $(1-\alpha)\%$ dove α è la probabilità che il parametro si trovi al di fuori dell'intervallo di confidenza
- Quindi la confidenza è il grado di fiducia che l'intervallo possa contenere effettivamente il parametro di interesse

Quantili

Dato $\beta \in (0,1)$, β -quantile di P (o di F d $R F$) il numero $r_\beta = \inf \{x \in \mathbb{R} \mid F(x) \geq \beta\}$

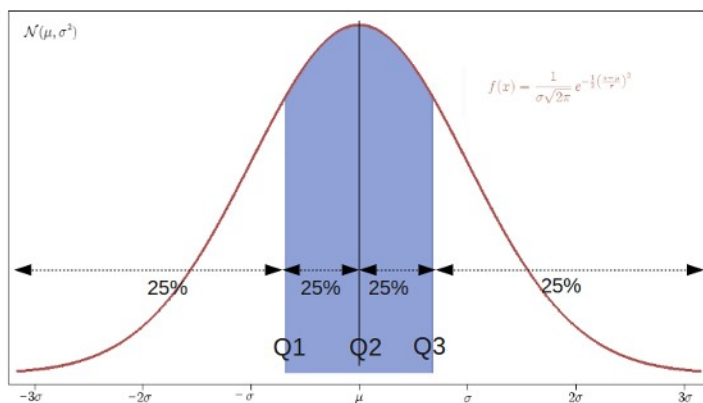
I quantili di $\mathcal{N}(0,1)$ si denotano con q_β, z_β

Per simmetria 1 $q_{1-\beta} = -q_\beta$

2 $P\{-q_{1-\frac{\alpha}{2}} \leq z \leq q_{1-\frac{\alpha}{2}}\} = 1-\alpha$ per $Z \sim \mathcal{N}(0,1), 0 < \alpha < 1$

Queste due proprietà valgono in generale se Z ha densità f pari, è più in generale se $Z \stackrel{(d)}{=} -Z$

In statistica il **quantile** di ordine α o α -quantile (con α un numero reale nell'intervallo $[0,1]$) è un valore q_α che divide la popolazione in due parti, proporzionali ad α e $(1-\alpha)$ e caratterizzate da valori rispettivamente minori e maggiori di q_α . Per poter calcolare un quantile di ordine α è necessario che il carattere sia almeno ordinato, cioè sia possibile definire un ordinamento sui valori



Densità di probabilità di una distribuzione normale con quantili in evidenza. L'area sotto la curva rossa è la stessa negli intervalli $(-\infty, Q_1)$, (Q_1, Q_2) , (Q_2, Q_3) e $(Q_3, +\infty)$

Nel caso di una densità di probabilità la funzione di ripartizione F è continua e il quantile di ordine α è definito da $F(q_\alpha) = \alpha$. Questo quantile può non essere unico se la funzione di densità è nulla in un intervallo, ovvero se la funzione di ripartizione è costante ed assume il valore α per più di un valore q_α ; ciononostante per ognuno di questi valori la distribuzione viene correttamente divisa in due parti proporzionali ad α e $(1-\alpha)$, in quanto un intervallo a densità nulla non contribuisce al calcolo della probabilità, quindi non fa differenza quale punto dell'intervallo si scelga come q_α .

Nel caso di una densità discreta il quantile di ordine α è un valore q_α nel quale la somma delle probabilità discrete sia maggiore o uguale ad α , ovvero tale che la somma delle probabilità *fino a* quel valore incluso sia almeno α e che la somma delle probabilità discrete da quel valore in poi (incluso) sia maggiore o uguale a $1-\alpha$. Nel caso discreto, oltre alla non unicità del quantile, si può avere una divisione della distribuzione non proporzionale ad α e $1-\alpha$ (del resto una variabile discreta può essere divisa solo in un numero discreto di modi).

Il quantile di ordine $\alpha=0,1$ (anche detto *primo decile*) è quel valore della distribuzione per cui la probabilità cumulata fino a quel valore, incluso, sia maggiore o uguale a 0,1, e la probabilità cumulata da quel valore, incluso, in poi sia maggiore o uguale a 0,9. Il quantile di ordine $\alpha=0,5$ (la *mediana*) è quel valore della distribuzione per cui la probabilità cumulata fino a quel valore, incluso, sia maggiore o uguale a 0,5, e la probabilità cumulata da quel valore, incluso, in poi sia maggiore o uguale a 0,5.

