

## BIOS 736, Fall 2023: Notes for Class # 2

### WHAT IS REQUIRED FOR AN ACTUAL MEASUREMENT ERROR OR MISCLASSIFICATION ADJUSTMENT?

First, let's return to exposure misclassification

Recall the scenario for a case-control study (independent samples): *true exposure ("gold standard")*

True Cell Probabilities				"True" Cell Counts		
	X				X	
D	1	0	$\Rightarrow$	D	1	0
1	$\pi_1$	$1 - \pi_1$		1	a	b
0	$\pi_0$	$1 - \pi_0$		0	c	d

“Observed” Cell Probabilities			$\Rightarrow$	Observed Cell Counts		
D	W			D	W	
	1	0		1	0	
1	$\pi_1^*$	$1 - \pi_1^*$		1	A	B
0	$\pi_0^*$	$1 - \pi_0^*$		0	C	D

*error-prone version of X*

Again,  $\pi_1 = \Pr(X=1 \mid D=1)$  prob of exposure among cases

$\pi_0 = \Pr(X=1 \mid D=0)$  prob of exposure among controls

Parameter of interest:  $OR = \pi_1(1-\pi_0) / [\pi_0(1-\pi_1)]$

Note that it is  $\pi_1^* = \Pr(W=1 \mid D=1)$  and  $\pi_0^* = \Pr(W=1 \mid D=0)$  that are directly estimable from the observed (D, W) table, and these estimates will generally be biased for  $\pi_1$  and  $\pi_0$  due to misclassification

Allowing for possible differential misclassification, recall these definitions:

Sensitivities:  $SE_d = \Pr(W=1 | X=1, D=d)$

(d=0,1)

Specificities:  $SP_d = \Pr(W=0 | X=0, D=d)$

4 ~~parameters~~  
parameters

If we knew these values, we could directly implement a correction of the “naïve”  $\pi_1^*$  and  $\pi_0^*$  estimates:

$$\begin{aligned}\pi_d^* &= \Pr(W=1 | D=d) = \Pr(W=1, X=1 | D=d) + \Pr(W=1, X=0 | D=d) \\ &= \Pr(W=1 | X=1, D=d) \times \Pr(X=1 | D=d) \\ &\quad + \Pr(W=1 | X=0, D=d) \times \Pr(X=0 | D=d) \\ &= SE_d \pi_d + (1 - SP_d)(1 - \pi_d) \quad (d=0,1)\end{aligned}$$

By solving for  $\pi_d$ , an equivalent identity becomes:

$$\pi_d = \frac{\hat{\pi}_d^* + SP_d - 1}{SE_d + SP_d - 1} \quad (d=0, 1)$$

- NOTE: These two identities are the basis of the so-called “matrix method”, which has roots in classic papers on misclassification (e.g., Bross 1954 *Biometrics*; Barron 1977 *Biometrics*)

- An important point (e.g., Bross 1954) is that in non-differential case,  $H_0: OR = 1 \Rightarrow OR^* = 1$ . Thus, standard inference remains valid based only on the (D,W) table!

Test of  
OR=1 based  
on W vs. D  
table is valid

$$\begin{aligned}SE_1 &= SE_0 = \Pr(W=1 | X=1) \\ SP_1 &= SP_0 = \Pr(W=0 | X=0)\end{aligned}$$

ASIDE: Why do they call it the "matrix method" (e.g., Barron, 1977)?

