

Università degli Studi di Padova  
Facoltà di scienze statistiche  
Corso di laurea in statistica popolazione e società

# INTERVALLI DI CONFIDENZA APPROSSIMATI PER PROPORZIONI

Laureando:  
Basso Sebastiano

Relatore:  
Prof.ssa Alessandra Salvan

Anno Accademico 2008/2009

Ringrazio: Umberto, Maurizia, Chiara, Elena,  
Silvia P., Silvia M., Taddeo, Lucia, Silvio, Arianna, Irene, Federico,  
Martina.

# Indice

1 INFERENZA SULLA PROBABILITA' DI SUCCESSO IN UN MODELLO BINOMIALE 5	
1.1 Stima puntuale di $\theta$ nel modello binomiale . . . . .	5
1.2 Inferenza di verosimiglianza per il modello binomiale . . . . .	6
1.2.1 Test di verosimiglianza e intervalli di copertura . . . . .	7
1.3 Intervalli di copertura per la proporzione Binomiale . . . . .	11
1.4 Simulazione . . . . .	13
1.4.1 Simulazione con R . . . . .	14
1.5 Probabilità di copertura . . . . .	16
1.5.1 Cenni teorici . . . . .	16
1.5.2 Espressione per la probabilità di copertura degli intervalli di copertura per la proporzione binomiale . . . . .	16
1.6 Apporto di questa tesi . . . . .	17
2 RISULTATI 19	
2.1 Oscillazioni . . . . .	19
2.2 Dipendenza dal valore del parametro . . . . .	19
2.3 Dipendenza dalla numerosità campionaria . . . . .	24
3 CONCLUSIONI 33	
3	
4	

## Capitolo 1

# INFERENZA SULLA PROBABILITA' DI SUCCESSO IN UN MODELLO BINOMIALE

## 1.1 Stima puntuale di $\theta$ nel modello binomiale

In un esperimento casuale si osserva il numero complessivo di successi  $y$  in  $n$  prove indipendenti con costante probabilità di successo,  $\theta \in (0; 1)$  ignoto. Allora  $y$  è una realizzazione di  $Y \sim \text{Bi}(n; \theta)$ . Una stima naturale del parametro  $\theta = \hat{\theta}_n$ , ossia della probabilità di successo in una singola prova elementare è la frequenza relativa di successi ottenuti nelle  $n$  prove,

$$\hat{\theta}_n = \frac{y}{n}$$

che coincide con la stima di massima verosimiglianza

Lo stimatore  $\hat{\theta}_n$  è una variabile casuale discreta con supporto

$$S_{\hat{\theta}_n} = \{0, \frac{1}{n}, \dots, 1\}$$

con

$$P_{\hat{\theta}_n}(\hat{\theta}_n = \frac{y}{n}) = \binom{n}{y} \theta^y (1-\theta)^{n-y}$$

è funzione di probabilità sotto  $\theta$

$$P_{\hat{\theta}_n}(\hat{\theta}_n = \frac{y}{n}) = \binom{n}{y} \theta^y (1-\theta)^{n-y}$$

5

$$P_{\hat{\theta}_n}(\hat{\theta}_n = \frac{y}{n}) = \binom{n}{y} \theta^y (1-\theta)^{n-y}$$

5

Se  $n$  è sufficientemente grande, lo stimatore  $\hat{\theta}_n$  produce realizzazioni che con elevata probabilità non si discostano eccessivamente dal vero e ignoto valore del parametro.

Infatti, quale sia  $\theta \in (0; 1)$ , la variabile casuale  $\hat{\theta}_n$  ha valore atteso

$$E(\hat{\theta}_n) = E\left(\frac{Y}{n}\right) = \frac{E(Y)}{n}$$

5

1

n

$$E(Y) = n\theta$$

$n$   
 $= \frac{1}{n}$   
 e varianza  
 $V(\hat{\theta}_n) = \frac{1}{n} V(Y)$   
 $n \hat{\theta}_n = \sum_{i=1}^n Y_i$   
 $n^2 V(\hat{\theta}_n) = n V(Y)$   
 $n^2 = \frac{1}{n} V(Y)$   
 $n$   
 :

Per  $n$  estremamente grande, la distribuzione dello stimatore è quasi degenere attorno al vero valore del parametro.

Per  $n$  moderatamente grande, come è consueto nelle applicazioni, per un orientamento di massima, se sia  $n \hat{\theta}_n$  sia  $n(1 - \hat{\theta}_n)$  sono maggiori di 5, il teorema del limite centrale dà, sotto  $\hat{\theta}_n$ ,  $Y \sim N(\hat{\theta}_n, \frac{1}{n})$

$N(0; 1)$ :  
 Ciò fornisce per la distribuzione dello stimatore l' approssimazione  
 $\hat{\theta}_n \sim N(\theta, \frac{1}{n})$   
 $n \hat{\theta}_n$