Modello della Statistica pivotale

Data una statistica pivotale cioè $ST = g(\theta, X, n)$

1) ST è invertibile come funzione di θ

2) ST ha legge che non dipende da θ cioè $P\{ST \in A\}$ non dipende da θ

Allora ha una regione di fiducia di livello $1-\alpha$ per θ è data da

$$D = g(\cdot, x_1, \dots, x_n)^{-1}(A) \text{ dove } A \text{ è t.c. } P\{ST \in A\} = 1-\alpha$$

$$\text{infatti } P\{\theta \in D\} = P\{ST \in g(\cdot, x_1, \dots, x_n)(\theta)\} = 1-\alpha$$

regione di fiducia asintotica

$D \in (q, 1)$, una regione di fiducia asintotica di livello $(1-\alpha) \in (0, 1)$

è una successione di insiemi aleatori $D_n(\omega) \subseteq \Theta, \omega \in \Omega, n \in \mathbb{N}$ t.c.

$$\liminf P_{\theta}\{\theta \in D_n\} \geq 1-\alpha$$

distribuzioni legate alle v.a. Gaussiane

(Ω, \mathcal{F}, P)

Proposizione

a) $Z \sim N(0, 1) \Rightarrow Z^2 \sim \Gamma\left(\frac{1}{2}, \frac{1}{2}\right) = \chi_1^2$

b) $Z_1, \dots, Z_n \sim N(0, 1) \Rightarrow Z_1^2 + \dots + Z_n^2 \sim \Gamma\left(\frac{n}{2}, \frac{1}{2}\right) = \chi_n^2$

c) $Z = (z_1, \dots, z_n)^T$ vettore di $N(0, 1)$ indep, $\|Z\|^2 \sim \chi_n^2$ vettore di $N(0, 1)$ indep

d) $Z_1, \dots, Z_n \sim N(0, 1)$ indep. $\Rightarrow \bar{Z}, S^2$ indep
 $\bar{Z} \sim N\left(0, \frac{1}{n}\right), (n-1)S^2 \sim \chi_{n-1}^2$

distribuzione t - di student

Date $Z \sim N(0, 1), Y \sim \chi_n^2, Z, Y$ indep

distribuzione t - di student a mgrodi di libertà (t.n) la legge di $T = \frac{\sqrt{n}Z}{Y}$

$$f_{t_n}(x) = \frac{\Gamma\left(\frac{n+1}{2}\right)}{\sqrt{n\pi} \Gamma\left(\frac{n}{2}\right)} \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}}$$

Corollario

$$Z_1, \dots, Z_n \sim N(0, 1) \text{ i.i.d.} \Rightarrow \sqrt{n} \frac{\bar{Z}}{S} \sim t_{n-1}$$

Prop

$$X_1, \dots, X_n \text{ i.i.d.} \sim N(m, \sigma^2) \quad \text{con } \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i, S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

a) \bar{X}, S^2 indep

b) $\bar{X} \sim N\left(m, \frac{\sigma^2}{n}\right), (n-1) \frac{S^2}{\sigma^2} \sim \chi_{n-1}^2$

c) $\sqrt{n} \frac{\bar{X} - m}{S} \sim t_{n-1}$

Esempi notevoli di stimatori corretti

- La media campionaria $\bar{X} = \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ è un buon stimatore di $E(X_1)$

$$\text{infatti } E[\bar{X}] = \frac{1}{n} \sum_{i=1}^n E[X_i] = E^*[X_1]$$

- La varianza $S^2 = S_{n-1}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{1}{n-1} \left(\sum_{i=1}^n X_i^2 - n \bar{X}^2 \right)$ è media campionaria degli scarti quadratici ed è un buon stimatore di $\text{Var}(X_1)$

$$\text{infatti } E[\bar{X}^2] = \frac{1}{n^2} \sum_{i,j=1}^n E^*[X_i X_j] = \frac{1}{n^2} \sum_{i=1}^n E^*[X_i^2] + \frac{1}{n^2} \sum_{j \neq i} E^*[X_i] E^*[X_j] = \frac{1}{n} E^*[X_1^2] + \frac{n-1}{n} E^*[X_1]^2$$

$$E^*[S^2] = \frac{1}{n-1} \left(\sum_{i=1}^n E^*[X_i^2] - n E^*[\bar{X}^2] \right) = \frac{1}{n-1} \left(n E^*[X_1^2] - E^*[X_1^2] - (n-1) E^*[X_1]^2 \right) = E^*[X_1^2] - E^*[X_1]^2 = \text{Var}^*(X_1)$$

- Il coeff. di correl. campionario $r^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\left(\sum_{i=1}^n (X_i - \bar{X})^2 \right)^{1/2} \left(\sum_{i=1}^n (Y_i - \bar{Y})^2 \right)^{1/2}}$

