

## Stima puntuale

- **Stimare**: attribuire un valore plausibile a un parametro incognito.

- **Stimatore del parametro**  $\theta$ : ogni statistica

$$T = t(X_1, X_2, \dots, X_n)$$

utilizzata per stimare  $\theta$ . Come ogni altra statistica, lo stimatore è una variabile aleatoria in quanto funzione della variabile aleatoria multipla  $(X_1, X_2, \dots, X_n)$  che rappresenta il campione.

- **Stima**: la singola determinazione dello stimatore,

$$t = t(x_1, x_2, \dots, x_n)$$

cioè il valore che lo stimatore assume in corrispondenza di un campione osservato  $(x_1, x_2, \dots, x_n)$ .

- Per stimare un parametro  $\theta$  possiamo utilizzare *diversi stimatori* (es. per stimare la media della popolazione potremmo usare la mediana, anziché la media, campionaria). E' allora importante trovare dei *criteri per scegliere* lo stimatore più appropriato.

## Criteri di valutazione di uno stimatore

- Uno dei criteri principali è l'**Errore quadratico medio** (MSE) che per lo stimatore

$$T = t(X_1, X_2, \dots, X_n)$$

di  $\theta$  è definito come

$$MSE(T) = E[(T - \theta)^2]$$

- **Esempio**: si assuma che  $X \sim \text{Bin}(1, p)$  e che i campioni hanno dimensione  $n = 5$ . Se  $p = 0,3$  si ha

$\bar{x}$	$f(\bar{x})$	$(\bar{x} - 0,3)^2$	$(\bar{x} - 0,3)^2 f(\bar{x})$	$MSE(\bar{X}) = 0,042$
0,0	0,168	0,09	0,0151	
0,2	0,360	0,01	0,0036	
0,4	0,309	0,01	0,0031	
0,6	0,132	0,09	0,0119	
0,8	0,028	0,25	0,0071	
1,0	0,002	0,49	0,0012	
	1,000		0,0420	

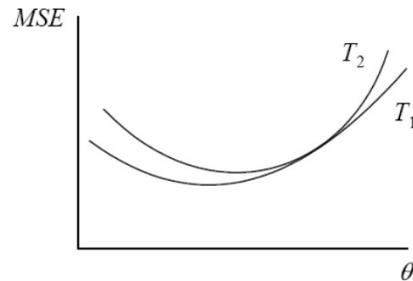
- MSE è una funzione di  $\theta$  che misura l'*errore che si commette* in media utilizzando T per stimare  $\theta$ . Tra due stimatori è preferibile quello che ha *uniformemente* (per ogni  $\theta$ ) l'MSE più piccolo.

$$T = t(X_1, X_2, \dots, X_n)$$

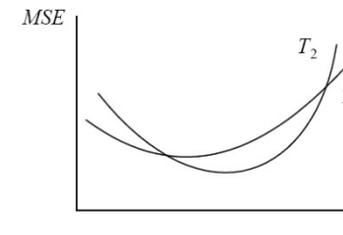
- Dati due stimatori  $T_1$  e  $T_2$  di  $\theta$ ,  $T_1$  è **più efficiente** di  $T_2$  se

$$MSE_{\theta}(T_1) \leq MSE_{\theta}(T_2), \quad \forall \theta$$

e la disuguaglianza vale in senso stretto per almeno un valore di  $\theta$ .



- In generale non esiste uno stimatore con MSE *uniformemente minimo* quindi usualmente ci si restringe alla classe degli stimatori non distorti.



- Lo stimatore  $T = t(X_1, X_2, \dots, X_n)$  di  $\theta$  si dice **non distorto** se

$$E_{\theta}(T) = \theta, \quad \forall \theta$$

- Si chiama **distorsione** (o *bias*) di uno stimatore T di  $\theta$ :

$$B_{\theta}(T) = E_{\theta}(T) - \theta$$

- La distorsione può essere positiva (in media T sovrastima  $\theta$ ) o negativa (in media T sottostima  $\theta$ ).