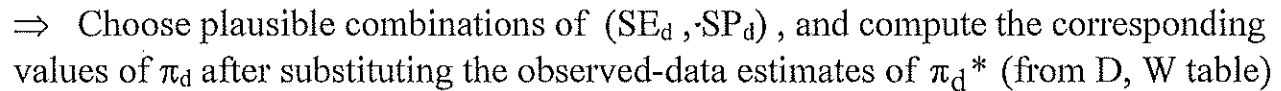
In this case-control setting, we can write:

$$\underbrace{\begin{pmatrix} \pi_1^* \\ 1 - \pi_1^* \\ \pi_0^* \\ 1 - \pi_0^* \end{pmatrix}}_{\tilde{\pi}^*} = \underbrace{\begin{pmatrix} SE_1 & 1 - SP_1 & 0 & 0 \\ 1 - SE_1 & SP_1 & 0 & 0 \\ 0 & 0 & SE_0 & 1 - SP_0 \\ 0 & 0 & 1 - SE_0 & SP_0 \end{pmatrix}}_{\tilde{A}} \underbrace{\begin{pmatrix} \pi_1 \\ 1 - \pi_1 \\ \pi_0 \\ 1 - \pi_0 \end{pmatrix}}_{\tilde{\pi}}$$

$$\Rightarrow \tilde{\pi}^* = \tilde{A} \tilde{\pi}$$

$$\Rightarrow \boxed{\tilde{\pi} = \tilde{A}^{-1} \tilde{\pi}^*}$$

“sensitivity” analysis:



⇒ Provides an idea of the range of possible OR values (but, should consider variability due to uncertainty in the  $\hat{\pi}_d^*$ 's)

But let's assume we have external validation data, of the following typical form:

"gold standard"

$$\begin{aligned}\hat{SE} &= \hat{P}_N(w=1 | x=1) \\ \hat{SP} &= \hat{P}_N(w=0 | x=0)\end{aligned}$$

$$\hat{Var}(\hat{SE}) = \frac{\hat{SE}(1-\hat{SE})}{E+F}, \quad \hat{Var}(\hat{SP}) = \frac{\hat{SP}(1-\hat{SP})}{G+H}, \quad \hat{Cov}(\hat{SE}, \hat{SP}) \cong 0$$

to assume non-differentiability!

Now we can get an actual "corrected" OR estimate:

$$\hat{\pi}_d = \frac{\hat{\pi}_d^* + SP_d - 1}{SE_d + SP_d - 1} \quad (d=0, 1) \Rightarrow \hat{OR} = \frac{\hat{\pi}_1(1-\hat{\pi}_0)}{\hat{\pi}_0(1-\hat{\pi}_1)}$$

And, we could obtain an estimate of  $\text{Var}[\ln(\hat{OR})]$  using (you guessed it) the delta method

Review origin and application of delta method here:

$$\begin{aligned} \ln(\hat{OR}) &= \ln \left\{ \frac{\hat{\pi}_1(1-\hat{\pi}_0)}{\hat{\pi}_0(1-\hat{\pi}_1)} \right\} \\ \Rightarrow \text{var}[\ln(\hat{OR})] &= \text{var} \left[ \ln \left( \frac{\hat{\pi}_1}{1-\hat{\pi}_1} \right) - \ln \left( \frac{\hat{\pi}_0}{1-\hat{\pi}_0} \right) \right] \\ &= \sum_{d=0}^1 \text{var} \left[ \ln \left( \frac{\hat{\pi}_d}{1-\hat{\pi}_d} \right) \right] - 2 \text{cov} \left[ \downarrow, \downarrow \right] \end{aligned}$$

Univariate  $\Delta$ -method results:

$$\text{If } Y = g(X) \Rightarrow \text{var}(Y) \doteq \left( \frac{dg}{dx} \Big|_{x=\mu_X} \right)^2 \text{var}(X)$$

Also, if  $Y_1 = g_1(X_1)$  and  $Y_2 = g_2(X_2)$

$$\begin{aligned} \Rightarrow \text{cov}(Y_1, Y_2) &= \text{cov}[g_1(X_1), g_2(X_2)] \\ &\doteq \left( \frac{dg_1}{dx_1} \Big|_{x_1=\mu_{X_1}} \right) \left( \frac{dg_2}{dx_2} \Big|_{x_2=\mu_{X_2}} \right) \text{cov}(X_1, X_2) \end{aligned}$$