è un buon stimatore del coeff. di correl. tra le X_i e Y_i $\rho(X_i, Y_i)$

momento secondo $z+p$ $E^*[X_i^2]$ $z+p$ b. stim per $E(X)$

quarto $E^*[X_i^4]$ $z+p$ " $\text{Var}^*(X)$

esempi di modelli esponenziali

leggi esponenziali $\rightarrow f_{\lambda}(x) = \lambda e^{-\lambda x} \mathbb{1}_{(0,+\infty)}(x)$

$$\begin{aligned}\theta &= \lambda \\ T(x) &= -x \\ g(x) &= \mathbb{1}_{(0,+\infty)}(x) \\ \zeta &= \lambda\end{aligned}$$

leggi di Poisson $\rightarrow p_{\lambda}(k) = \frac{\lambda^k}{k!} e^{-\lambda}$

$$\begin{aligned}\theta &= \log(\lambda) \\ T(k) &= k \\ g(k) &= \frac{1}{k!} \\ \zeta &= e^{-\lambda} \\ \lambda^k &= e^{k \log(\lambda)}\end{aligned}$$

Leggi Gaussiane $\mathcal{N}(m, \sigma^2)$ $\sigma^2 = 1$ per semplicità $f_m(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-m)^2}{2}}$

$$\begin{aligned}\theta &= m \\ T(x) &= 2x \\ g(x) &= e^{-x^2} \\ \zeta &= \frac{1}{\sqrt{2\pi}} e^{-m^2}\end{aligned}$$

leggi geometriche $p_p(k) = p(1-p)^{k-1} \mathbb{1}_{\mathbb{N}_{k \geq 1}}$

$$\begin{aligned}\theta &= \log(1-p) \\ T(k) &= k-1 \\ g(k) &= \mathbb{1}_{\mathbb{N}_{k \geq 1}} \\ \zeta &= p\end{aligned}$$

Stimatore di max verosimiglianza di media e var Gaussiana

$(\mathbb{R}, \mathcal{B}(\mathbb{R}), \mathcal{N}(m, \sigma^2))_{m \in \mathbb{R}, \sigma^2 \in (0, +\infty)}$ $\theta = \mathbb{R} \times (0, +\infty)$

$L_{(m, \sigma^2)}(x_1, \dots, x_n) = f_{(m, \sigma^2)}(x_1) \dots f_{(m, \sigma^2)}(x_n) = \frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - m)^2}$

passo al log.

$$\begin{cases} \frac{\partial}{\partial m} = 0 \\ \frac{\partial}{\partial \theta} = 0 \end{cases}$$

e ottengo che $(\bar{x}, \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2)$ è MLE per (m, σ^2)

se σ^2 è nota (ovè $\theta = m$) $\Rightarrow \bar{x}$ è MLE per m

se m è nota (ovè $\theta = \sigma^2$) $\Rightarrow \frac{1}{n} \sum_{i=1}^n (x_i - m)^2$ è MLE per σ^2

INTERVALLI DI FIDUCIA

media di una popolazione normale (varianza nota)

media di una popolazione Bernoulli (grandi campioni)

media di una popolazione, varianza non nota (grandi campioni)

media di una popolazione normale (varianza non nota)

varianza di una popolazione Gaussiana

Intervalli di fiducia per la media di una popolazione normale (varianza nota)

$$(\mathbb{R}, \mathcal{B}(\mathbb{R}), (\mathcal{N}(m, \sigma^2))_{m \in \mathbb{R}})$$

Poiché $\bar{X} = \frac{1}{n} \sum x_i$ è uno stimatore di m è ragionevole cercare D nell'intervallo centrato in \bar{X}

$$D = [\bar{X} - d, \bar{X} + d]$$

Sappiamo che $\bar{X} \sim \mathcal{N}(m, \frac{\sigma^2}{n})$ cioè $\frac{\sqrt{n}}{\sigma}(\bar{X} - m) \sim \mathcal{N}(0, 1)$

calcoliamo $P_n\{m \in [\bar{X} \pm d]\}$

$$P_n\{m \in [\bar{X} \pm d]\} = P_n\{|\bar{X} - m| \leq d\} = P_n\left\{\left|\frac{\sqrt{n}}{\sigma}(\bar{X} - m)\right| \leq \frac{d\sqrt{n}}{\sigma}\right\} = P(|Z| \leq \frac{d\sqrt{n}}{\sigma}) = \Phi\left(\frac{d\sqrt{n}}{\sigma}\right) - \Phi\left(-\frac{d\sqrt{n}}{\sigma}\right) = 2\Phi\left(\frac{d\sqrt{n}}{\sigma}\right) - 1 \quad \text{con } Z \sim \mathcal{N}(0, 1)$$

Imponiamo $P_n\{m \in [\bar{X} \pm d]\} = 1 - \alpha$ per avere il più piccolo intervallo possibile

$$2\Phi\left(\frac{d\sqrt{n}}{\sigma}\right) - 1 = 1 - \alpha \Rightarrow \Phi\left(\frac{d\sqrt{n}}{\sigma}\right) = 1 - \frac{\alpha}{2} \Rightarrow \frac{d\sqrt{n}}{\sigma} = q_{1 - \frac{\alpha}{2}}$$

