#### Statistical Methods, part 1 Module 2: Latent Class Analysis of Survey Error Models for measurement errors, continued

Dan Hedlin Stockholm University November 2012

Acknowledgement: Paul Biemer



Regression coefficients in the presence of measurement errors

• 
$$Var(\bar{y}) = \frac{1}{R} \frac{\sigma^2_{\tau}}{n}$$

- 'Effective sample size' is not *n*, it is *Rn*
- Recall that the estimated variance is largely ok, despite the true variance is enlarged by measurement errors
- In regression analysis measurement errors in y will enlarge residuals but not introduce bias

### Bias in regression

- Measurement error in x will:
  - Attenuate slope
  - Attenuate correlation
  - Make intercept larger, if all variables are non-negative and slope positive
- Often hard to guess if measurement error affects analysis a lot or only modestly
- Sometimes hard even to guess direction, let alone how severely measurement errors influence estimates

# From now on focus on dichotomous variables

- Y is either 1 or 0
- Interested in proportions
- True score is  $prob(y_i=1)$  (i.e. =  $E(y_i)$ )
- m=2 yields a 2x2 table:



#### **Shortcut Computation**



$$\frac{\sum_{i=1}^{n} (y_{i1} - y_{i2})^2}{n} = \frac{b+c}{n} = g \qquad ("Gross difference rate")$$
$$\boxed{s_2^2 = \frac{b+c}{2n} = \frac{g}{2}}$$

#### Inconsistency (or "Unreliability") Ratio

Thus, when 
$$\frac{(N-1)}{N} \approx 1$$
,

Inconsistency Ratio is

$$I = \frac{SRV}{SV + SRV}$$

Reliability Ratio

$$R = 1 - I$$

Notes on *I*:

- 1. Estimation of *I* or *R* requires two measurements or observations from the response distribution
- 2. 0<*I*<1

small I (or R is large) ==> small measurement error
large I (or R is small) ==> large measurement error

3. U.S. Census Bureau "Rule of Thumb"  $0.0 \le I < 0.2 \Longrightarrow$  good reliability ( $R \ge .8$ )  $0.2 \le I < 0.5 \Longrightarrow$  moderate reliability ( $.5 \le R < .8$ )  $0.5 \le I < 1.0 \Longrightarrow$  poor reliability ( $0 \le R < .5$ )

U.S. Bureau of the Census (1985)

#### Estimation of I from Reinterview

Recall

$$SRV = \frac{b+c}{2n} = \frac{g}{2}$$
$$SV + SRV = p_1q_1 \text{ or } p_2q_2$$
$$or \frac{p_1q_1 + p_2q_2}{2} \text{ or } \frac{p_1q_2 + p_2q_1}{2}$$

Therefore, we can estimate *I* (or *R*) in a number of ways:

or  

$$\hat{I}_{1} = \frac{g}{2p_{1}q_{1}}$$
or  

$$\hat{I}_{2} = \frac{g}{p_{1}q_{1} + p_{2}q_{2}}$$

$$\hat{I}_{3} = \frac{g}{p_{1}q_{2} + p_{2}q_{1}}$$
Best (Census Bureau  
1985)  
Where  $g = \frac{b+c}{n}$  (gross difference rate)

The estimate of *I* denoted by  $\hat{I}_3$  is called the "index of inconsistency"

It can be shown (see Hess, Singer, and Bushery, 1999) that  $\kappa = 1 - \hat{I}_3$  is identical to Cohen's kappa statistic

$$\kappa = \frac{P_0 - P_e}{1 - P_e}$$
  
where  $P_0 = \frac{a + d}{n}$  and  $P_e = p_1 p_2 + q_1 q_2$   
Why is this result remarkable?

### NOTE:

$$E(g) = E\left[\frac{\sum_{i=1}^{n}(y_{i1} - y_{i2})^{2}}{n} | s\right]$$
  
= Var  $y_{i1} | s$ ) + Var  $y_{i2} | s$ ) - 2Cov $(y_{i1}, y_{i2} | s)$   
Thus, if Cov $(y_{i1}, y_{i2} | s)$  < 0, then  $\frac{g}{2}$  will overestimate *SRV*.  
Further, if Cov $(y_{i1}, y_{i2} | s)$  > 0,  $\frac{g}{2}$  will underestimate *SRV*.  
What might lead to Cov $(y_{i1}, y_{i2} | s)$  > 0? Cov $(y_{i1}, y_{i2} | s)$  < 0?