## 1.2 Inferenza di verosimiglianza per il modello binomiale

Data  $y$ , un'osservazione di  $Y \sim \text{Bi}(n; \theta)$ , con  $0 < \theta < 1$  la funzione di verosimiglianza è:  
 $L(\theta) = \binom{n}{y} \theta^y (1 - \theta)^{n-y}$

$n$   
 $y! \binom{n}{y} (1 - \theta)^{n-y}$ :  
 In  $L(\theta)$  i fattori che non dipendono da  $\theta$  si possono trascurare e questo è il caso del coefficiente binomiale. Quindi la funzione di log-verosimiglianza risulta essere  
 $l(\theta) = y \log \theta + (n - y) \log(1 - \theta)$ : (1.1)

La funzione punteggio risulta

$$l'(\theta) = y \frac{1}{\theta} - \frac{n - y}{1 - \theta}$$

l'informazione osservata risulta

$$j(\theta) = y \frac{1}{\theta^2} + \frac{n - y}{(1 - \theta)^2}$$

La stima di massima verosimiglianza è la radice dell' equazione di verosimiglianza

$$y \frac{1}{\theta} - \frac{n - y}{1 - \theta} = 0$$

che risulta, per  $0 < y < n$ ,

$$\hat{\theta} = \frac{y}{n}$$

e la (1.2) calcolata in  $\hat{\theta}$  da

$$j(\hat{\theta}) = \frac{y}{\hat{\theta}^2} + \frac{n - y}{(1 - \hat{\theta})^2}$$

Per  $y = 0$  e  $y = n$  la stima di massima verosimiglianza non esiste e si procede allargando lo spazio parametrico da  $(0; 1)$  a  $[0; 1]$  con la sicurezza che le proprietà asintotiche non vengono alterate in quanto:

$\Pr(0 < Y < n) \rightarrow 1$  per  $n \rightarrow \infty$  per ogni  $\theta \in (0; 1)$ .

### 1.2.1 Test di verosimiglianza e intervalli di confidenza

La funzione di log-verosimiglianza (1.1) calcolata nella stima di massima verosimiglianza (1.2) risulta essere:

$$l(\hat{\theta}) = y \log y + (n - y) \log(1 - y) \quad (1.5)$$

Con questa e con la (1.1) si possono calcolare i test di verosimiglianza, unilaterali e bilaterali, di seguito riportati per l'ipotesi semplice  $H_0: \theta = \theta_0$ :

Il log-rapporto di verosimiglianza è

$$W(\theta_0) = 2 [l(\hat{\theta}) - l(\theta_0)]$$

La radice con segno di  $W(\theta_0)$  è

$$r(\theta_0) = \text{sgn}(y - \theta_0) \sqrt{2 [l(\hat{\theta}) - l(\theta_0)]}$$

con  $W(\theta_0) = r(\theta_0)^2$ .

Sotto  $H_0$ , per  $n$  sufficientemente grande,

$$W(\theta_0) \xrightarrow{d} \chi^2_1$$

e  $r(\theta_0) \xrightarrow{d} N(0, 1)$ :

Un intervallo di confidenza basato su  $W(\theta)$  è

$$f_{1-\alpha/2}(0; 1) \leq W(\theta) \leq f_{\alpha/2}(0; 1)$$

o, equivalentemente,

$$f_{1-\alpha/2}(0; 1) \leq l(\theta) - l(\hat{\theta}) \leq f_{\alpha/2}(0; 1)$$

Dal grafico di  $l(\theta)$  si può ricavare un intervallo di confidenza tracciando la linea orizzontale di

$$l(\hat{\theta}) - f_{\alpha/2}(0; 1)$$

ordinata  $l(\hat{\theta}) - f_{\alpha/2}(0; 1)$ . Tutti i valori di  $\theta$  per cui  $l(\theta)$  è superiore alla linea orizzontale formano

l'intervallo di confidenza a livello nominale 0.95 basato su  $W(\theta)$ . Da Pace, Salvani (2001).

Una miglioramento della statistica  $r(\theta)$

Questa nuova statistica si propone di dare un miglioramento della statistica  $r(\theta)$  utilizzando una quantità  $u$  definita come nel seguito. L'espressione generale è:

$$r_u(\theta_0) = r(\theta_0) +$$

$1$

$$r(\theta_0)$$

$$\log u(\theta_0)$$

$$r(\theta_0) u(\theta_0)$$

Per la definizione di  $r_u$  si fa riferimento a Severini (2000, capitolo 10).