↳ quantile

Quindi $[\bar{X} \pm \frac{\sigma}{\sqrt{n}} q_{1 - \frac{\alpha}{2}}]$ è i.c. di fiducia per m di livello $1 - \alpha$

Per TCL $\lim_n P\left\{-q_{1 - \frac{\alpha}{2}} \leq \frac{\sqrt{n}}{\sigma}(\bar{X}_n - m) \leq q_{1 + \frac{\alpha}{2}}\right\} = 1 - \alpha$, l'inter. di prima è asintotico

Intervalli di fiducia per la media di una popolazione normale (varianza non nota)

$$(\mathbb{R}, \mathcal{B}(\mathbb{R}), (\mathcal{N}(m, \sigma^2))_{m \in \mathbb{R}})$$

$$St = \frac{\sqrt{n}}{S}(\bar{X} - m) \sim t_{n-1} \quad (\text{t-di-Student, } n-1 \text{ gradi di libertà})$$

chiamo $t_{\beta, n-1}$ il quantile di t_{n-1} di ordine β ($P\left\{\frac{T}{t_{n-1}} \leq t_{\beta, n-1}\right\} = \beta$)

$$P\{St \in [-t_{1 - \frac{\alpha}{2}, n-1}, t_{1 - \frac{\alpha}{2}, n-1}]\} = 1 - \alpha$$

Otengo $[\bar{X} \pm \frac{S}{\sqrt{n}} t_{1 - \frac{\alpha}{2}, n-1}]$ è inter. di fiducia esatto di livello $1 - \alpha$ per m

proprietà.

$t_{\beta, n-1} > q_{\beta} \quad \forall n, \forall \beta > \frac{1}{2} \Rightarrow$ incertezza nel caso var non nota > incertezza caso var nota

Intervalli di fiducia asintotici per la media di una popolazione (varianza non nota), grandi campioni

$$(\mathbb{R}, \mathcal{B}(\mathbb{R}), (\mathcal{Q}_m)_m), \quad m = E_m(X), \quad X \sim \mathcal{Q}_m$$

$$P_{1 - \alpha} \left(\frac{\sqrt{n}(\bar{X}_n - m)}{\sqrt{V_n}} \right) \leq \dots$$

$$(\mathbb{R}, \mathcal{B}(\mathbb{R}), (\mathcal{Q}_m)_m), \quad m = E_m(X), \quad X \sim \mathcal{Q}_m$$

Problema $\frac{\sqrt{n}}{\sigma} (\bar{X}_n - m)$ è funzione anche di σ non nota

Però vale il corollario

$$E[X_i^2] < +\infty \Rightarrow \frac{\sqrt{n}}{S_n} (\bar{X}_n - m) \xrightarrow{\text{in legge}} Z \sim N(0,1)$$

Uso $ST = \frac{\sqrt{n}}{S_n} (\bar{X}_n - m) \Rightarrow \left[\bar{X}_n \pm \frac{\sqrt{n}}{S_n} q_{1-\frac{\alpha}{2}} \right]$ è i. fiducia asintotico (di livello $1-\alpha$) per m .

Intervalli di fiducia per la varianza di una popolazione normale

$$(\mathbb{R}, \mathcal{B}(\mathbb{R}), (\mathcal{N}(m, \sigma^2))_{\sigma^2 \in (0, +\infty)})$$

m non nota

$$\text{Uso } ST = (n-1) \frac{S^2}{\sigma^2} \sim \chi^2_{n-1}$$

chiamo $\chi^2_{\beta, n-1}$ il quantile di ordine β di χ^2_{n-1} ($P\{Y \leq \chi^2_{\beta, n-1}\} = \beta$)

$$P\{ST \in [\chi^2_{\frac{\alpha}{2}, n-1}, +\infty)\} = 1-\alpha$$

$\left[0, \frac{(n-1)S^2}{\chi^2_{\frac{\alpha}{2}, n-1}} \right]$ è int. di fiducia unilatero per σ^2 di livello $1-\alpha$

Intervallo di fiducia per una proporzione (media di una popol. di Bernoulli)

$$(\mathbb{R}, \mathcal{B}(\mathbb{R}), (\mathcal{B}(p))_{p \in (0,1)})$$

$$p = E_p(X) \quad X \sim \mathcal{B}(p)$$

$$\text{Uso } ST = \frac{n}{\bar{X}_n(1-\bar{X}_n)} (\bar{X}_n - p)$$

$\left[\bar{X}_n \pm \sqrt{\frac{\bar{X}_n(1-\bar{X}_n)}{n}} q_{1-\frac{\alpha}{2}} \right]$ è int. di fiducia

Test Statistici

Cosa sono i test statistici?