Let's use these results:

$$g(\hat{\pi}_d) = \ln[\hat{\pi}_d(1-\hat{\pi}_d)^{-1}] \quad , \quad d=0, 1$$

$$\Rightarrow \frac{dg}{d\pi_d} = \frac{1}{\pi_d(1-\pi_d)}$$

$$\Rightarrow \text{var}[\ln(\hat{OR})] = \sum_{d=0}^1 \frac{\text{var}(\hat{\pi}_d)}{[\hat{\pi}_d(1-\hat{\pi}_d)]^2} - \frac{2 \cdot \text{cov}(\hat{\pi}_1, \hat{\pi}_0)}{\hat{\pi}_1(1-\hat{\pi}_1)\hat{\pi}_0(1-\hat{\pi}_0)}$$

\* To get the final result, use multivariate  $\Delta$ -method: to get  $\text{var}(\hat{\pi}_d)$  and  $\text{cov}(\hat{\pi}_1, \hat{\pi}_0)$

$$\hat{\pi}_d = g(\underbrace{\hat{\pi}_d^*}_{\text{main study}}, \underbrace{\hat{SE}, \hat{SP}}_{\text{external validity}})$$

$$\hat{\text{var}}(\hat{\pi}_d) = \hat{D}_d \hat{\Sigma} \hat{D}_d'$$

where  $\hat{D}_d$  is estimated vector of 1st derivatives of  $g$  w/respect  $\pi_d^*, SE, SP$ .

$$\text{Also, } \hat{\text{cov}}(\hat{\pi}_1, \hat{\pi}_0) = \hat{D}_1 \hat{\Sigma} \hat{D}_0'$$

where  $\hat{\Sigma} = \text{var} \begin{pmatrix} \hat{\pi}_d^* \\ \hat{SE} \\ \hat{SP} \end{pmatrix}$   
 $3 \times 3$

$\Rightarrow$  Now can get  $\hat{SE}[\ln(\hat{OR})]$ ,  
 $\sim 95\% \text{ CI}$

Now, what if the validation data are internal:

⇒ X measured (in addition to D and W) on a random subsample of your subjects

*gold standard*

Let's consider the pros/cons of external vs. internal validation designs:

*SE, SP that apply in external pop. n are the same as those that apply in your study!*

Validation Study Characteristics

	Advantages	Disadvantages
<u>External</u>	CHEAP, EASY	MAY NOT BE AVAILABLE AND/OR "TRANSPORTABLE"; NOT STATISTICALLY EFFICIENT; LIKELY REQUIRES NON-DIFFERENTIAL ERROR ASSUMPTION
<u>Internal</u>	CAN BE COSTLY, TIME CONSUMING	SOLVES "TRANSPORTABILITY" PROBLEM; STATISTICALLY EFFICIENT; NON-DIFFERENTIALITY NOT REQUIRED

⇒ Different data layout than for the external case:

			Internal Validation subsample				
				D=1		D=0	
Main study sample			W	X=1	X=0	X=1	X=0
D	W=1	W=0					
1	n <sub>11</sub>	n <sub>12</sub>	1	n <sub>13</sub>	n <sub>14</sub>	n <sub>03</sub>	n <sub>04</sub>
0	n <sub>01</sub>	n <sub>02</sub>	0	n <sub>15</sub>	n <sub>16</sub>	n <sub>05</sub>	n <sub>06</sub>

$$n_{dj} \quad (d=0,1), j=1,\dots,6$$

\* For this study design, we could specify and maximize the likelihood function accounting for all 12 types of observations (12 different cell counts)

ex) What is the likelihood contribution for each of the  $n_{11}$  subjects in the main study with  $(D=1, W=1)$ ?

$$\pi_d^* = \Pr(W=1 | D=1) = SE_1 \pi_1 + (1 - SP_1)(1 - \pi_1) \quad (\text{see pg. 2})$$

↓ could do for all 12 cell counts

- Note: formulating the likelihood based on such “matrix-method” identities [i.e., parameterizing in terms of  $(\pi_d, SE_d, SP_d)$ ] leads to a function that can be maximized numerically, but does not appear to yield closed-form MLEs.