Da  $l(\theta) = l(\theta; \hat{\theta})$  ricaviamo la derivata rispetto a  $\hat{\theta}$

$$l'_{\hat{\theta}}(\theta) = n \log \theta - n \log(1 - \theta)$$

da cui si ricava

$$l'_{\hat{\theta}}(\hat{\theta}) = n \log \hat{\theta} - n \log(1 - \hat{\theta})$$

da cui

$$\hat{\theta} - n \log \hat{\theta} + n \log(1 - \hat{\theta})$$

da cui

$$1 - n \log \hat{\theta} + n \log(1 - \hat{\theta})$$

da cui

$$= n \log \hat{\theta} - n \log(1 - \hat{\theta})$$

da cui

$$\hat{\theta} - n \log \hat{\theta} + n \log(1 - \hat{\theta})$$

da cui

$$\hat{\theta} - n \log \hat{\theta} + n \log(1 - \hat{\theta})$$

Calcolando la (1.2) in  $\hat{\theta}$  si trova

$$u = j(\hat{\theta})^{-1} \frac{d}{d\theta} \ln l(\hat{\theta}) = n \hat{\theta}^{-1} \ln \hat{\theta}$$

$$u(\hat{\theta}) = n \hat{\theta}^{-1} \ln \hat{\theta}$$

$$u(\hat{\theta}) = n \hat{\theta}^{-1} \ln \hat{\theta}$$

o in alternativa si può utilizzare il dato da

$$u_L =$$

$$1 \cdot \exp \left[ \ln \left( \frac{1}{\hat{\theta}} \right) \right] = \frac{1}{\hat{\theta}}$$

$$j(\hat{\theta}) = \frac{1}{\hat{\theta}^2}$$

dove

$$m(\hat{\theta}) = (g_0(t))^{-1}$$

con

$$g(t) = nt$$

per il caso binomiale e di conseguenza

$$g_0(t) = n$$

ne consegue

$$m(\hat{\theta}) =$$

$$1$$

$$n$$

$$9$$

ora abbiamo

$$u_L = \frac{1}{\hat{\theta}}$$

$$u(\hat{\theta}) = n \hat{\theta}^{-1} \ln \hat{\theta}$$

$$j(\hat{\theta}) = \frac{1}{\hat{\theta}^2}$$

Con le precedenti formule possiamo ora calcolare la statistica  $r$

$$r = r +$$

$$1$$

$$r$$

$$\log \frac{u}{u_L}$$

$$r =$$

$$0$$

$$r_0 = r +$$

$$1$$

$$r$$

$$\log \frac{u_L}{u}$$

$$r =$$

Questa statistica non può essere utilizzata per  $\hat{\theta} = 1$  poichè in questo caso  $u = 0$  mentre  $r$  è finito. Per l'argomentazione di ciò si veda Severini (2000, capitolo 10).

10

### 1.3 Intervalli di confidenza per la proporzione Binomiale

In questo paragrafo si riportano schematicamente le espressioni di diversi intervalli di confidenza per  $\theta$ . Con  $z_\alpha$  si indica il quantile  $\alpha$  di una normale standard.

1. Wald Standard

$$\hat{\theta} \pm z_{\alpha/2} \sqrt{\frac{\hat{\theta}(1-\hat{\theta})}{n}}$$

$$\hat{\theta} \pm z_{\alpha/2} \sqrt{\frac{\hat{\theta}(1-\hat{\theta})}{n}}$$

$$n$$

Questo è il classico intervallo di confidenza approssimato che viene presentato in tutti i corsi di statistica. Esso è simmetrico rispetto al valore stimato di  $\theta$  e la sua ampiezza dipende dall'errore standard asintotico di  $\hat{\theta}$ . Questo metodo fallisce nel produrre un intervallo nel caso in cui il numero di successi sia uguale a 0 o a 1. Questo tipo di intervallo è preso in considerazione ad esempio in Agresti e Coull (1998), Vollset (1993), Brown, Cai, DasGupta (2001), Borkowf (2006), Wang (2007)

2. Wald con correzione di continuità

$$\hat{\theta} \pm z_{\alpha/2} \sqrt{\frac{\hat{\theta}(1-\hat{\theta})}{n}}$$

$$\hat{\theta} \pm z_{\alpha/2} \sqrt{\frac{\hat{\theta}(1-\hat{\theta})}{n}}$$

$$n$$

$$+$$

$$1$$

2n#

E' simile simile all'intervallo alla Wald standard a meno di una quantità 1  
2n che è la correzione  
di continuità poiché la distribuzione binomiale è discreta. Vollset (1993).

### 3. Wald Blyth Still

^\_ \_2664

Z\_

2 r n □ □ Z\_

2

2 □ □

2z\_

2 p n □ □ 1

n p ^ \_ (1 □ □ ^ \_ ) +

1

2n

3775

Questo intervallo è stato proposto da Blyth e Still (1983), esso è una modi\_cazione della  
correzione di continuità per l'intervallo Wald.

### 4. Wald aggiustato (Agresti-Coull)

11

Aggiunge due successi e due insuccessi per la stima di \_, ottenendo ~\_ = y+2

n+4

~\_ \_ Z\_

2 r ~\_ (1 □ □ ~\_ )

n + 4

Questo intervallo di con\_denza è stato proposto da Agresti e Coull (1998) ed è ottenuto  
sommano due successi e due insuccessi al campione in esame. Citato anche in , Blaker  
(2000), Brown, Cai, DasGupta (2001), Reiczigel (2003), Borkowf (2006), Wang (2007).