Un test statistico è una procedura che consente di verificare con un elevato grado di fiducia una nostra ipotesi iniziale H_0 , che si chiama anche ipotesi di lavoro o ipotesi nulla, inerente a un fenomeno che stiamo studiando. Il test essenzialmente è una procedura di calcolo che si basa sui risultati dei dati numerici che abbiamo a disposizione, dati che sono interpretati come valori osservati di una certa variabile aleatoria. Tale procedura di calcolo si conclude fornendo un giudizio, un consenso, in inglese *test statistic* [1], che permette di confermare oppure di respingere l'ipotesi di lavoro, in favore di un'altra ipotesi iniziale H_1 , detta ipotesi alternativa. Essendo la procedura aleatoria, possono comunque presentarsi quattro casi diversi:

	H_0 è vera	H_0 è falsa
Accetto H_0	H_0 è vera ed accetto H_0	H_0 è falsa ma accetto H_0
Rifiuto H_0	H_0 è vera ma rifiuto H_0	H_0 è falsa e rifiuto H_0

Le celle evidenziate con lo sfondo colorato rappresentano i casi in cui commettiamo degli errori: commettiamo un errore di primo tipo quando l'ipotesi H_0 è vera ma noi la rifiutiamo. È possibile anche a valutare la probabilità di commettere tale errore: tale probabilità di commettere si indica con α , ed essa viene chiamata livello di significatività [2] del test. Esiste anche una probabilità β di commettere un errore di secondo tipo, ossia l'errore di accettare l'ipotesi H_0 nonostante essa sia falsa. In linea di principio, uno sperimentatore desidera abbassare il livello α il più possibile verso zero; sfortunatamente, si può dimostrare che quanto compari il fatto che più abbassiamo di conseguenza la probabilità β a di riflesso questo diminuirà il numero [3], che viene definito con il termine di potenza del test (costa, la probabilità di rifiutare l'ipotesi di lavoro giacché essa è falsa). Riflettendosi a quanto avevamo detto in [2] a proposito dei test (diagnostici), la potenza di un test [4] è legata alla sensibilità ed ai valori predittivi positivi, mentre la quantità [5] si collega alla specificità ed ai valori predittivi negativi. Per illustrare ancora meglio la situazione, possiamo osservare questi due schemi (avendo in mente un processo diagnostico):

risoluzione costante	innocente	colpevole
	innocente risultato innocente condannato	colpevole risultato colpevole condannato
negativo positivo	malato	malato
	vero negativo falso positivo	falso negativo vero positivo

Quello che si vorrebbe realizzare è di non commettere l'errore di primo tipo ("condannare un innocente", "curare un sano"), inoltre si vuole anche abbassare la possibilità di commettere l'errore di secondo tipo ("assolvere il colpevole", "non curare un malato") allo scopo di massimizzare la potenza del test ("condannare il colpevole", "curare un malato"). Si noti però che c'è una sorta di "asimmetria" nei due errori. Tra i molti test statistici esistenti, il più celebre in ambito medico sicuramente è il test t di Student. Ci occuperemo di trattare in una prossima dispensa anche un altro importante test, il test χ^2 di indipendenza di Karl Pearson.

Test statistici

Un test statistico è una procedura per verificare, sulla base dei dati del campione se una certa H_0 è plausibile o no
affermazione sul parametro

$(S, J, (Q_\theta)_{\theta \in \Theta})$ modello statistico, (X_1, \dots, X_n) iid di legge Q_θ , def. su $(\Omega, \mathcal{F}, (P_\theta)_{\theta \in \Theta})$

Elementi principali di un test (parametrico)

1. L' H_0 da verificare: $\Theta = \Theta_0 \cup \Theta_1$ partizione di Θ ($\Theta_0 \cap \Theta_1 = \emptyset$)
di carattere

- $H_0: \theta \in \Theta_0$ ipotesi nulla, di cui verificare se è plausibile o no
- $H_1: \theta \in \Theta_1$ ipotesi alternativa

Osservazione

- Se l'esito del test, sulla base dei dati, è il rifiuto di H_0 , in favore di H_1 , allora H_0 è poco plausibile, cioè c'è evidenza statistica per H_1
- Se l'esito del test è l'accettazione di $H_0 \Rightarrow H_0$ è plausibile, ma non c'è per forza evidenza per H_0

2. La regione critica o di rifiuto

Un evento $C \in \mathcal{F}$, solit. dipendente da (X_1, \dots, X_n) è:

se l'esito w cade in C rifiutiamo l' H_0 in favore di H_1

se l'esito w non cade in C , accettiamo l' H_0

La scelta di C deve essere tale che: C sia poco prob. sotto H_0
se $w \in C$ allora questa è un' evidenza per H_1

esempio

n lanci di moneta con $p = \text{prob. testa}$

Il banco dichiara che $p = \frac{1}{2}$

Sulla base degli esiti degli n lanci, vogliamo verificare se l' $H_0: p = \frac{1}{2}$ è plausibile o no