- L'errore quadratico medio può essere decomposto come segue:

$$MSE(T) = Var(T) + [B(T)]^2$$

**Dimostrazione:**

$$\begin{aligned} MSE(T) &= E[(T - \theta)^2] = E[(T - E(T) + E(T) - \theta)^2] \\ &= E\left[(T - E(T))^2 + 2(T - E(T))(E(T) - \theta) + (E(T) - \theta)^2\right] \\ &= E(T - E(T))^2 + 2E(T - E(T))(E(T) - \theta) + E(E(T) - \theta)^2 \\ &= Var(T) + 2(E(T) - \theta)E(T - E(T)) + [B(T)]^2 = Var(T) + [B(T)]^2 \end{aligned}$$

- Se  $T$  è uno stimatore non distorto, si ha

$$MSE(T) = Var(T)$$

- Dati due stimatori non distorti  $T_1$  e  $T_2$  di  $\theta$ ,  $T_1$  è **più efficiente** di  $T_2$  se

$$Var_{\theta}(T_1) \leq Var_{\theta}(T_2), \quad \forall \theta$$

e la disuguaglianza vale in senso stretto per almeno un valore di  $\theta$ .

### 1. Stime come tiri al bersaglio

Stimare esattamente un parametro è come centrare il bersaglio da parte di un tiratore. Stime ripetute del parametro si possono immaginare come tiri ripetuti. Nella Figura 16.4a vengono rappresentate stime ripetute prodotte da uno stimatore non distorto: le stime sono "disperse" attorno al valore del parametro (centro) senza che si manifestino deviazioni in una particolare direzione; nella Figura 16.4b, invece, le stime tendono a concentrarsi in una particolare zona al di sotto del centro, manifestando, così, la presenza di una distorsione.

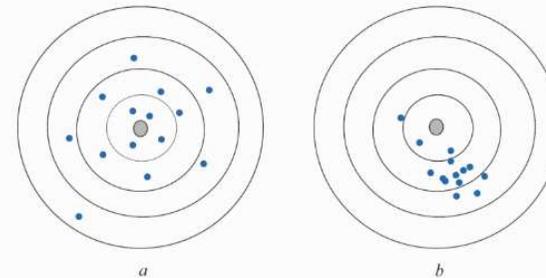


Figura 16.4 Rappresentazione del comportamento di due stimatori, uno non distorto (a) e uno distorto (b).

### Esempio

- A prescindere dalla popolazione di riferimento, la **media campionaria** è uno stimatore non distorto della media della popolazione:

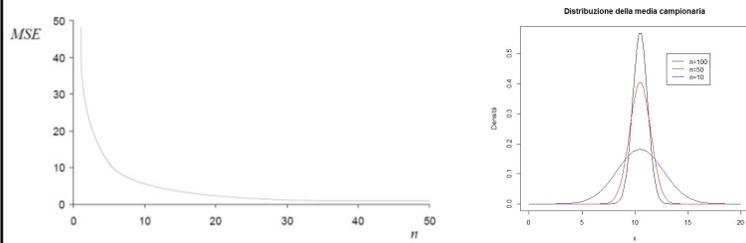
$$E(\bar{X}) = \mu$$

- Di conseguenza l'errore quadratico medio coincide con la varianza:

$$MSE(\bar{X}) = Var(\bar{X}) = \frac{\sigma^2}{n}$$

e quindi diminuisce (la sua precisione aumenta) al crescere della numerosità del campione.

- Ad esempio, se la varianza della popolazione è  $\sigma^2 = 48,5$  si ha



### Esempio

- Si consideri lo **stimatore della varianza**

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

- Questo stimatore è proporzionale alla varianza campionaria

$$\hat{\sigma}^2 = \frac{n-1}{n} S^2$$

di conseguenza  $E(\hat{\sigma}^2) = \frac{n-1}{n} E(S^2) = \frac{n-1}{n} \sigma^2$

e la sua distorsione è pari a

$$B(\hat{\sigma}^2) = E(\hat{\sigma}^2) - \sigma^2 = \frac{n-1}{n} \sigma^2 - \sigma^2 = -\frac{\sigma^2}{n} < 0$$

Lo stimatore **sottostima** la varianza della popolazione. La distorsione tende a 0 al crescere della dimensione campionaria.

