

# INTERVALLI di FIDUCIA (CONFIDENCE INTERVAL)

Consideriamo un campione  $X_1, \dots, X_n$  (variabili aleatorie i.i.d.),  
con legge dipendente da un parametro  $\theta \in \Theta \subseteq \mathbb{R}$ ,  
scriveremo  $P_\theta$  per indicare la legge e di conseguenza  
 $F_\theta(t) = P_\theta(X_1 \leq t)$ ,  $\int \theta \sigma p_\theta$  se le v.a. hanno densità o sono discrete

Cerchiamo, a partire da dati  $x_1, \dots, x_n$ , un intervallo numerico in cui  
ci aspettiamo di trovare  $\theta$

Definizione: un intervallo di fiducia per  $\theta$  di LIVELLO  $1-\alpha$  ( $\alpha \in (0,1)$ , tipicamente piccolo)  
è un intervallo aleatorio della forma  $I = [a(X_1, \dots, X_n), b(X_1, \dots, X_n)]$   
tale che  $P_\theta(\theta \in I) = P(a(X_1, \dots, X_n) \leq \theta \leq b(X_1, \dots, X_n)) \geq 1-\alpha$

Osservazione: nella ricerca di un intervallo di fiducia abbiamo due richieste in antitesi:

- Vogliamo che  $I$  sia piccolo (stima precisa di  $\theta$ )
- Vogliamo che  $\alpha$  sia piccolo (probabilità bassa che  $I$  non contenga affatto  $\theta$ )

## INTERVALLI di FIDUCIA per GAUSSIANE

Sia  $X_1, \dots, X_n$  campione di legge  $N(m, \sigma^2)$  che dipende solo dai parametri  $(m, \sigma)$ .

- $\sigma$  NOTO: vogliamo trovare un intervallo di fiducia per  $m \in \mathbb{R}$

sia  $\bar{X}_n = \frac{X_1 + \dots + X_n}{n}$  uno stimatore consistente per  $m$

Nota: se  $X_1, \dots, X_n$  è ragionevole cercare  $I$  nella forma  $[\bar{X}_n - d, \bar{X}_n + d]$  con  $d$  da determinare  
Imponiamo della definizione:

$$\begin{aligned} 1-\alpha &\leq P_m(m \in I) = P_m(\bar{X}_n - d \leq m \leq \bar{X}_n + d) \\ &= P_m(|\bar{X}_n - m| \leq d) = P_m\left(\left|\sqrt{n} \frac{\bar{X}_n - m}{\sigma}\right| \leq \frac{\sqrt{n} d}{\sigma}\right) = \\ &= P\left(|Z| \leq \frac{\sqrt{n} d}{\sigma}\right) = 2\Phi\left(\frac{\sqrt{n} d}{\sigma}\right) - 1 \end{aligned}$$

allora  $\bar{X}_n$  ha legge  
 $N(m, \sigma^2/n)$  e  
 $Z = \frac{\bar{X}_n - m}{\sigma/\sqrt{n}}$  è  $N(0,1)$

Note: il  $\beta$ -quantile  
della distribuzione  $N(0,1)$   
è il numero  $q_\beta$  tale che  
 $\Phi(q_\beta) = \beta$

Abbiamo quindi  $1-\alpha \leq \Phi\left(\frac{\sqrt{n}d}{\sigma}\right) - 1$   
e prendiamo l'uguaglianza  $\Phi\left(\frac{\sqrt{n}d}{\sigma}\right) = \frac{2-\alpha}{2} = 1 - \frac{\alpha}{2}$   
 $\frac{\sqrt{n}d}{\sigma} = q_{1-\frac{\alpha}{2}} \Rightarrow d = \frac{\sigma}{\sqrt{n}} \cdot q_{1-\frac{\alpha}{2}}$

Dunque

$$I = \left[ \bar{X}_n - \frac{\sigma}{\sqrt{n}} q_{1-\frac{\alpha}{2}}, \bar{X}_n + \frac{\sigma}{\sqrt{n}} q_{1-\frac{\alpha}{2}} \right]$$

è un intervallo di fiducia per  $\mu$  a livello  $1-\alpha$

Esempio: descriviamo l'altezza di un italiano scelto a caso con una  
variabile aleatoria  $N(\mu, \sigma^2)$  di cui supponiamo di conoscere  
 $\sigma = 5$  cm. Date misurazioni  $X_1, \dots, X_n$  con  $\bar{X}_n = 170$  cm  
a livello 95% ( $\alpha = 0.05$ ).

Ho un intervallo di fiducia per  $\mu$  del tipo:

$$\left[ 170 - \frac{5}{\sqrt{n}} q_{0.975}, 170 + \frac{5}{\sqrt{n}} q_{0.975} \right]$$

Osservazione:  $d$  è anche detta "precisione della stima":

(\*)

- all'aumentare di  $n$ ,  $d$  decresce
- all'aumentare di  $\sigma$ ,  $d$  cresce
- al diminuire di  $\alpha$ ,  $q_{1-\frac{\alpha}{2}}$  aumenta e così  $d$

- $\sigma$  NON NOTO: stimiamo  $\sigma$  con  $S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$  che è uno stimatore  
CORRETTO e CONSISTENTE. Ci serve:

Fatto: se  $X$  e  $C_n$  sono variabili aleatorie indipendenti con densità  
rispettivamente  $N(0,1)$  e  $\chi^2(n)$

allora  $T_n = \sqrt{n} \frac{X}{\sqrt{C_n}}$  ha densità  $f_{T_n}(t) = \frac{1}{\sqrt{n\pi}} \frac{\Gamma(\frac{n+1}{2})}{\Gamma(\frac{n}{2})} \left(1 + \frac{t^2}{n}\right)^{-\frac{n}{2}-\frac{1}{2}}$   
che chiamiamo densità di Student con  $n$  gradi di libertà