$H_0: p = \frac{1}{2}$

mi aspetto una regione critica della forma $C = \{|\bar{X} - \frac{1}{2}| > d\}$ per un d opportuno, per cui $P_{\frac{1}{2}}(C)$ sia piccola

$H_0: p = \frac{1}{2}$

$C = \{|\bar{X} - \frac{1}{2}| > d\}$ per cui $P_p(C)$ sia piccola $\forall p < \frac{1}{2}$

Spesso per C si usa una statistica ST la cui legge non dipende da σ .

Errori

errore di 1^a specie: H_0 vera ma l'esito del test è il rifiuto detto

errore di 2^a specie: H_0 falsa ma l'esito del test è l'accettazione di H_0

Prob di commettere questi errori

$d \in (0, 1)$. Un test è di livello d se $\sup_{\theta \in \Theta_0} P_\theta(C) \leq d$

Cioè la prob. dell'errore di 1^a specie è al massimo d

Più d è piccolo, maggiore è l'evidenza statistica fornita dal test contro H_0 in caso di rifiuto

potenza di un test

$\pi_C: \Theta \rightarrow [0, 1]$ $\pi_C(\theta) = P_\theta(C)$
↓

$$\pi_c: \Theta_1 \rightarrow [0,1] \quad \pi_c(\theta) = P_\theta(C)$$

↓
prob. di rifiutare H_0 quando il v. del parametro è $\theta \in \Theta_1$

$\pi_c(\theta)$ rappresenta la capacità di accorgersi che l'ip. H_0 è falsa

Più $\pi_c(\theta)$ è piccolo più commetti errore di 2^a specie
più α è piccolo più " " " 1^a specie

Approccio classico.

Fissa α piccolo e cerco C in modo da massimizzare la potenza

Legame tra test e intervalli di fiducia

caso ip semplice $\Theta_0 = \{\theta_0\}$ e $\Theta_1 = \Theta \setminus \{\theta_0\}$

- dato $\alpha \in (0,1)$ se D è una regione di fiducia di livello $1-\alpha$ per $\theta \Rightarrow$

$C = \{\theta_0 \notin D\} = \{\omega \in \Omega \mid \theta_0 \notin D(\omega)\}$ è una regione critica di livello α
perché
 $P_{\theta_0}(C) = 1 - P_{\theta_0}(D) = 1 - (1-\alpha) = \alpha$

Vi versa

- dato $\alpha \in (0,1)$ se $\forall \theta_0 \in \Theta$, C_{θ_0} è una regione critica di livello α per il test $H_0: \theta = \theta_0 \Rightarrow$

$D(\omega) = \{\theta \in \Theta \mid \omega \notin C_\theta\}$ è una regione di fiducia per θ di livello $1-\alpha$

perché
 $P_{\theta_0}(\theta_0 \in D) = P_{\theta_0}(C_{\theta_0}^c) = 1 - P_{\theta_0}(C_{\theta_0}) \geq 1 - \alpha$

Metodo del p-value / soglie di accettazione

Abbiamo una famiglia di regioni critiche $(C_\alpha)_{\alpha \in (0,1)}$

C_α regione critica a livello $\alpha \forall \alpha$

$$\alpha < \alpha' \Rightarrow C_\alpha \subset C_{\alpha'}$$

$$\bigcap_{\alpha} C_\alpha = \emptyset \quad \bigcup_{\alpha} C_\alpha = \Omega$$

Dato $\omega \in \Omega$ esiste $\exists \bar{\alpha} = \bar{\alpha}(\omega) \in (0,1)$ t.c.

- $\alpha < \bar{\alpha} \Rightarrow \omega \notin C_\alpha \Rightarrow$ accetto H_0
- $\alpha > \bar{\alpha} \Rightarrow \omega \in C_\alpha \Rightarrow$ rifiuto H_0

$\bar{\alpha}(\omega)$ fornisce il grado di plausibilità di H_0 se si verifica ω

più $\bar{\alpha}$ è basso meno plausibile è H_0 , minore è il livello α per cui rifiuto

Quando si effettua un test d'ipotesi si fissa un'ipotesi nulla e un valore soglia α (per convenzione di solito 0,05) che indica il **livello di significatività** del test. Calcolato il *p-value* relativo ai dati osservati è possibile comportarsi come segue:

- se valore $p > \alpha$ l'evidenza empirica non è sufficientemente contraria all'ipotesi nulla che quindi non può essere rifiutata;
- se valore $p \leq \alpha$ l'evidenza empirica è fortemente contraria all'ipotesi nulla che quindi va rifiutata. In tal caso si dice che i dati osservati sono statisticamente significativi.