In this case, an alternative parameterization [in terms of  $(\pi_d^*, PPV_d, NPV_d)$ ] turns out to be useful

- Can show (Lyles 2002 *Biometrics*) that the ML estimates for this setting are identical to the so-called “inverse matrix” estimators (Marshall 1990 *J Clinical Epidemiology*)

The “inverse matrix” method is based on a similar but different identity to that underlying the “matrix” method:

$$\begin{aligned} PPV_d &= \Pr(X=1 | W=1, D=d) \\ NPV_d &= \Pr(X=1 | W=0, D=d) \end{aligned}$$

$$\begin{aligned} \pi_d &= \Pr(X=1 | D=d) = \Pr(X=1, W=1 | D=d) + \Pr(X=1, W=0 | D=d) \\ &= \Pr(X=1 | W=1, D=d) \times \Pr(W=1 | D=d) \\ &\quad + \Pr(X=1 | W=0, D=d) \times \Pr(W=0 | D=d) \\ &= PPV_d \pi_d^* + (1 - NPV_d)(1 - \pi_d^*) \quad (d=0,1) \end{aligned}$$

where  $PPV_d = \Pr(X=1 | W=1, D=d)$  and  $NPV_d = \Pr(X=1 | W=0, D=d)$  are the “positive and negative predictive values”



Let's enumerate the likelihood contributions for each of the separate types of observations based on this second parameterization: ( $D=1$ )

# obs	obs. type	Probability (likelihood contribution)
$n_{11}$	$W=1, D=1$	$P_1(W=1 D=1) = \pi_1^*$
$n_{12}$	$W=0, D=1$	$(1 - \pi_1^*)$
$n_{13}$	$W=1, X=1, D=1$	$PPV_1 \pi_1^*$
$n_{14}$	$W=1, X=0, D=1$	$(1 - PPV_1) \pi_1^*$
$n_{15}$	$W=0, X=1, D=1$	$(1 - NPV_1) (1 - \pi_1^*)$
$n_{16}$	$W=0, X=0, D=1$	$NPV_1 (1 - \pi_1^*)$
		$\downarrow$ $NPV_1 = P_1(X=0 W=0, D=1)$

question  
(sum to 1)?

$\Rightarrow$  up to a constant,

$\log L$  is

$$l = \sum_{d=0}^1 \left\{ n_{d1} \ln(\pi_d^*) + n_{d2} \ln(1 - \pi_d^*) + \dots \text{" (4 terms)} \right\}$$

$\Rightarrow$  take

$$\begin{cases} \frac{\partial l}{\partial PPV_d} = \dots = 0 \\ \frac{\partial l}{\partial NPV_d} = \dots = 0 \\ \frac{\partial l}{\partial \pi_d^*} = \dots = 0 \end{cases}$$

Can take those 3 derivatives easily if you want (results on next pg.)

Based on the study design in the previous table and assuming differential misclassification, the MLE's for the parameters involved in the "inverse matrix method" identity are:

$$\hat{\pi}_d^* = \frac{n_{d1} + n_{d3} + n_{d4}}{n_d}, \quad \hat{PPV}_d = \frac{n_{d3}}{n_{d3} + n_{d4}}, \quad \text{and} \quad \hat{NPV}_d = \frac{n_{d6}}{n_{d5} + n_{d6}} \quad (d=0,1),$$

where  $n_d = n_{d1} + n_{d2} + n_{d3} + n_{d4} + n_{d5} + n_{d6}$

\* We can then use the preceding inverse matrix identity to get the MLEs for  $\pi_d$  and for the OR, and again can use the delta method (or the observed information matrix) to estimate  $\text{Var}[\ln(\text{OR})]$ :

$$\hat{\pi}_d = \hat{PPV}_d \hat{\pi}_d^* + (1 - \hat{NPV}_d)(1 - \hat{\pi}_d^*)$$

\* The MLE for  $\pi_d$  is same as Marshall's "inverse matrix" estimator.