- L'errore quadratico medio di questo stimatore non coincide con la sua varianza.

### Esempio

- Si considerino due diversi stimatori della media  $\mu$  di una popolazione, basati su un campione di  $n=4$  estrazioni

$$T_1 = \frac{1}{4} \sum_{i=1}^4 X_i \quad \text{e} \quad T_2 = 0,1 \cdot X_1 + 0,4 \cdot X_2 + 0,4 \cdot X_3 + 0,1 \cdot X_4$$

- I due stimatori sono corretti?

$$E(T_1) = \frac{1}{4} E\left(\sum_{i=1}^4 X_i\right) = \frac{1}{4} \sum_{i=1}^4 E(X_i) = \frac{1}{4} \sum_{i=1}^4 \mu = \frac{1}{4} 4\mu = \mu$$

$$\begin{aligned} E(T_2) &= E(0,1 \cdot X_1 + 0,4 \cdot X_2 + 0,4 \cdot X_3 + 0,1 \cdot X_4) \\ &= 0,1 \cdot E(X_1) + 0,4 \cdot E(X_2) + 0,4 \cdot E(X_3) + 0,1 \cdot E(X_4) \\ &= 0,1 \cdot \mu + 0,4 \cdot \mu + 0,4 \cdot \mu + 0,1 \cdot \mu = \mu \end{aligned}$$

- Quale stimatore è più efficiente?

$$Var(T_1) = \frac{1}{16} V\left(\sum_{i=1}^4 X_i\right) = \frac{1}{16} \sum_{i=1}^4 V(X_i) = \frac{1}{16} \sum_{i=1}^4 \sigma^2 = \frac{1}{16} 4\sigma^2 = \frac{1}{4} \sigma^2$$

$$\begin{aligned} Var(T_2) &= V(0,1 \cdot X_1 + 0,4 \cdot X_2 + 0,4 \cdot X_3 + 0,1 \cdot X_4) \\ &= 0,01 \cdot V(X_1) + 0,16 \cdot V(X_2) + 0,16 \cdot V(X_3) + 0,01 \cdot V(X_4) \\ &= 0,01 \cdot \sigma^2 + 0,16 \cdot \sigma^2 + 0,16 \cdot \sigma^2 + 0,01 \cdot \sigma^2 = 0,34\sigma^2 \end{aligned}$$

### Proprietà asintotiche di uno stimatore

- Le *proprietà asintotiche* riguardano il comportamento di uno stimatore  $T$  quando la dimensione del campione tende a infinito.
- In questo caso si usa la notazione  $T_n$  per indicare che lo stimatore è applicato a un campione di dimensione  $n$ .
- Uno stimatore  $T_n$  di  $\theta$  è **consistente** se per ogni  $\varepsilon > 0$

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

- Uno stimatore  $T_n$  di  $\theta$  è **consistente in media quadratica** se

$$\lim_{n \rightarrow \infty} E[(T_n - \theta)^2] = 0, \quad \forall \theta$$

La consistenza in media quadratica implica quella semplice.