Note: per  $n \rightarrow \infty$  la variabile di Student con  $n$  gradi di  
libertà converge in distribuzione alla Gaussiana  $N(0,1)$

Fatto: se  $X_1, \dots, X_n$  sono un campione  $N(\mu, \sigma^2)$

1)  $\bar{X}_n$  e  $S_n^2$  sono INDIPENDENTI

2) La densità  $N(\mu, \frac{\sigma^2}{n})$  e  $\frac{n-1}{\sigma^2} S_n^2$  ha densità  $\chi^2(n-1)$

3)  $T_n = \sqrt{n} \frac{\bar{X}_n - \mu}{S_n}$  è di Student con  $n-1$  gradi di libertà

Cerchiamo l'intervallo nella forma  $[\bar{X}_n - d, \bar{X}_n + d]$

$$1 - \alpha \leq P_{m, \sigma^2} (m \in I) = P_{m, \sigma^2} (|\bar{X}_n - m| \leq d) = \\ = P_{m, \sigma^2} \left( \left| \sqrt{n} \frac{\bar{X}_n - m}{S_n} \right| \leq \frac{\sqrt{n} d}{S_n} \right)$$

Indichiamo con  $\tau_{\beta, m}$  il  $\beta$ -quantile delle variabili di Student  $T_m$ , ovvero  $F_{T_m}(\tau_{\beta, m-1}) = \beta$ .

Siccome la densità  $f_{T_m}$  è simmetrica vale  $F_{T_m}(t) + F_{T_m}(-t) = 1$

$$= F_{T_m} \left( \frac{\sqrt{n} d}{S_n} \right) - F_{T_m} \left( -\frac{\sqrt{n} d}{S_n} \right) = \\ = 2 F_{T_m} \left( \frac{\sqrt{n} d}{S_n} \right) - 1$$

$$\Downarrow \\ F_{T_m} \left( \frac{\sqrt{n} d}{S_n} \right) = 1 - \frac{\alpha}{2} \quad \frac{\sqrt{n} d}{S_n} = \tau_{1-\frac{\alpha}{2}, m-1}$$

Nota: quando  $n$  è generale  $\tau_{\beta, m}$  è indistinguibile da  $q_{\beta}$

$$d = \frac{S_n}{\sqrt{n}} \tau_{1-\frac{\alpha}{2}, m-1}$$

Dunque  $I = \left[ \bar{X}_n - \frac{S_n}{\sqrt{n}} \tau_{1-\frac{\alpha}{2}, m-1}, \bar{X}_n + \frac{S_n}{\sqrt{n}} \tau_{1-\frac{\alpha}{2}, m-1} \right]$

Esempio: analogamente al precedente se  $\bar{X}_n = 170 \text{ cm}$

$$\text{e } \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{X}_n)^2 = (5 \text{ cm})^2$$

allora

$$\left[ 170 \pm \frac{5}{\sqrt{n}} \tau_{0.975, n-1} \right] \text{ è di fiducia} \\ \text{con } \alpha = 0,05$$

Valgono le osservazioni di (\*)

## INTERVALLO di FIDUCIA per VARIABILI ALEATORIE BERNOULLI ( $p$ )

Consideriamo un campione  $X_1, \dots, X_n$  con legge Bernoulli ( $p$ ).

Sappiamo che  $\bar{X}_n$  è uno stimatore corretto e costante di  $p$ , quindi cerchiamo

$$I = [\bar{X}_n \pm d]$$

$$1 - \alpha = P_p(p \in I) = P_p(|\bar{X}_n - p| \leq d)$$

Sappiamo per TCL che  $\sqrt{n} \frac{\bar{X}_n - p}{\sqrt{p(1-p)}}$  è approssimativamente  $N(0, 1)$

ma non possiamo dividere per  $\sqrt{p(1-p)}$  poiché  $p$  non è noto. Si può però mostrare che:

$$\sqrt{n} \frac{\bar{X}_n - p}{\sqrt{p(1-p)}} = \sqrt{n} \frac{\bar{X}_n - p}{\sqrt{\bar{X}_n(1-\bar{X}_n)}} \xrightarrow{n \rightarrow \infty} N(0, 1)$$

$$\hookrightarrow P_p\left(\left|\sqrt{n} \frac{\bar{X}_n - p}{\sqrt{\bar{X}_n(1-\bar{X}_n)}}\right| \leq \sqrt{n} \frac{d}{\sqrt{\bar{X}_n(1-\bar{X}_n)}}\right)$$

$$\approx P_p(|Z| \leq \frac{\sqrt{n} d}{\sqrt{\bar{X}_n(1-\bar{X}_n)}})$$

Abbiamo quindi la condizione approssimata:

$$1 - \alpha \approx 2 \Phi\left(\frac{\sqrt{n} d}{\sqrt{\bar{X}_n(1-\bar{X}_n)}}\right) - 1 \quad d \approx \frac{\sqrt{\bar{X}_n(1-\bar{X}_n)}}{\sqrt{n}} q_{1 - \frac{\alpha}{2}}$$

dunque

$$[\bar{X}_n \pm \frac{\sqrt{\bar{X}_n(1-\bar{X}_n)}}{\sqrt{n}} q_{1 - \frac{\alpha}{2}}]$$

Esempio: intervistiamo  $n$  persone che hanno votato  $X_1, \dots, X_n$  ( $= 0, 1$ ) con  $\bar{X}_n = 0,4$

allora  $[0,4 \pm \frac{\sqrt{0,4 \cdot 0,6}}{\sqrt{n}} q_{0,975}]$  è di fiducia per  $p$  e rappresenta la percentuale di Sì su tutta la popolazione a livello 95%