

# L'errore di misura negli studi retrospettivi in epidemiologia <sup>(1)</sup>

*Measurement errors in Epidemiological Retrospective Studies*

Monica Ferraroni

Dip. di Med. Chirurgia e Odontoiatria - Università di Milano Via A. Di Rudinì, 8 I-20142

Adriano Decarli

Istituto di Statistica Medica e Biometria - Università di Milano Via Venezian, 1 I-20133

Maura Mezzetti

Istituto di Metodi Quantitativi, Università Bocconi Via Gobbi, 5 I-20136 Milano

**Abstract:** The relationships between a set of covariates and a response variable  $Y$  are usually quantified by means of regressive methods. The accuracy of this quantification depends on the precision of covariate measurements. However, in many studies only a subset of explanatory variables  $Z$  is observable without error, while for some other variables  $X$  the information is not directly available. Instead of  $X$ , mis-measured surrogate variables  $W$  related to  $X$  are observable, and consequently the parameters relating  $Y$  and  $X$  cannot be directly estimated. Measurement errors (ME) related to surrogate variables may produce biased results. We propose a procedure which adjusts the estimated association for the effects of misclassification ascribable to ME. The performance of the procedure is evaluated by simulation.

**Parole chiave:** Measurement error, case-control study.

## 1. Premesse

In questo lavoro verranno trattati problemi relativi alla correzione per l'errore di misura (ME) utilizzati nell'analisi di studi di tipo caso-controllo, in cui il legame fra la variabile di risposta  $Y$  e i suoi predittori sia definibile tramite una relazione funzionale  $Y = f(\beta_z'Z, \beta_x'X)$  con  $Z$  insieme dei predittori misurati senza errore e  $X$  insieme dei predittori la cui misurazione diretta senza errore non risulta possibile. In questi studi la variabile  $Y$  è un indicatore che generalmente assume valore 1 per i "casi" (soggetti malati) e 0 per i "controlli" (soggetti sani). La misura vera di  $X$  non è usualmente nota, mentre è osservata  $W$ , associata a  $X$ , chiamata "surrogato" di  $X$ . I parametri del modello non sono direttamente stimabili e la correzione per ME passa attraverso le stime indirette di tali parametri, ottenute tramite un modello in cui  $Y$  risulta funzione di  $(Z, W)$ . La procedura di stima tramite la massima verosimiglianza prevede la determinazione della distribuzione congiunta di  $Y$  e  $W$  dato  $Z$ , che per  $\Pr(Y|W, X, Z) = \Pr(Y|X, Z)$  è possibile scrivere come:

$$f_{Y, W|Z}(y, w|z) = \int \Pr(y|\beta, x, z) \times f_{W, X|Z}(w, x|z) dx \quad (1)$$

<sup>(1)</sup> Il presente lavoro è stato finanziato dal generoso contributo dell'AIRC.

dove  $\beta$  è il vettore dei parametri che esprime la relazione tra la risposta  $Y$  e l'esposizione e  $f_{w,x|z}(w,x|z)$  necessita di essere definita attraverso un modello che specifichi ME<sup>1</sup>.

## 2. Modelli per la definizione dell'errore di misura

I modelli per la definizione di ME vengono generalmente distinti in due classi:

'*Error Models*': in cui  $W$  rappresenta la vera osservazione contaminata da un errore assunto indipendente da  $(X,Z)$  e la distribuzione congiunta di  $W$  e  $X$  dato  $Z$  è definita da  $f_{w,x|z}(w,x|z) = f_{w|x,z}(w|x,z) \times f_{x|z}(x|z)$ . La (1) diventa quindi:

$$f_{Y,W|Z}(y,w|z) = \int_x \Pr(y|\beta, x, z) \times f_{W|X,Z}(w|x,z) \times f_{X|Z}(x|z) dx \quad (2)$$

Le tre componenti della (2) rappresentano rispettivamente: la relazione tra malattia e vera esposizione; il modello per ME, la distribuzione della variabile  $X$ , non osservabile. In caso di errore non differenziale, la distribuzione condizionata  $f_{w|x,z}(w|x,z)$  non dipende da  $Y$ . Questo tipo di modello è sovente utilizzato negli studi epidemiologici quando sia disponibile uno studio di validazione esterno, cioè uno studio in cui per un gruppo di soggetti "sani", non inclusi nello studio principale siano contemporaneamente disponibili la misura "vera"  $X$  e la misura "surrogato"  $W$ .