\* In differential case,  $\hat{\pi}_1 \perp \hat{\pi}_0 \rightarrow$  because of case-control sampling design

$$\Rightarrow \hat{\text{var}}[\ln(\hat{OR}_{ML})] = \sum_{d=0}^1 \frac{\hat{\text{var}}(\hat{\pi}_d)}{[\hat{\pi}_d(1 - \hat{\pi}_d)]^2}$$

↓  
Marshall (1990),  
Morrissey & Spiegelman (1999)  
Lyles (2002) → HW 1

no cov( $\hat{\pi}_1, \hat{\pi}_0$ ) here.

Interesting note: For this main study / internal validation design, there is no known closed-form MLE for the OR when the misclassification is assumed to be *non-differential*

⇒ Obtain MLE by numerical maximization of the likelihood function

OR

⇒ Inverse variance-weighted  $\ln(\text{OR})$  estimators – weighting the estimators for  $\ln(\text{OR})$  based on i) the internal validation (D,X) data and ii) the main study + validation study data where the latter are treated “as external”

(e.g., Greenland 1987 *Stats in Medicine*; Thurston et al. 2005 *J Stat Planning and Inference*; Lyles et al. 2007 *Epidemiology*)

Basic idea of Greenland's (1987) closed-form “weighted” estimator:

$$\hat{\theta} = \hat{w} \hat{\theta}_I + (1 - \hat{w}) \hat{\theta}_E$$

↓

$\log(\hat{OR})$

(almost as efficient as MLE!)

(w is an inverse-variance weight)

$\hat{\theta}_I = \ln(\hat{OR})$  based only on (D,X) pairs from internal validation sample

$\hat{\theta}_E = \ln(\hat{OR})$  based on (D,w) pairs in the main study, combined with the (X,w) pairs in the validation sample, treating it like an external validation sample.

(see Hw # 1)

### Real Data Example

Consider the following cell counts from a case-control study of SIDS, where exposure (X) is maternal use of antibiotics during pregnancy:

			Validation subsample				
			D=1		D=0		
Main study sample			W	X=1	X=0	X=1	X=0
D	W=1	W=0					
1	122	442	1	29	22	21	12
0	101	479	0	17	143	16	168

\*Refs: (Drews et al. 1990 *Int J Epidemiol*; Greenland 1988 *Stats in Med*; Marshall 1990 *J Clin Epidemiol*; Morrissey and Spiegelman 1999 *Biometrics*)

Note: "Naïve" OR estimate is  $\frac{122 \times 479}{442 \times 101} = 1.31$

⇒ "Naïve":  $\ln(\hat{OR}) = 0.269$ ,  $SE[\ln(\hat{OR})] = 0.150$

usual Woolf's method  

$$\sqrt{\frac{1}{A} + \frac{1}{B} + \frac{1}{C} + \frac{1}{D}}$$

In contrast, the MLE ("inverse matrix") estimator accounting for (assumed differential) misclassification is based on:

$$\hat{\pi}_1^* = \frac{122 + 29 + 22}{775} = 0.223, \quad \hat{PPV}_1 = \frac{29}{29 + 22} = 0.568, \quad \hat{NPV}_1 = \frac{143}{17 + 143} = 0.894,$$

$$\hat{\pi}_0^* = \frac{101 + 21 + 12}{797} = 0.168, \quad \hat{PPV}_0 = \frac{21}{21 + 12} = 0.636, \quad \hat{NPV}_0 = \frac{168}{16 + 168} = 0.913$$

$$\Rightarrow \quad \hat{\pi}_1 = 0.568 \times 0.223 + (1 - 0.894) \times (1 - 0.223) = 0.209$$

$$\hat{\pi}_0 = 0.636 \times 0.168 + (1 - 0.913) \times (1 - 0.168) = 0.179$$

$$\Rightarrow \quad \text{vs. } 1.50 \quad \text{adjusted (vs. } 1.31 \text{)} \quad \text{unadjusted}$$

$$\hat{OR} = 1.21, \quad \hat{\ln(OR)} = 0.192$$

Also get  $SE[\ln(OR)] = 0.221$  via delta method or observed information matrix  
Likelihood ratio test for nondifferentiality can also be conducted ( $p=0.14$  here)

w/ Carey Denny-Botsch

(See, e.g., Lyles et al. 2007 *Epidemiology* for further analysis of this example)