### 5. Wald logistico

1 □ □ "1 + exp(ln\_ y

n □ □ y\_

Z\_

2 p n ^ \_ (1 □ □ ^ \_ ) # □ □ 1

### 6. Intervallo Score

^\_ +

Z\_2

\_2

2n \_ Z\_

2 r ^ \_ (1 □ □ ^ \_ ) +

Z\_2

\_2

4n

n

1 +

Z\_2

\_2

n

Questo intervallo è basato anch'esso sull'approssimazione normale ma si di\_erenza dall'intervallo  
Wald standard poiché usa l'errore standard sotto l'ipotesi nulla al posto di quello  
stimato. Questo intervallo fu discusso in principio da Wilson (1927). Esso è citato in  
Vollset (1993), Agresti e Coull (1998), Blaker (2000), Reiczigel (2003).

### 7. Score con correzione di continuità

(y + 1=2) +

Z\_2

\_2

2\_ Z\_

2 r (y \_ 1=2) □ □ (y+1=2)2

n +

Z\_2

\_2

4

n + Z\_2

\_2

E' simile all' intervallo score a meno della costante 1=2 che è la correzione di continuità

poichè la binomiale è una distribuzione discreta. Vollset (1993).

12

Intervalli esatti

Sono basati sull'inversione del test esatto:

$$P = 2[fPr_{f_0}(Y = y) + \min Pr_{f_0}(Y \leq y); Pr_{f_0}(Y \geq y)]$$

dove  $f_0$  è il parametro sotto l'ipotesi nulla e  $0 \leq f \leq 1$ . Con  $f = 1$  si ottiene il p-value massimo.

Con  $f = 1/2$  si ottiene il MID-p-value. Invertendo il test si ottengono due intervalli di confidenza.

1. Max-P

$$MAX_l = 8<$$

$f_0$ :

$$y \sum_{t=0}^n X_t$$

0@

n

$$t1A_{-t}(1 - f_0)^n \sum_{t=0}^y = 1 - f_0$$

$f_0 =$

$$MAX_u = 8<$$

$f_0$ :

$$yX$$

t=0

0@

n

$$t1A_{-t}(1 - f_0)^n \sum_{t=0}^y = f_0$$

$f_0 =$

2. Mid-P

E' formato dalle stesse equazioni del Max-P con la differenza che viene aggiunta metà della probabilità assegnata al risultato osservato:

$$MID_l = 8<$$

$f_0$ :

1

20@

n

$$y1A_{-y}(1 - f_0)^n \sum_{t=0}^y +$$

$$y \sum_{t=0}^n X_t$$

0@n

$$t1A_{-t}(1 - f_0)^n \sum_{t=0}^y = 1 - f_0$$

$f_0 =$

$$MID_u = 8<$$

$f_0$ :

1

20@

n

$$y1A_{-y}(1 - f_0)^n \sum_{t=0}^y +$$

$$yX$$

t=0

0@

n

$$t1A_{-t}(1 - f_0)^n \sum_{t=0}^y = f_0$$

$f_0 =$

## 1.4 Simulazione

Per valutare l'adattabilità di un metodo per determinare un intervallo di confidenza, in molti articoli viene considerata la probabilità di copertura empirica confrontata con quella nominale.

Con R vengono generati una serie di campioni casuali dalla distribuzione binomiale, per ciascun campione viene costruito un intervallo di confidenza a livello predefinito (nel nostro caso al 95%) e

13



viene calcolata la percentuale di volte in cui il "vero" valore del parametro  $\theta$  (cioè quello utilizzato per generare i campioni) cade dentro l'intervallo generato. Per un'intervallo al 95% mi aspetto, se il metodo fosse esatto, che nel 95% degli esperimenti il valore del parametro utilizzato per generare i campioni sia compreso negli intervalli generati di volta in volta con lo stesso metodo. Nella pratica, non sempre il livello di copertura è uguale a quello nominale. Adirittura, alcuni intervalli tendono ad avere un livello effettivo sempre maggiore di quello nominale e per questo sono detti "conservativi", mentre altri hanno oscillazioni anche considerevoli al di sotto e al di sopra del livello nominale. E' anche vero che non sempre sia migliore un intervallo "conservativo" che uno meno conservativo. Le procedure che considereremo migliori saranno allora quelle che avranno livello di copertura più vicino al valore nominale.

### 1.4.1 Simulazione con R

Per effettuare le simulazioni si è usato l'ambiente R. Per ogni intervallo di copertura analizzato si sono utilizzate 10000 simulazioni. Di seguito è riportato un esempio di un semplice programma in R per la simulazione dell'intervallo Wald aggiustato (Agresti e Coull (1998)):

```
simu3<-function(n,p,conf.l)
{cont<-0;for (i in 1:nsim)
{sample<-rbinom(n,1,p);
wald.adj<-((sum(sample)+2)/(n+4))+c(-1,1)*qnorm(1-(1-conf.level)/2)
*sqrt(((sum(sample)+2)/(n+4))*(1-((sum(sample)+2)/(n+4))))/(n+4));
if (wald.adj[1]<p & p<wald.adj[2])cont=cont+1};
print(cont/nsim)}
```

la funzione così generata viene poi vettorizzata rispetto  $\theta$  o rispetto  $n$ , nell'esempio rispetto  $\theta$ :

```
simuvect3<-Vectorize(simu3,vectorize.args="p")
```

da questa viene fatto il grafico con variabile indipendente  $\theta$  (o  $n$ ) e nelle ordinate la probabilità di copertura:

14

```
plot(function(x) simuvect3(5,x,0.95),from=0,to=1,xlab="p",ylab="coverage", main="c.i.
waldadj 95 n=5")
```

ed il risultato ottenuto è:

```
0.0 0.2 0.4 0.6 0.8 1.0
0.90 0.92 0.94 0.96 0.98 1.00
c.i. waldadj 95% n=5
```

```
p
coverage
```

Figura 1.1: Probabilità di copertura al variare di  $\theta$  per l'intervallo Wald aggiustato

15

## 1.5 Probabilità di copertura

### 1.5.1 Cenni teorici

Si dice probabilità di copertura della regione di copertura  $\hat{C}_n(Y)$  la funzione di

$\Pr(\hat{C}_n(Y))$ :

La regione di copertura ideale ha probabilità di copertura pari a 1 per ogni possibile valore di  $\theta$ . Ovviamente, di solito solo la regione banale,  $\hat{C}_n(Y) = \mathcal{I}$  per ogni  $y$ , ha la proprietà idealmente richiesta. Si desidera in concreto che la regione aleatoria  $\hat{C}_n(Y)$  contenga il bersaglio  $\theta_0$ , quale sia il suo valore in  $\mathcal{I}$ , con una probabilità di copertura assegnata  $1 - \alpha$ , detta livello di copertura, ossia che valga:

$\Pr(\hat{C}_n(Y)) = 1 - \alpha$  per ogni  $\theta$ :

Il livello di copertura va scelto prossimo a 1. L'opzione più corrente è

$1 - \alpha = 0.95$ . Spesso inoltre  $\Pr(\hat{C}_n(Y))$  è solo in via approssimata uguale a  $1 - \alpha$  al variare di  $\theta$ . Si parla allora di regione di copertura con livello approssimato  $1 - \alpha$ . Da Pace, Salvani (2001).

### 1.5.2 Espressione per la probabilità di copertura degli intervalli di copertura per la proporzione binomiale

Per tutti i possibili intervalli di copertura la probabilità di copertura per un dato valore di  $\theta$  è data da:

$C_n(\theta) =$

$\sum_{k=0}^n X_{k=0}$

$I(k; n, \theta)$

n

$I(k; \underline{y}) = \sum_{j=k}^n \mathbb{1}_{\{Y_j \leq k\}}$

dove  $I(k; \underline{y})$  è pari a 1 se l'intervallo contiene  $\underline{y}$  quando  $Y = k$  ed è uguale a 0 se l'intervallo non contiene  $\underline{y}$ .

Oltre alla probabilità di copertura, in molti articoli, viene presa in considerazione l'ampiezza o lunghezza dell'intervallo di confidenza. L'idea che sta alla base è che un intervallo più stretto sia migliore di uno più ampio. Questa ampiezza dipende però dalla parametrizzazione

16 del modello, cioè, facendo scelte diverse per il parametro, si possono avere conclusioni diverse. Per questo l'ampiezza verrà tralasciata per il confronto tra intervalli di confidenza.

## 1.6 Apporto di questa tesi

In letteratura sono disponibili risultati di simulazione per valutare la probabilità di copertura principalmente per gli intervalli Wald, Wald aggiustato e score. Non sembrano invece essere disponibili risultati di simulazione per valutare la probabilità di copertura per gli intervalli basati sulle statistiche  $r$  e  $r_{\text{adj}}$ . Nel seguito verranno presentate, a tal proposito, delle analisi grafiche relative a simulazioni eseguite per gli intervalli basati sulle statistiche  $r$  e  $r_{\text{adj}}$ .

17

18

# Capitolo 2 RISULTATI

## 2.1 Oscillazioni

Come si potrà vedere dai grafici successivi assistiamo ad un fenomeno di oscillazione della probabilità di copertura, sia nel caso di variabile indipendente  $n$  che  $\underline{y}$ . Questo fenomeno è conseguenza del fatto che la distribuzione usata per generare i dati è discreta. Brown, Cai e DasGupta (2001) parlano di valori "fortunati" e "sfortunati" di  $n$  e  $\underline{y}$ . Ci sono coppie, "fortunate", di  $n$  e  $\underline{y}$  per cui la probabilità di copertura è molto vicina o addirittura superiore a quella nominale e altre coppie, "sfortunate", di  $n$  e  $\underline{y}$  la cui probabilità di copertura è molto più bassa del valore nominale.

## 2.2 Dipendenza dal valore del parametro

Le seguenti sono alcune analisi grafiche riguardanti la probabilità di copertura al variare del parametro con  $n$  fissato. Si riportano di seguito i grafici delle probabilità di copertura stimate tramite simulazione con 10000 replicazioni, in funzione di  $\underline{y}$  con  $n$  fissato. Si vedano le Figure 2.1-2.10.