- se valore $p \leq \alpha$ l'evidenza empirica è fortemente contraria all'ipotesi nulla che quindi va rifiutata. In tal caso si dice che i dati osservati sono statisticamente significativi.

Caso H_p composta

$$\left. \begin{array}{l} \theta_0 = (-\infty, \theta_0] \quad H_0: \theta \leq \theta_0 \\ \theta_1 = (\theta_0, +\infty) \quad H_1: \theta > \theta_0 \end{array} \right\} \text{ o viceversa}$$

def

Dati $(S, J, (Q_\theta)_{\theta \in \Theta})$ modello statistico con $\Theta \subset \mathbb{R}$ intervallo

(X_1, \dots, X_n) campione i.i.d. di legge Q_θ def su $(\Omega, \mathcal{F}, (P_\theta)_{\theta \in \Theta})$

$T: \Omega \rightarrow \mathbb{R}$ v.a.

diciamo che il modello è a rapporto di verosimiglianza crescente rispetto a T

se $\forall \theta_1 < \theta_2 \exists f_{\theta_1, \theta_2}: \mathbb{R} \rightarrow (0, +\infty)$ strett. cresc. fr.

$$\frac{L(\theta_2, X_1, \dots, X_n)}{L(\theta_1, X_1, \dots, X_n)} = f_{\theta_1, \theta_2}(T)$$

Teorema

Nelle H_p sopra considero il Test $H_0: \theta \leq \theta_0$ contro $H_1: \theta > \theta_0$

sia $C = \{T \geq d\}$

Allora

a) $\sup_{\theta \leq \theta_0} P_\theta(C) = P_{\theta_0}(C)$

b) il Test di regione critica C è il più potente tra i test di livello $P_{\theta_0}(C)$

Lemma

$\Theta = [\theta_0, \theta_1]$ con $H_0: \theta = \theta_0$ e $H_1: \theta = \theta_1$,

$$C = \{L_{\theta_0}(X_1, \dots, X_n) \leq c L_{\theta_1}(X_1, \dots, X_n)\} \quad c > 0$$

Allora

a) il Test di regione critica C è il più potente tra i test di livello $P_{\theta_0}(C)$

b) $P_{\theta_0}(C) \leq P_{\theta_1}(C)$

TEST 2

Test Bilatero per la media di una popol. gaussiana (var nota)

Dati ↓

$(\mathbb{R}, \mathcal{B}(\mathbb{R}), \mathcal{N}(m, \sigma^2)_{m \in \mathbb{R}})$, $\theta = m \in \Theta = \mathbb{R}$, σ^2 nota (X_1, \dots, X_n) campione i.i.d. di $\mathcal{N}(m, \sigma^2)$

Test ↓

dato $m_0 \in \mathbb{R}$, considero il Test $H_0: m = m_0$ vs $H_1: m \neq m_0$, a livello α

Regione critica C ↓

Poiché \bar{X} è uno stimatore di m e $H_1: m \neq m_0$, cerco C della forma

$$C = \{|\bar{X} - m_0| \geq d\}$$

TEST 2

Test unilatero per la media di una popol. gaussiana (var nota)

Dati ↓

$(\mathbb{R}, \mathcal{B}(\mathbb{R}), \mathcal{N}(m, \sigma^2)_{m \in \mathbb{R}})$, $\theta = m \in \Theta = \mathbb{R}$, σ^2 nota (X_1, \dots, X_n) campione i.i.d. di $\mathcal{N}(m, \sigma^2)$

Test ↓

dato $m_0 \in \mathbb{R}$, considero il Test $H_0: m \leq m_0$ vs $H_1: m > m_0$, a livello α

Regione critica C ↓

Per il Teorema $C = \{\bar{X} > d\}$ è ottimale d opportuno

Poiché \bar{X} è uno stimatore di m e $H_1: m \neq m_0$, cerco C della forma

$$C = \{ |\bar{X} - m_0| > d \} \text{ d opportuno}$$

Procedura ↓

Uso $ST = \frac{\sqrt{n}}{\sigma} (\bar{X} - m) \sim N(0,1)$ Sotto H_0

$$\text{Calcolo } P_{m_0}(C) = P_{m_0} \left\{ \left| \frac{\sqrt{n}}{\sigma} (\bar{X} - m_0) \right| > \frac{\sqrt{n}}{\sigma} d \right\} = P \left\{ |Z| > \frac{\sqrt{n}}{\sigma} d \right\}$$

Impongo $P_{m_0}(C) = \alpha$ cioè $\frac{\sqrt{n}}{\sigma} d = q_{1-\frac{\alpha}{2}}$

$$\text{Ottengo } C = \left\{ \left| \frac{\sqrt{n}}{\sigma} (\bar{X} - m_0) \right| > q_{1-\frac{\alpha}{2}} \right\} = \left\{ |\bar{X} - m_0| > \frac{\sigma}{\sqrt{n}} q_{1-\frac{\alpha}{2}} \right\}$$