#### The Classification Probability Model

(Bross 1954, Biemer and Stokes, 1991)

Next, we consider a better specified model for measurement error in categorical data, referred to as the *classification probability error model*.

Recall that

 $E_j(\bullet | s)$  = expectation across the "columns" of responses for a given unit on the row of the response probability matrix.

 $E_s E_j (\Box i)$  = expectation across "rows" or all possible selections of units or samples of units from the rows.

### Terminology and notation

- \$\overline{\phi\_i}\$ referred to as probability of a false positive
- $1-\phi_i$  referred to as "specificity"
- $\theta_i$  is probability of a false negative
- $1-\theta_i$  referred to as "sensitivity"
- $\pi = P(\mu_i = 1) =$  true population proportion

$$P_i = E_j(y_{ij} | s) = \Pr(y_{ij} = 1 | i)$$
  
$$\mu_i = \text{true value}$$

We now extend that notation

Again, consider  

$$y_{ij} = \begin{cases} 1 \text{ if } yes \text{ or "+"} \\ 0 \text{ if } no \text{ or "-"} & \text{"phi"} = p \text{ for positive} \end{cases}$$

$$P(y_{ij} = 1 | \mu_i = 0) = \phi_i \text{ false positive probability}$$

$$P(y_{ij} = 0 | \mu_i = 1) = \theta_i \text{ false negative probability}$$

Thus, we have the following for unit *i* 

 $y_{ij}$  $\mu_i$  $P(y_{ij} | \mu_i)$ 00 $1-\phi_i$  specificity10 $\phi_i$ false positive probability01 $\theta_i$ false negative probability11 $1-\theta_i$  sensitivity

 $E_{j}(y_{ij} | s) = P(y_{ij} = 1 | s)$ 

$$= P(y_{ij} = 1 | \mu_i = 1) P(\mu_i = 1) + P(y_{ij} = 1 | \mu_i = 0) P(\mu_i = 0)$$

 $= (1 - \theta_i)\mu_i + \phi_i(1 - \mu_i)$ 

Let 
$$E_s(\theta_i) = \theta$$
  
 $E_s(\phi_i) = \phi$ 

#### Then

$$E(y_{ij}) = E_s E_j(y_{ij} \mid s)$$
$$= E_s [1 - \theta_i)\pi + \phi_i (1 - \pi)]$$
$$= (1 - \theta)\pi + \phi(1 - \pi)$$

### Bias of p

From these results,

Bias
$$(y_{ij}) = E(y_{ij}) - \pi = (1 - \theta)\pi + \phi(1 - \pi) - \pi$$
  

$$= -\theta\pi + \phi(1 - \pi)$$
Let  $p = \sum_{i=1}^{n} \frac{y_i}{n}$ 
Bias $(p) = \sum_{i=1}^{n} \frac{[E(y_i) - \pi]}{n} = -\theta\pi + \phi(1 - \pi)$ 

#### When is the bias 0?

Stockholm University, autumn semester

٦

#### Bias = 0 if and only if either:

$$\theta \pi = \phi(1 - \pi)$$

or

$$\theta = \phi = 0$$

### Variance of p

• Note that  $P_i$  from Census Model can be written as  $P_i = (1 - \theta_i)$  if  $\mu_i = 1$ 

$$=\phi_i$$
 if  $\mu_i=0$ 

or

$$P_i = (1 - \theta_i)\mu_i + \phi_i(1 - \mu)$$

#### Variance of *p* Under the Classification Probability Model

Assumptions: Same as Census Bureau Model Recall, under Census Bureau Model for m = 1

$$\operatorname{Var}(p) = (1-f) \frac{S_1^2}{n} + \frac{S_2^2}{n}$$

where

$$S_{1}^{2} = \sum_{i=1}^{N} \frac{(P_{i} - P)^{2}}{N - 1}$$
$$S_{2}^{2} = \sum_{i=1}^{N} \frac{P_{i}Q}{N}$$