- Uno stimatore  $T_n$  di  $\theta$  è **asintoticamente non distorto** se

$$\lim_{n \rightarrow \infty} B(T_n) = 0, \quad \forall \theta$$

### Metodi per la determinazione di uno stimatore

- Metodo della massima verosimiglianza
- Metodo dei momenti
- Metodo dei minimi quadrati

#### Metodo della massima verosimiglianza

- Dato un campione  $(x_1, x_2, \dots, x_n)$  proveniente da  $X \sim f(x; \theta)$ , si definisce **funzione di verosimiglianza** la probabilità (densità) del campione interpretata come funzione del parametro  $\theta$ ,

$$L(\theta) = f(x_1; \theta) \cdot f(x_2; \theta) \cdot \dots \cdot f(x_n; \theta) = \prod_{i=1}^n f(x_i; \theta)$$

- Quanto più è elevata la funzione di verosimiglianza per un certo valore di  $\theta$ , tanto più questo valore è plausibile. Ciò permette di stabilire qual è il valore più plausibile del parametro.

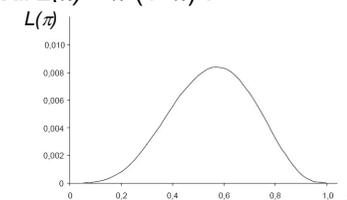
### Esempio

- Si consideri un campione  $(x_1, x_2, \dots, x_n)$  estratto da una *popolazione Bernoulliana* con parametro  $\pi$ , cioè  $X \sim \text{Bin}(1, \pi)$ .
- La *funzione di verosimiglianza* è pari a

$$L(\pi) = \prod_{i=1}^n \pi^{x_i} (1 - \pi)^{1 - x_i} = \pi^{\sum_{i=1}^n x_i} (1 - \pi)^{\sum_{i=1}^n (1 - x_i)} = \pi^s (1 - \pi)^{n-s}$$

dove  $s = \sum_{i=1}^n x_i$  è il numero di successi e  $n-s$  di insuccessi

- Se il campione osservato è  $(1, 0, 1, 0, 0, 1, 1)$  si ha  $n = 7$ ,  $s = 4$  e  $n-s = 3$  da cui  $L(\pi) = \pi^4 (1 - \pi)^3$ .



- La *stima di massima verosimiglianza* è il valore di  $\theta$  che massimizza la funzione di verosimiglianza, cioè

$$\hat{\theta} = t(x_1, x_2, \dots, x_n) \text{ tale che } L(\hat{\theta}) = \max_{\theta \in \Theta} L(\theta)$$

- Lo *stimatore di massima verosimiglianza*, come ogni altro stimatore è una variabile casuale:

$$t(X_1, X_2, \dots, X_n)$$

- Generalmente risulta conveniente massimizzare la *log-verosimiglianza*

$$l(\theta) = \log[L(\theta)]$$

risolvendo l'equazione

$$\frac{d l(\theta)}{d \theta} = 0$$

e controllando che in corrispondenza del valore trovato di  $\theta$  la derivata seconda di  $l(\theta)$  sia negativa.

### Esempio

- Nel caso della *popolazione Bernoulliana*  $X \sim \text{Bin}(1, \pi)$ , la *funzione di log-verosimiglianza* è pari a

$$l(\pi) = \log[L(\pi)] = s \log(\pi) + (n - s) \log(1 - \pi)$$

- La *derivata prima* di  $l(\pi)$  rispetto a  $\pi$  è

$$\frac{d l(\pi)}{d \pi} = \frac{s}{\pi} - \frac{n - s}{1 - \pi}$$

che posta pari a 0 dà:

$$\frac{s}{\pi} - \frac{n - s}{1 - \pi} = 0 \rightarrow s(1 - \pi) - (n - s)\pi = 0 \rightarrow s - s\pi - n\pi + s\pi = 0$$

$$s - n\pi = 0 \rightarrow s = n\pi \text{ da cui}$$

$$\hat{\pi} = \frac{s}{n} = \frac{\sum_{i=1}^n X_i}{n}$$

è la *stima di massima verosimiglianza* di  $\pi$ . Si noti che questa coincide con la media campionaria. Quindi lo *stimatore di massima verosimiglianza* di  $\pi$  è:

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

- Per il campione osservato (1,0,1,0,0,1,1), si ha  $\hat{\pi} = 4/7 = 0,571$

### Esempio

- Si assuma che il campione è estratto da una *popolazione Normale* con media e varianza incognite,  $X \sim N(\mu, \sigma^2)$ .
- Lo *stimatore di massima verosimiglianza* della media,  $\mu$ , coincide con la media campionaria

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

- Lo *stimatore di massima verosimiglianza* della varianza,  $\sigma^2$ , è dato da

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n}$$

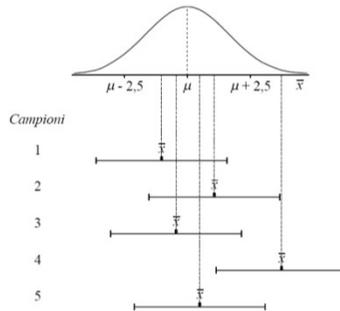
### Proprietà di uno stimatore di massima verosimiglianza

- Se  $\hat{\theta}$  è lo stimatore di massima verosimiglianza di  $\theta$ , allora  $g(\hat{\theta})$  è lo stimatore di massima verosimiglianza di  $g(\theta)$ . Ad esempio lo stimatore di  $\sigma$  nel caso di popolazione Normale è  $\hat{\sigma} = \sqrt{\hat{\sigma}^2}$ .
- In genere permette di ottenere *stimatori efficienti* nella classe di quelli non distorti.
- Sotto certe condizioni di regolarità, lo stimatore di massima verosimiglianza gode delle seguenti *proprietà asintotiche*:
  - è *consistente*;
  - è *asintoticamente efficiente*;
  - è *asintoticamente non distorto*;
  - ha *distribuzione asintotica normale*.

## Stima per intervallo

- In ogni stima è ovviamente insito un certo *marginale di errore*, ad esempio la probabilità  $P(\bar{X} = \mu)$  è nulla quando  $X$  è continua.
- Nell'ambito della *stima per intervallo* un parametro viene stimato tramite un *intervallo di confidenza* anziché con un singolo valore come nell'ambito della *stima puntuale*.

• Esempio: Se  $X \sim N(\mu, 44)$ ,  $n = 27$  e si usano intervalli del tipo  $(\bar{X} - 2,5; \bar{X} + 2,5)$  anziché  $\bar{X}$  per stimare  $\mu$ , si ottiene



## Intervallo di confidenza

- Dato un campione  $(X_1, X_2, \dots, X_n)$  estratto da una popolazione di cui interessa stimare il parametro  $\theta$ , e date due statistiche campionarie  $L_1 = L_1(X_1, X_2, \dots, X_n)$  e  $L_2 = L_2(X_1, X_2, \dots, X_n)$  con  $L_1 \leq L_2$  per ogni possibile campione, l'intervallo  $(L_1, L_2)$  è un *intervallo di confidenza* per  $\theta$  al  $(1-\alpha)100\%$  se

$$P(L_1 \leq \theta \leq L_2) = 1 - \alpha$$

- Il *coefficiente fiduciario*  $1-\alpha$  è solitamente pari al 90%, 95% o 99%.
- L'*ampiezza di un intervallo di confidenza* è data da

$$A = L_2 - L_1$$

## Intervalli di confidenza per la media (popolazione normale, varianza nota)

- Se la popolazione ha *distribuzione normale*,  $X \sim N(\mu, \sigma^2)$ , si ha

$$\frac{\bar{X} - \mu}{\sqrt{\sigma^2/n}} \sim N(0, 1)$$

- Di conseguenza

$$P\left(-z_{\alpha/2} \leq \frac{\bar{X} - \mu}{\sqrt{\sigma^2/n}} \leq z_{\alpha/2}\right) = 1 - \alpha$$

dove  $z_{\alpha/2}$  è il quantile tale che  $P(Z \geq z_{\alpha/2}) = \alpha/2$