'*Regression Calibration Models*': in cui si specifica la distribuzione condizionale di  $X$  dato  $(W,Z)$ , ad es:  $X = W + Z + U_*$  con  $E(U_* | W, Z) = 0$ . In questo caso la distribuzione congiunta diviene  $f_{w,x|z}(w,x|z) = f_{x|w,z}(x|w,z) \times f_{w|z}(w|z)$ , e la (1) risulta:

$$f_{Y,W|Z}(y,w|z) = \int_x \Pr(y|\beta, x, z) \times f_{X|W,Z}(x|w,z) \times f_{W|Z}(w|z) dx \quad (3)$$

La terza componente dell'espressione (3) non dipende da  $X$ . E' quindi possibile stimare  $\beta$  dalla (3) massimizzando la verosimiglianza di:

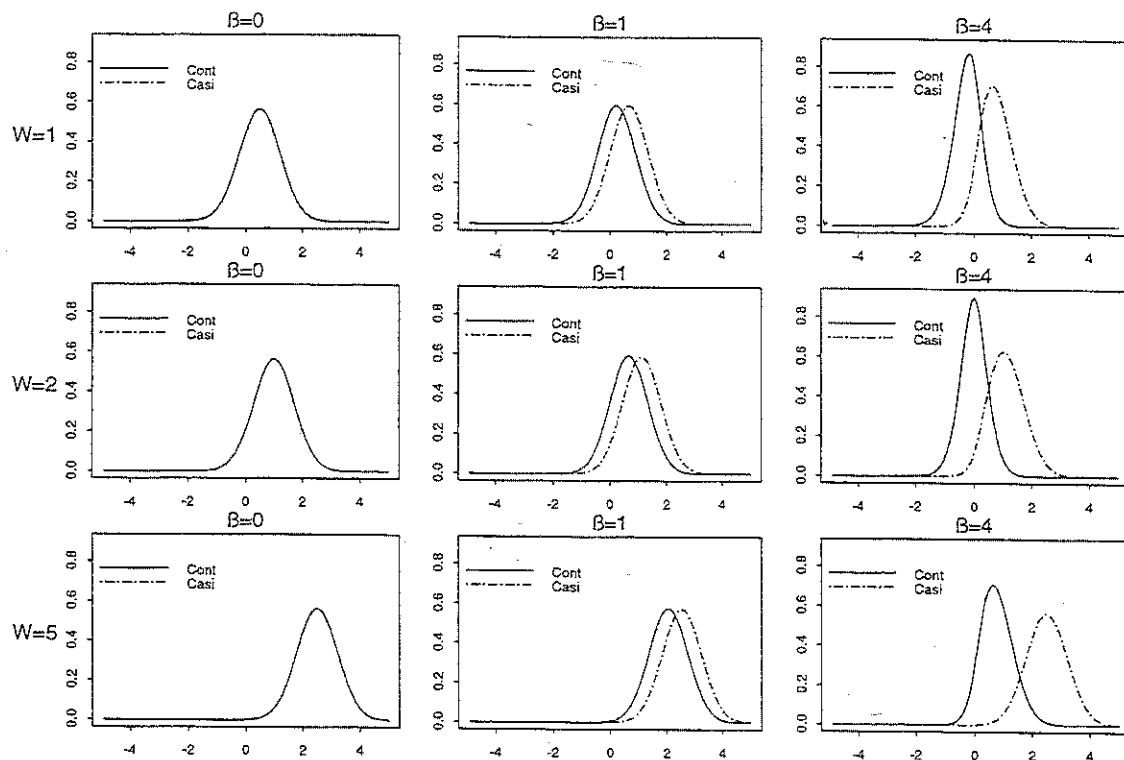
$$f_{Y|W,Z}(y|w,z) = \int_x \Pr(y|\beta, x, z) \times f_{X|W,Z}(x|w,z) dx \quad (4)$$

Con questo tipo di approccio nello studio principale si stimano i parametri di interesse con gli usuali modelli sostituendo  $W$  con  $E(X | W, Z)$  stimato con l'ausilio di uno studio di validazione. In questa classe di modelli la distribuzione  $f_{x|w,z}(x|w,z)$  dipende in realtà anche da  $Y$ . La distribuzione congiunta di  $W$  e  $X$  in assenza di covariate  $Z$  e con una sola variabile  $X$ , può essere scritta come  $f_{w,x}(w,x) = f_{w,x}(w,x = w + u_*) = f_w(w) \times f_{u_*}(x - w)$  per cui la distribuzione condizionata è:

$$f_{X|W,Y}(x|w,y) = \frac{P(y|x) \times f_{U_*}(x-w)}{\int_x P(y|x) \times f_{U_*}(x-w) dx} \quad (5)$$

Nell'ipotesi in cui  $f_{U_i}(x-w) \sim N(0,1)$ , con  $P(Y=1|x)$  espressa da una funzione logistica, la rappresentazione grafica della (5) per alcuni valori del parametro  $\beta$  e della variabile surrogato  $W$  è riportata nella figura 1. Solo per  $\beta=0$  la distribuzione condizionata è la stessa tra i due gruppi mentre al crescere di  $\beta$  aumenta la diversità tra le due distribuzioni. Ciò rende più difficile utilizzare uno studio di validazione esterno, nella fase di correzione per ME. Lo studio di validazione esterno permette infatti una stima di  $f_{X|W,Y}(x|w,y)$  solo per  $y=0$  (soggetti sani).