Come si può notare dalle Figure 2.1, 2.2, 2.7, 2.8, 2.9 e 2.10, gli intervalli Wald,  $r$  e  $r_{\text{adj}}$  offrono un buon livello di copertura per valori centrali del parametro, mentre l'intervallo Wald aggiustato e score, Figure 2.3-2.6 sono migliori anche per valori vicini alle estremità dell'intervallo di variazione del parametro, ciò è dovuto al metodo utilizzato.

19

0.0 0.2 0.4 0.6 0.8 1.0  
0.0 0.2 0.4 0.6 0.8

**c.i. wald 95% n=5**

p

coverage

Figura 2.1: Probabilità di copertura per intervallo Wald standard al 95 % al variare di  $\underline{y}$ ,  $n = 5$

0.0 0.2 0.4 0.6 0.8 1.0 0.0 0.2 0.4 0.6 0.8

**c.i. wald 95% n=10**

p

coverage

Figura 2.2: Probabilità di copertura per intervallo Wald standard al 95 % al variare di  $\underline{y}$ ,  $n = 10$

20

0.0 0.2 0.4 0.6 0.8 1.0

0.90 0.92 0.94 0.96 0.98 1.00

**c.i. waldadj 95% n=5**

p

coverage

Figura 2.3: Probabilità di copertura per intervallo Wald aggiustato al 95 % al variare di  $\underline{y}$ ,  $n = 5$

0.0 0.2 0.4 0.6 0.8 1.0

0.94 0.96 0.98 1.00

**c.i. waldadj 95% n=10**

p

coverage

Figura 2.4: Probabilità di copertura per intervalloWald aggiustato al 95 % al variare di  $\theta$ , n = 10

21

0.0 0.2 0.4 0.6 0.8 1.0  
0.80 0.85 0.90 0.95 1.00

**c.i. score 95% n=5**

p  
coverage

Figura 2.5: Probabilità di copertura per intervallo score al 95 % al variare di  $\theta$ , n = 5

0.0 0.2 0.4 0.6 0.8 1.0 0.90 0.92 0.94 0.96 0.98 1.00

**c.i. score 95% n=10**

p  
coverage

Figura 2.6: Probabilità di copertura per intervallo score al 95 % al variare di  $\theta$ , n = 10

22

0.0 0.2 0.4 0.6 0.8 1.0  
0.0 0.2 0.4 0.6 0.8

**c.i. r al 95% n=5**

p  
coverage

Figura 2.7: Probabilità di copertura per intervallo basato sulla statistica r al 95 % al variare di

$\theta$ , n = 5

0.0 0.2 0.4 0.6 0.8 1.0  
0.0 0.2 0.4 0.6 0.8 1.0

**c.i. r al 95% n=11**

p  
coverage

Figura 2.8: Probabilità di copertura per intervallo basato sulla statistica r al 95 % al variare di

$\theta$ , n = 10

23

0.0 0.2 0.4 0.6 0.8 1.0  
0.0 0.2 0.4 0.6 0.8

**c.i. r\* al 95% n=5**

p  
coverage

Figura 2.9: Probabilità di copertura per intervallo basato sulla statistica r\*(ch. paragrafo 1.2.1)

al 95 % al variare di  $\theta$ , n = 5

0.0 0.2 0.4 0.6 0.8 1.0  
0.0 0.2 0.4 0.6 0.8 1.0

**c.i. r\* al 95% n=10**

p  
coverage

Figura 2.10: Probabilità di copertura per intervallo basato sulla statistica r\*(ch. paragrafo 1.2.1)

al 95 % al variare di  $\theta$ , n = 10

## 2.3 Dipendenza dalla numerosità campionaria

Analisi gra\_ che riguardanti la numerosità campionaria e la probabilità di co- pertura al variare della numerosità campionaria. Si riportano di seguito i gra\_ci delle probabilità di copertura

24

stimate tramite simulazione con 10000 replicazioni, in funzione di n con  $\theta$  ssato. Si vedano le

Figure 2.11-2.25.

0 20 40 60 80 100  
0.0 0.2 0.4 0.6 0.8

**c.i. wald 95% p=0.3**

n  
coverage

Figura 2.11: Probabilità di copertura per l'intervallo Wald standard al 95% al variare di n,  $\theta=0.3$

Per tutti gli intervalli si nota chiaramente che il livello di copertura si stabilizza intorno al livello nominale al crescere di n. Per le numerosità meno elevate, gli intervalli che portano più rapidamente a livelli di copertura vicini a quello nominale sono l'intervallo Wald aggiustato proposto da Agresti e Coull (1998), Figura 2.13-2.15 e l'intervallo score, Figura 2.17-2.19, che è, adirittura, conservativo per valori di n moderati e valori del parametro bassi. Gli altri hanno bisogno di numerosità più elevate per arrivare al livello di copertura nominale, per valori centrali del parametro  $\theta$  la numerosità può essere leggermente meno elevata per raggiungere il livello di copertura nominale.