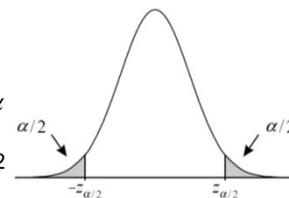
- Esplicitando rispetto a  $\mu$

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

- L'*intervallo di confidenza* per  $\mu$  al  $(1-\alpha)100\%$  è dato da

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

- L'*ampiezza* dell'intervallo di confidenza è  $A = 2z_{\alpha/2}\sqrt{\sigma^2/n}$

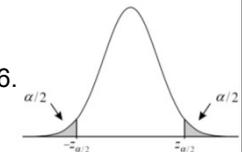


## Esempio

- Si supponga di voler costruire l'*intervallo di confidenza* al 95% per  $\mu$  sulla base del campione  
170,75 186,14 173,39 185,12 173,39  
sotto l'assunzione  $X \sim N(\mu, \sigma^2)$ , con varianza  $\sigma^2 = 50$ .

- La *media campionaria* è pari a  $\bar{x} = 177,758$

- Dato che  $1-\alpha = 0,95$ , si ha che  $\alpha = 0,05$  e il corrispondente *quantile* della distribuzione Normale standard è pari a  $z_{\alpha/2} = z_{0,025} = 1,96$ . Si noti che  $\Phi(z_{0,025}) = 0,975 = 1-\alpha/2$ .



- L'*intervallo di confidenza* è quindi

$$(177,758 - 1,96\sqrt{50/5}, 177,758 + 1,96\sqrt{50/5}) = (171,560, 183,956)$$

che ha ampiezza pari a

$$A = 2 \cdot 1,96\sqrt{50/5} = 12,396$$

**Esempio (determinazione della numerosità campionaria)**

- Nell'esempio precedente abbiamo ottenuto l'*intervallo di confidenza* al 95%

$$(177,758 - 1,96\sqrt{50/5}, 177,758 + 1,96\sqrt{50/5}) = (171,560, 183,956)$$

che ha ampiezza pari a

$$A = 2 \cdot 1,96\sqrt{50/5} = 12,396$$

- Immaginiamo di voler ridurre l'ampiezza dell'intervallo. Che strumenti abbiamo a disposizione?
  - Possiamo ridurre il livello di confidenza
  - Possiamo aumentare la numerosità campionaria.
- Quanto deve essere grande il campione se vogliamo ottenere un intervallo di confidenza di ampiezza pari a 5 cm?

$$5 = 2 \cdot 1,96\sqrt{50/n} \rightarrow 5^2 = 2^2 \cdot 1,96^2 \cdot 50/n \rightarrow n = \frac{2^2 \cdot 1,96^2 \cdot 50}{5^2} \approx 31$$

- In generale

$$A = 2z_{\alpha/2}\sqrt{\sigma^2/n} \rightarrow A^2 = 2^2 \cdot z_{\alpha/2}^2 \cdot \sigma^2/n \rightarrow n = \frac{2^2 \cdot z_{\alpha/2}^2 \cdot \sigma^2}{A^2}$$

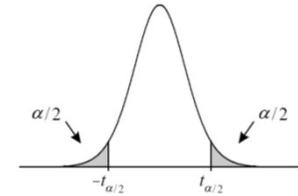
**Intervallo di confidenza per la media  
(popolazione normale, varianza non nota)**

- Se  $X \sim N(\mu, \sigma^2)$  ma  $\sigma^2$  non è nota, si utilizza il risultato

$$\frac{\bar{X} - \mu}{\sqrt{S^2/n}} \sim t(n-1)$$

- Di conseguenza

$$P\left(-t_{\alpha/2} \leq \frac{\bar{X} - \mu}{\sqrt{S^2/n}} \leq t_{\alpha/2}\right) = 1 - \alpha$$



dove  $t_{\alpha/2}$  è il quantile tale che  $P(T \geq t_{\alpha/2}) = \alpha/2$  e  $T \sim t(n-1)$

- L'*intervallo di confidenza* per  $\mu$  al  $(1-\alpha)100\%$  è dato da

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

la cui ampiezza è  $A = 2t_{\alpha/2}\sqrt{S^2/n}$

**Esempio**

- Si supponga di voler costruire l'*intervallo di confidenza* al 95% per  $\mu$  sulla base del campione  
170,75 186,14 173,39 185,12 173,39  
sotto l'assunzione  $X \sim N(\mu, \sigma^2)$ , con varianza  $\sigma^2$  non nota.