Figura 1:  $f_{X|W,Y}(x|w,y)$  al variare del valore del parametro  $\beta$  per fissi valori di  $W$ .



### 3. Un modello iterativo

Per stimare in modo corretto gli effetti di interesse ( $\beta_x$ ) in generale sono stati proposti metodi appartenenti alla classe dei 'Regression Calibration Models' che ipotizzano  $f_{X|W}(x|w)$  uguale tra i due gruppi. La procedura proposta non ipotizza  $f_{X|W,Y}(x|w,y) = f_{X|W}(x|w)$  e consente di giungere ad una stima di  $\beta$  in modo iterativo quando sia disponibile uno studio di validazione esterno in cui  $W$  e  $X$  siano congiuntamente rilevati. La procedura è costituita dai seguenti passi:

1. Stima di  $\beta$  utilizzando i dati relativi allo studio principale, attraverso la regressione logistica, come avviene usualmente in questi studi.

2. Stima per ogni soggetto di  $E(X|W,Y) = \int_x x f_{X|W,Y}(x|w,y) dx$ , utilizzando le infor=

mazioni relative allo studio di validazione, con  $f_{X|W,Y}(x|\bar{w},y)$  espressa dalla (5) in cui  $P(y|x)$  è sostituito dalle stime ottenute in 1.

3. Stima di  $\beta_x$  come in 1. utilizzando i valori  $E(X|W,Y)$  ottenuti in 2.:

$$P(Y|E(X|W,Y)) = \frac{[\exp(\alpha + \beta_x^{(i)} E(X|W,Y))]^y}{1 + \exp(\alpha + \beta_x^{(i)} E(X|W,Y))} \quad [y=0,1]$$

4. Ripetizione dei punti 2. e 3 finché  $|\beta_x^{(i)} - \beta_x^{(i-1)}| \sim 0$

#### 4. Validazione dei modelli di correzione per l'errore di misura.

Il metodo proposto è stato saggiato analizzando una serie di 500 studi tipo caso-controllo generati da una popolazione in cui erano noti la distribuzione di  $X$ , la relazione funzionale tra  $X$  e  $Y$  e la struttura dell'errore che lega  $W$  a  $X$ . I risultati ottenuti sono stati confrontati con quelli ottenibili con altri metodi presenti in letteratura.

In tabella 1 si riportano il valor vero di  $\beta$ , quello stimato senza correzione per ME ( $\hat{\beta}_w$ ), quello ottenuto con la correzione proposta da Rosner<sup>2</sup> ( $\hat{\beta}_1$ ), con la correzione proposta da Reeves<sup>3</sup> ( $\hat{\beta}_2$ ) e quello ottenuto con il metodo iterativo esposto nella sezione 3 ( $\hat{\beta}_3$ ). Il modello proposto compensa solo in parte la sottostima dovuta a ME, ma sembra fornire risultati migliori del modello proposto da Rosner. Il modello di Reeves fornisce una stima più prossima al valor vero, benché sovrastimata e con valore più elevato della deviazione standard. Il metodo proposto sarà applicato a dati relativi ad uno studio caso-controllo su abitudini alimentari e tumore alla mammella.

**Tabella 1:** Confronto tra valor vero ( $\beta$ ), valore stimato senza correzione ( $\hat{\beta}_w$ ), valore stimato con il metodo di Rosner ( $\hat{\beta}_1$ ), il metodo di Reeves ( $\hat{\beta}_2$ ) e con il metodo proposto ( $\hat{\beta}_3$ ). Risultati ottenuti su 500 simulazioni (deviazioni standard in parentesi).

$\beta$	$\hat{\beta}_w$	$\hat{\beta}_1$	$\hat{\beta}_2$	$\hat{\beta}_3$
0.01790	0.0111 (0.0018)	0.0132 (0.0021)	0.0184 (0.0024)	0.0152 (0.0019)

#### Riferimenti bibliografici

- <sup>1</sup>Carroll R.J., Gail M.H., Lubin J.H. (1993) Case-control studies with errors in covariates, *Journal of American Statistical Association*, 88, 185-199.
- <sup>2</sup>Reeves G.K., Cox D.R., Darby S.C. and Whitley E. (1999) Some aspects of measurement error in explanatory variables for continuous and binary regression models, *Statistics in Medicine*, 17, 2157-2177.
- <sup>3</sup>Rosner B.A. (1996) Measurement error models for ordinal exposure variables measured with error. *Statistics in Medicine*, 15, 293-303.

**SIS** 2000  
@ds.unifi.it

Società Italiana di Statistica

**XL RIUNIONE SCIENTIFICA**

Firenze, 26-28 aprile 2000

Sessioni plenarie e specializzate: sintesi  
Sessioni comunicazioni spontanee  
Sessioni satellite