25

0 20 40 60 80 100  
0.0 0.2 0.4 0.6 0.8

**c.i. wald 95% p=0.5**

n  
coverage

Figura 2.12: Probabilità di copertura per l'intervallo Wald standard al 95% al variare di n,  $\theta=0.5$

0 20 40 60 80 100 0.0 0.2 0.4 0.6 0.8 1.0

**c.i. wald 95% p=0.8**

n  
coverage

Figura 2.13: Probabilità di copertura per l'intervallo Wald standard al 95% al variare di  $n$ ,  $\alpha=0.8$

26

0 20 40 60 80 100  
0.94 0.95 0.96 0.97 0.98 0.99 1.00

**c.i. waldadj 95% p=0.3**

n  
coverage

Figura 2.14: Probabilità di copertura per l'intervallo Wald aggiustato al 95% al variare di  $n$ ,

$\alpha=0.3$

0 20 40 60 80 100  
0.94 0.95 0.96 0.97 0.98 0.99 1.00

**c.i. waldadj 95% p=0.5**

n  
coverage

Figura 2.15: Probabilità di copertura per l'intervallo Wald aggiustato al 95% al variare di  $n$ ,

$\alpha=0.5$

27

0 20 40 60 80 100  
0.65 0.70 0.75 0.80 0.85 0.90 0.95 1.00

**c.i. waldadj 95% p=0.8**

n  
coverage

Figura 2.16: Probabilità di copertura per l'intervallo Wald aggiustato al 95% al variare di  $n$ ,

$\alpha=0.8$

0 20 40 60 80 100 0.94 0.95 0.96 0.97 0.98 0.99 1.00

**c.i. score 95% p=0.3**

n coverage

Figura 2.17: Probabilità di copertura per l'intervallo score al 95% al variare di  $n$ ,  $\alpha=0.3$

28

0 20 40 60 80 100  
0.88 0.90 0.92 0.94 0.96 0.98 1.00

**c.i. score 95% p=0.5**

n  
coverage

Figura 2.18: Probabilità di copertura per l'intervallo score al 95% al variare di  $n$ ,  $\alpha=0.5$

0 20 40 60 80 100  
0.65 0.70 0.75 0.80 0.85 0.90 0.95

**c.i. score 95% p=0.8**

n  
coverage

Figura 2.19: Probabilità di copertura per l'intervallo score al 95% al variare di  $n$ ,  $\alpha=0.8$

29

0 20 40 60 80 100  
0.0 0.2 0.4 0.6 0.8 1.0

**c.i. r al 95% p=0.3**

n  
coverage

Figura 2.20: Probabilità di copertura per l'intervallo basato sulla statistica  $r$  al 95% al variare

di  $n$ ,  $\alpha=0.3$

0 20 40 60 80 100  
0.0 0.2 0.4 0.6 0.8 1.0

**c.i. r al 95% p=0.5**

n  
coverage

Figura 2.21: Probabilità di copertura per l'intervallo basato sulla statistica  $r$  al 95% al variare

di  $n$ ,  $\alpha=0.5$

30

0 20 40 60 80 100  
0.0 0.2 0.4 0.6 0.8 1.0

**c.i. r al 95% p=0.8**

n  
coverage

Figura 2.22: Probabilità di copertura per l'intervallo basato sulla statistica  $r$  al 95% al variare

di  $n$ ,  $\alpha=0.8$

0 20 40 60 80 100  
0.0 0.2 0.4 0.6 0.8 1.0

**c.i. r\* al 95% p=0.3**

p  
coverage

Figura 2.23: Probabilità di copertura per l'intervallo basato sulla statistica  $r_{\alpha}$  (ch. paragrafo 1.2.1) al 95% al variare di  $n$ ,  $\alpha=0.3$

31

0 20 40 60 80 100  
0.0 0.2 0.4 0.6 0.8 1.0

**c.i. r\* al 95% p=0.5**

p  
coverage

Figura 2.24: Probabilità di copertura per l'intervallo basato sulla statistica  $r_{\underline{}}$  (ch. paragrafo 1.2.1) al 95 % al variare di  $n$ ,  $\underline{=} = 0.5$

0 20 40 60 80 100  
0.0 0.2 0.4 0.6 0.8 1.0  
c.i.  $r^*$  al 95%  $p=0.8$   
p  
coverage

Figura 2.25: Probabilità di copertura per l'intervallo basato sulla statistica  $r_{\underline{}}$  (ch. paragrafo 1.2.1) al 95 % al variare di  $n$ ,  $\underline{=} = 0.8$