- La *media campionaria* e la *varianza campionaria* sono pari a

$$\bar{x} = 177,758 \quad e \quad s^2 = 52,932$$

- Dato che  $1-\alpha = 0,95$ , si ha che  $\alpha = 0,05$  e il corrispondente *quantile* della distribuzione  $t(4)$  è pari a  $t_{\alpha/2} = t_{0,025} = 2,776$ .

- L'*intervallo di confidenza* è quindi

$$(177,758 - 2,776\sqrt{52,932/5}, 177,758 + 2,776\sqrt{52,932/5}) = (168,726, 186,790)$$

che ha ampiezza pari a  $A = 2 \cdot 2,776\sqrt{52,932/5} = 18,064$  maggiore di quella calcolata per il caso della varianza nota.

**Intervallo di confidenza per la media**

(popolazione qualsiasi, varianza non nota, grandi campioni)

- Per una *popolazione qualsiasi*, nel caso di *grandi campioni* ( $n \geq 30$ ) si ha che

$$\frac{\bar{X} - \mu}{\sqrt{S^2/n}}$$

ha approssimativamente distribuzione Normale standard.

- Di conseguenza

$$P\left(-z_{\alpha/2} \leq \frac{\bar{X} - \mu}{\sqrt{S^2/n}} \leq z_{\alpha/2}\right) = 1 - \alpha$$

dove  $z_{\alpha/2}$  è il quantile tale che  $P(Z \geq z_{\alpha/2}) = \alpha/2$  e  $Z \sim N(0, 1)$

- L'*intervallo di confidenza* per  $\mu$  al  $(1-\alpha)100\%$  è dato da

$$\left(\bar{X} - z_{\alpha/2}\sqrt{S^2/n}, \bar{X} + z_{\alpha/2}\sqrt{S^2/n}\right)$$

la cui ampiezza è  $A = 2z_{\alpha/2}\sqrt{S^2/n}$

**Esempio**

- Si supponga di voler costruire l'*intervallo di confidenza* al 90% per  $\mu$  sulla base di un campione di dimensione  $n = 100$  che ha media e varianza campionarie pari rispettivamente a

$$\bar{x} = 15,52 \quad e \quad s^2 = 24,47$$

- Dato che  $1-\alpha = 0,90$ , si ha che  $\alpha = 0,10$  e il corrispondente *quantile* della distribuzione Normale standard è pari a  $z_{\alpha/2} = z_{0,05} = 1,645$ . Si noti che  $\Phi(z_{0,05}) = 0,95 = 1-\alpha/2$ .

- L'*intervallo di confidenza* è quindi

$$(15,52 - 1,645\sqrt{24,47/100}, 15,52 + 1,645\sqrt{24,47/100}) = (14,706, 16,334)$$

che ha ampiezza pari a

$$A = 2 \cdot 1,645\sqrt{24,47/100} = 1,627$$

**Intervalli di confidenza per la media**

(popolazione Bernoulliana, grandi campioni)

- Per una *popolazione Bernoulliana*, nel caso di *grandi campioni* ( $n \geq 30$ ) si ha che

$$\frac{\bar{X} - \pi}{\sqrt{\bar{X}(1-\bar{X})/n}}$$

ha approssimativamente distribuzione Normale standard.