### Rewriting SV and SRV

Let's find  $S_1^2$  (or *SV*) and  $S_2^2$  (or *SRV*) Under this Model

1. 
$$S_1^2 = \sum_{i=1}^N \frac{(P_i - P)^2}{N - 1}$$
  
 $P_i = (1 - \theta_i) \mu_i + \phi_i (1 - \mu_i)$   
 $P = (1 - \theta) \pi + \phi(1 - \pi)$ 

#### Rewriting SV and SRV (cont'd)

$$\boldsymbol{\Theta} = \frac{1}{N_1} \sum_{i=1}^{N_1} \boldsymbol{\Theta}_i, N_1 = \sum_{i=1}^{N} \boldsymbol{\mu}_i$$

$$\phi = \frac{1}{N_0} \sum_{i=1}^N \phi_i, N_0 = N - N_1$$

$$\pi = \frac{1}{N} \sum_{i=1}^{N} \mu_i$$

### Rewriting SV and SRV (cont'd)

Therefore, after some algebra

$$\left(\frac{N-1}{N}\right)SV = \pi(1-\pi)(1-\theta-\phi)^2 + \gamma_{\theta\phi}$$

where

$$\gamma_{\theta\phi} = \pi \sigma_{\theta}^2 + (1 - \pi) \sigma_{\phi}^2$$

(Note: Often,  $\gamma_{\theta\phi}$  is assumed to be negligible.)

#### Rewriting SV and SRV (cont'd)

2. 
$$S_2^2 = \frac{1}{N} \sum_{i=1}^{N} P_i Q_i = SRV$$

Where

$$P_i = (1 - \theta_i)\mu_i + \phi_i(1 - \mu_i)$$

After some algebra,

$$SRV = \pi \theta (1 - \theta) + (1 - \pi) \phi (1 - \phi) - \gamma_{\theta \phi}$$

#### Remarks:

1. If f can be ignored,

$$nVar(p) = SV + SRV$$
  
=  $\pi(1 - \pi)(1 - \theta - \phi)^2$   
+  $\pi\theta(1 - \theta) + (1 - \pi)\phi(1 - \phi)$ 

2. This is useful for studying how false +'s and false -'s affect variance

#### Remarks: (cont'd)

- 3. *R* (and *I*) vary as  $\pi$  varies.
- 4. When  $\pi$  is very small even a small false positive error rate can be important.

Eg. 
$$\pi = .01$$
  $\theta = 0$ 

$$\varphi = .01 \qquad \gamma_{\theta \varphi} = 0$$

Then I = .50 or 50%

#### Summary

$$\left(\frac{N-1}{N}\right)SV = \pi(1-\pi)(1-\theta-\phi)^2 + \gamma_{\theta\phi}$$

$$SRV = \pi\theta(1-\theta) + (1-\pi)\phi(1-\phi) - \gamma_{\theta\phi}$$

$$I = \frac{\pi\theta(1-\theta) + (1-\pi)\phi(1-\phi) - \gamma_{\theta\phi}}{PQ}$$

## This expression illustrates the difficulty of interpreting *I* for categorical data.

### Estimation of SRV and I

Note: As we showed before,  $\frac{g}{2} = \frac{b+c}{2n}$  is <u>unbiased</u>

for *SRV* provided assumptions are satisfied. Recall the key assumptions are:

- Independence of classification errors
- Equal classification probabilities across trials or repeated measurements

### **Estimation of Bias**

Suppose reinterview is truth. This is usually assumed for reconciled reinterview



#### Estimation of False Positive and False Negatives





$$\hat{\phi} = \frac{b}{b+d}$$
 estimates  $\phi$ , false positive

$$\hat{\theta} = \frac{c}{a+c}$$
 estimates  $\theta$ , false negative

Stockholm University, autumn semester

#### Example



#### Problems with Reconciled Reinterview

- 1. Assumes agreements are correct.
- 2. Assumes disagreements can be reconcile accurately.
- 3. Biemer & Forsman; Sinclair & Gastwirth show this is not true for many surveys.

See Forsman and Schreiner (1991) for a good discussion of the issues with reinterview

### Question

• For test-retest reinterview, what does a significant net difference rate indicate?

 How should the net difference rate be interpreted if the second measurement is from a "preferred" survey process?