32

## Capitolo 3

# CONCLUSIONI

Dalla Figura 3.1 si nota il diverso livello di copertura al variare di  $\underline{}$  dei cinque intervalli approssimati presi in considerazione. Come riportato in Agresti e Coull (1998), Vollset (1993), Brown, Cai, DasGupta (2001) e altri si nota chiaramente come l'intervallo Wald standard sia quello con probabilità di copertura minore, in generale ha una bassa probabilità di copertura per i valori di  $\underline{}$  vicini a 0 e 1 ed è sempre minore del livello nominale di 0:95 anche per valori centrali di  $\underline{}$ . Già passando alla modi\_ ca proposta da Agresti e Coull (1998) si nota un notevole incremento della probabilità di copertura anche con piccoli campioni. In questo caso l'intervallo ha probabilità di copertura molto più stabile intorno al valore nominale rispetto a quello Wald standard. L'intervallo score è migliore in termini di probabilità di copertura strettamente vicino al valore nominale. I due intervalli basati sulla statistica  $r$  e  $r_{\underline{}}$ , per valori vicini a 0 e 1 hanno una probabilità di copertura molto distante dal valore nominale e si avvicinano  $\underline{}$  a toccarlo, ed in qualche caso superarlo, solo per valori compresi tra 0:2 e 0:8 di  $\underline{}$ .

33

0.0 0.2 0.4 0.6 0.8 1.0  
0.0 0.2 0.4 0.6 0.8 1.0  
p  
coverage

Figura 3.1: Probabilità di copertura al variare di  $p$  per gli intervalli: Wald standard (verde), Wald aggiustato (Agresti-Coull) (giallo), score (blu),  $r$  (nero),  $r_{\underline{}}$  (rosso).

In Figura 3.2 e 3.3 c'è il dettaglio del confronto tra  $r$  e  $r_{\underline{}}$  in cui si vede che l'intervallo basato su  $r_{\underline{}}$  non apporta signifi cativi miglioramenti all'intervallo basato su  $r$ . Adirittura per numerosità molto alte (nell'esempio  $n = 100$ ) l'intervallo basato su  $r_{\underline{}}$ , per valori centrali di  $\underline{}$ , con 10000 simulazioni, è peggiore dell'intervallo basato su  $r$ . Si veda Figura 3.3.

34

0.0 0.2 0.4 0.6 0.8 1.0  
0.80 0.85 0.90 0.95 1.00  
p  
coverage

Figura 3.2: Probabilità di copertura al variare di  $\underline{}$  di  $r$  (nero) e  $r_{\underline{}}$  (rosso),  $n = 10$

0.0 0.2 0.4 0.6 0.8 1.0  
0.80 0.85 0.90 0.95 1.00  
p  
coverage

Figura 3.3: Probabilità di copertura al variare di  $\underline{}$  di  $r$  (nero) e  $r_{\underline{}}$  (rosso),  $n = 100$

35

Per quanto riguarda la numerosità campionaria, dalla Figura 3.4 risulta che l'intervallo score è migliore di quello basato su  $r$  e  $r_{\underline{}}$  per le numerosità più basse e sono praticamente identici a partire da  $n = 40$ .

0 20 40 60 80 100  
0.70 0.75 0.80 0.85 0.90 0.95 1.00  
n  
coverage

Figura 3.4: Probabilità di copertura al variare di  $n$  per score (blu),  $r$  (nero) e  $r_{\underline{}}$  (rosso),  $\underline{=} = 0:5$

36

A conclusione di ciò si ha che: l'intervallo Wald standard, a dispetto della larga di\_ usione e semplicità, è il meno a\_ dabile e quindi può essere sostitu- to con l'intrevallo Wald aggiustato proposto da Agresti-Coull (1998). Anche l'intervallo score o\_ re un ottima alternativa sia al semplice intervallo Wald sia quello Wald aggiustato.  $r$  e  $r_{\underline{}}$  si inseriscono tra l'intervallo Wald standard e l'intervallo score, essi sono molto più soddisfacenti per numerosità campionarie elevate.

37

BIBLIOGRAFIA

- Agresti, A. e Coull, B.A. (1998). Approximate is better than exact for interval estimation of binomial proportion. *The American Statistician*, 52, 119-126.
- Blaker, H. (2000). Con\_dence curves and improved exact con\_dence intervals for discrete distribution. *The Canadian Journal of Statistics*, 28, 783-798.
- Blyth, C.R. e Still, H.A. (1983). Binomial con\_dence intervals. *Journal of the American Statistical Association*, 78, 108-116.
- Borkowf, B.C. (2006). Constructing binomial con\_dence intervals with near nominal coverage by adding a single imaginary failure or succes. *Statistics in Medicine*, 25, 3679-3695.
- Brown, L.D., Cai, T.T. e DasGupta, A. (2001). Interval estimation for a binomial proportion. *Statistical Science*, 16, 101-133.
- Pace, L. e Salvan, A. (2001). *Introduzione alla statistica*, CEDAM
- Reiczigel, J. (2003). Con\_dence intervals for binomial parameter: some new considerations. *Statistic in Medicine*, 22,611-621.
- Severini, T.A. (2000) *Likelihood Methods in Statistics*, Oxford University Press.
- 38
- Vollset, S.E. (1993). Con\_dence interval for a binomial proportion. *Statistics in Medicine*, 12, 809-824.
- Wang, H. (2007). Exact con\_dence coe\_cent of con\_dence interval for a binomial proportion. *Statistica Sinica*, 17, 361-368.
- 39