Allows using both internal and external validation data together, and testing for "transportability"

used variations on Greenland's wtd avg. estimator

**NOTE:** How could one extend the "matrix method" idea to the cross-sectional study design case, with both X and Y misclassified?

As in Barron (1977), interested in OR based on  $2 \times 2$  table, with cross-sectional sampling.

		X	
		1	0
Y	1	$\pi_{11}$	$\pi_{01}$
	0	$\pi_{10}$	$\pi_{00}$

X = true exposure

Y = true outcome

$$\pi_{xy} = \Pr(X=x, Y=y)$$

$$(x, y) = (0, 1)$$

Let  $X^*$  and  $Y^*$  be misclassified variables

		$X^*$	
		1	0
$Y^*$	1	$\pi_{11}^*$	$\pi_{01}^*$
	0	$\pi_{10}^*$	$\pi_{00}^*$

$$\pi_{xy}^* = \Pr(X^*=x, Y^*=y)$$

$$(x, y) = (0, 1)$$

Note:

$$\pi_{xy}^* = \sum_{i=0}^1 \sum_{j=0}^1 \Pr(X^*=x, Y^*=y, X=i, Y=j)$$

Suppose make these assumptions:

a)  $\Pr(X^*|X, Y^*, Y) = \Pr(X^*|X) = SE_x$

b)  $\Pr(Y^*|Y, X^*, X) = \Pr(Y^*|Y) = SE_y$

"independent and non-differential" misclassification.

Note, 1st term (of 4) for  $\pi_{11}^*$  is

$$\begin{aligned}
 & P_2(x^*=1, y^*=1, x=1, y=1) \\
 &= P_2(x^*=1 | x=1, y^*=1, y=1) P_2(y^*=1 | y=1, x=1) \\
 &= SE_x SE_y \pi_{11}
 \end{aligned}$$

Similarly, other 3 terms:

$$\begin{aligned}
 & SE_x (1 - SP_y) \pi_{10} \\
 & (1 - SP_x) SE_y \pi_{01} \\
 & (1 - SP_x) (1 - SP_y) \pi_{00}
 \end{aligned}$$

$\Rightarrow \pi_{11}^*$  is sum of these 4 terms.

Barron completed this exercise and wrote:

$$\text{let } \tilde{\pi}^* = \begin{pmatrix} \pi_{11}^* & \pi_{10}^* \\ \pi_{01}^* & \pi_{00}^* \end{pmatrix}, \quad \tilde{\pi} = \begin{pmatrix} \pi_{11} & \pi_{10} \\ \pi_{01} & \pi_{00} \end{pmatrix}$$

$$\begin{aligned}
 \Rightarrow \tilde{\pi}^* &= \tilde{A}' \tilde{\pi} \tilde{B} \\
 &= \begin{pmatrix} SE_x & 1 - SP_x \\ 1 - SE_x & SP_x \end{pmatrix} \begin{pmatrix} \pi_{11} & \pi_{10} \\ \pi_{01} & \pi_{00} \end{pmatrix} \begin{pmatrix} SE_y & 1 - SE_y \\ 1 - SP_y & SP_y \end{pmatrix} \\
 &\quad \tilde{A}' \quad \tilde{\pi} \quad \tilde{B}
 \end{aligned}$$

"matrix method"!

Tang et al. (2013)  
extended matrix to  
relax assumptions

$$\tilde{\pi}_{16} = (\tilde{A}')^{-1} \tilde{\pi}^* \tilde{B}^{-1}$$