↳ regione critica di livello $1-\alpha$

Test Bilatero per la media di una popol. Bernoulli, grandi camp

Dati ↓

$(\mathbb{R}, \mathcal{B}(\mathbb{R}), \mathcal{B}(\mathbb{P}))_{p \in (0,1)}$, $\theta = p \in \Theta = (0,1)$, (X_1, \dots, X_n) campione iid. di $B(p)$

Test ↓

dato $p_0 \in (0,1)$, considero il Test $H_0: p = p_0$ vs. $H_1: p \neq p_0$, a livello α

Regione critica C ↓

$$C = \{ |\bar{X} - p_0| > d \} \text{ d opportuno}$$

Procedura ↓

Uso $ST = \sqrt{\frac{n}{p_0(1-p_0)}} (\bar{X} - p_0) \stackrel{TC}{\sim} N(0,1)$ Sotto H_0

Impongo $P_{p_0}(C) = \alpha$ cioè $\frac{\sqrt{n}}{\sigma} d = q_{1-\frac{\alpha}{2}}$

$$\text{Ottengo } C = \left\{ \left| \frac{\sqrt{n}}{\sqrt{p_0(1-p_0)}} (\bar{X} - p_0) \right| > q_{1-\frac{\alpha}{2}} \right\} = \left\{ \frac{\sqrt{n} |\bar{X} - p_0|}{\sqrt{p_0(1-p_0)}} > q_{1-\frac{\alpha}{2}} \right\}$$

↳ regione critica di livello α

Per il Teorema $C = \{ |\bar{X} - m_0| > d \}$ è ottimale d opportuno

Procedura ↓

Uso $ST = \frac{\sqrt{n}}{\sigma} (\bar{X} - m) \sim N(0,1)$ Sotto H_0 uso statistica di Test.

Impongo $\sup_{m \in m_0} P_m(\bar{X} > d) = P_{m_0}(\bar{X} > d) = \alpha$

$$\alpha = P_{m_0}(\bar{X} > d) = P_{m_0} \left\{ \frac{\sqrt{n}}{\sigma} (\bar{X} - m_0) > \frac{\sqrt{n}}{\sigma} d \right\} = P(Z > \frac{\sqrt{n}}{\sigma} d)$$

$$\Leftrightarrow \frac{\sqrt{n}}{\sigma} d = q_{1-\alpha}$$

$$C = \left\{ \frac{\sqrt{n}}{\sigma} (\bar{X} - m_0) > q_{1-\alpha} \right\}$$

p-value

$$\bar{\alpha} = P \left(Z > \frac{\sqrt{n}}{\sigma} (\bar{X} - m_0) \right)$$

potenza

$$P_m(C) = P_m \left\{ \frac{\sqrt{n}}{\sigma} (\bar{X} - m_0) > q_{1-\alpha} - \frac{\sqrt{n}}{\sigma} (m - m_0) \right\} =$$

$$= P \left(Z > q_{1-\alpha} - \frac{\sqrt{n}}{\sigma} (m - m_0) \right) = 1 - \Phi \left(q_{1-\alpha} - \frac{\sqrt{n}}{\sigma} (m - m_0) \right)$$

Test unilatero per la media di una popol. gaussiana (var nota)

Dati ↓

$(\mathbb{R}, \mathcal{B}(\mathbb{R}), \mathcal{B}(\mathbb{P}))_{p \in (0,1)}$, $\theta = p \in \Theta = (0,1)$, (X_1, \dots, X_n) campione iid.

Test ↓

dato $p_0 \in \mathbb{R}$, considero il Test $H_0: p \leq p_0$ vs. $H_1: p > p_0$, a livello α

Regione critica C ↓

Per il Teorema $C = \{ \bar{X} > d \}$ è ottimale d opportuno

Procedura ↓

Uso $ST = \sqrt{\frac{n}{p_0(1-p_0)}} (\bar{X} - p_0) \stackrel{TC}{\sim} N(0,1)$ uso statistica di Test.

Ottengo come sopra

$$C = \left\{ \frac{\sqrt{n}}{\sqrt{p_0(1-p_0)}} (\bar{X} - p_0) > q_{1-\alpha} \right\} \text{ livello } \alpha$$

p-value

$$\bar{\alpha} = P \left(Z > \frac{\sqrt{n}}{\sqrt{p_0(1-p_0)}} (\bar{X} - p_0) \right)$$

In un test unilaterale la zona di rifiuto dell'ipotesi nulla è solamente in una coda della distribuzione; in un test bilaterale essa è equamente divisa nelle due code della distribuzione.