- Di conseguenza

$$P\left(-z_{\alpha/2} \leq \frac{\bar{X} - \pi}{\sqrt{\bar{X}(1-\bar{X})/n}} \leq z_{\alpha/2}\right) = 1 - \alpha$$

dove  $z_{\alpha/2}$  è il quantile tale che  $P(Z \geq z_{\alpha/2}) = \alpha/2$  e  $Z \sim N(0,1)$

- L'*intervallo di confidenza* per  $\pi$  al  $(1-\alpha)100\%$  è dato da

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

la cui ampiezza è

$$A = 2z_{\alpha/2}\sqrt{\bar{X}(1-\bar{X})/n}$$

**Esempio**

- Si supponga di voler costruire l'*intervallo di confidenza* al 99% per  $\pi$  sulla base di un campione di dimensione  $n = 50$  in cui ci sono  $\sum_{i=1}^n x_i = 32$  successi, da cui

$$\bar{x} = 32/50 = 0,64 \quad e \quad \bar{x}(1-\bar{x}) = 0,64 \cdot (1-0,64) = 0,23$$

- Dato che  $1-\alpha = 0,99$ , si ha che  $\alpha = 0,01$  e il corrispondente *quantile* della distribuzione Normale standard è pari a  $z_{\alpha/2} = z_{0,005} = 2,576$ . Si noti che  $\Phi(z_{0,005}) = 0,995 = 1-\alpha/2$ .

- L'*intervallo di confidenza* è quindi

$$(0,64 - 2,576\sqrt{0,23/50}, 0,64 + 2,576\sqrt{0,23/50}) = (0,465, 0,815)$$

che ha ampiezza pari a

$$A = 2 \cdot 2,576\sqrt{0,23/50} = 0,350$$

**Intervalli di confidenza per la varianza**

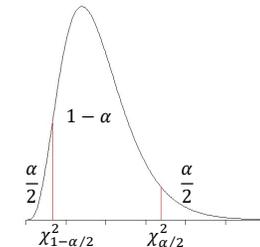
(popolazione normale, media e varianza non note)

- Se  $X \sim N(\mu, \sigma^2)$  con  $\mu$  non noto, si ha

$$\frac{(n-1)S^2}{\sigma^2} \approx \chi^2(n-1)$$

- Di conseguenza

$$P\left(\chi^2_{1-\alpha/2} \leq \frac{(n-1)S^2}{\sigma^2} \leq \chi^2_{\alpha/2}\right) = 1 - \alpha$$



dove  $\chi^2_{1-\alpha/2}$  e  $\chi^2_{\alpha/2}$  sono i quantili della distribuzione  $\chi^2(n-1)$  tali che

$$P(\chi^2 \geq \chi^2_{1-\alpha/2}) = 1 - \alpha/2 \quad e \quad P(\chi^2 \geq \chi^2_{\alpha/2}) = \alpha/2$$

- L'*intervallo di confidenza* per  $\sigma^2$  al  $(1-\alpha)100\%$  è dato da

$$\left((n-1)S^2 / \chi^2_{\alpha/2}, (n-1)S^2 / \chi^2_{1-\alpha/2}\right)$$

Passaggi per arrivare all'intervallo di confidenza:

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

### Esempio

- Si supponga di voler costruire l'*intervallo di confidenza* al 90% per  $\sigma^2$  sulla base del campione (già visto in precedenza)

170,75   186,14   173,39   185,12   173,39

sotto l'assunzione  $X \sim N(\mu, \sigma^2)$  con  $\mu$  non noto.

- La *varianza campionaria* è pari a

$$s^2 = 52,932;$$

dato che  $1-\alpha = 0,90$ , si ha che  $\alpha = 0,10$  e i corrispondenti *quantili* della distribuzione  $\chi^2(4)$  sono pari a

$$\chi_{1-\alpha/2}^2 = \chi_{0,95}^2 = 0,71 \quad e \quad \chi_{\alpha/2}^2 = \chi_{0,05}^2 = 9,49$$

- L'*intervallo di confidenza* è quindi

$$(4/9,49 \cdot 52,932, 4/0,71 \cdot 52,932) = (22,311, 298,208)$$

che ha ampiezza pari a  $A = 298,208 - 22,311 = 275,897$