36-303: Sampling, Surveys and Society

Quality in Surveys Brian Junker 132E Baker Hall brian@stat.cmu.edu

Handouts

- Lecture Notes
- Examples of I.1 Proposals [due Tue Jan 25]

20 January 2011

Outline

20 January 2011

- Quality in Surveys
- Project Proposals: I.1 on the "Project Schedule" handout.
- Reading:
 - □ Up to today: responsible for Groves Ch's 1, 2, 3
 - Next week:
 - Groves Ch 5
 - Groves, Ch 11 (sections 1-6)
 - Groves Ch 4 (sections 1-3; we will do more later) in that order
- Lecture notes online at

http://www.stat.cmu.edu/~brian/303

Quality in Surveys





Process Perspective on Surveys

Quality Perspective on Surveys

1

2

Quality Overview Representation leasuremen Total Survey Error Construct Target Populatio Each of the Quality Coverage Error Validity Components has a Sampling Frame VC verbal description and a statistical formulation Sampling Error Measuremen Error The Quality Components are properties of Respon onrespons Error individual survey design Respondent Processing Error and analysis decisions, Adjustment Error not of whole surveys Response Our job is to make Postsurve Adjustments decisions to minimize Survey Statistic error / maximize quality

Measurement Quality



- Working down the left side:
 - Validity
 - Measurement Error
 - Processing Error

20 January 2011

Some Notation...

- μ_i = value of the <u>construct</u>. E.g. # of doctor visits for ith person in population, i=1, ..., N
- Y_i = <u>ideal value</u> of the <u>measurement</u> for the ith person in the sample, i=1, ..., n
- y_i = <u>observed value</u> (reported number of doctor visits) for ith sample person
- y_{ip} = <u>observed value after editing/processing</u>
- y_{it} = value on the tth "trial" (tth time we run the survey)

Validity

20 January 2011

- $Y_i = \mu_i + \epsilon_i$
 - \square μ_i is the "true value" for the population
 - □ Y_i is the "ideal measured" value
 - ϵ_i is how much Y_i "deviates" from μ_i
 - Deviation/error is natural. We just have to account for it
- If there are T trials (repeats of the survey), t=1, ..., T, we might write

 $Y_{it} = \mu_i + \epsilon_{it}$

And expect that the errors $\epsilon_{\rm it}$ would "average out" over trials...

• A measure of the size of the errors ϵ_i is Corr(Y_i, μ_i)

This correlation is a measure of the Validity of the measurement

5

20 January 2011

Measurement Error

- y_i Y_i is the measurement error
 - □ Y_i is the ideal measurement
 - \Box y_i is the observed measurement
- There are two kinds of measurement error to worry about
 - □ <u>Variability</u>: $y_i = Y_i + error_i$, and the error "averages out" over repeated trials: $E_t[y_{it}] = Y_i$
 - □ <u>*Bias*</u>: $y_i = Y_i + \text{something that doesn't "average out": <math>E_t[y_{it}] \neq Y_i$

Processing Error

- y_{ip} y_i is the processing error
 - □ y_{ip} is the response after editing/processing
 - $\hfill\square$ y_i is the 'raw' response to the measurement
- These errors come in when you have to code, check, or fix survey responses, e.g.
 - Coding a verbal response
 - Range check can this person have been in High School for 7 years?
 - □ Clumping, e.g. "income between \$10,000 and \$30,000"
- These are generally <u>bias</u> and not <u>variability</u> issues

20 January 2011

Representation Quality



- Working down the right side:
 - Coverage Error
 - Sampling Error
 - Nonresponse
 Error (later lecture)
 - Adjustment Error

Coverage Error

- N = total Target Population (size)
- C = target population covered in frame
- U = target population missed by frame
- \overline{Y} = mean of target population
- \overline{Y}_C = mean of covered population
- \overline{Y}_U = mean of uncovered population
- $\overline{Y}_C \overline{Y}$ = <u>coverage error</u>
 - Also called <u>Coverage Bias</u>

20 January 2011

11

10

Ineligible unit



20 January 2011

Sampling Error

- How well does the sample represent the sampling frame?
 - Sampling bias
 - Best to try to anticipate and avoid
 - Can be looked at similarly to coverage bias
 - Another way to deal with is with weights, but this can introduce "adjustment error" (more in a couple pages)
 - Sampling variability this is a more familiar issue! (see next page)

Coverage Error/Coverage Bias

- Suppose we are interested in Monthy Mortgage Payment (\$0 if you rent)
 - □ Total population is all adults in (US/Pgh/...)
 - Data collection method is random digit dialling
 - Sampling frame is callable land-line phone #'s
- Renters may be more likely to have only a cell phone than homeowners
 - Renters are undercovered by our frame
 - Our estimate of mean mortgage payment will be too high
 - If we can get an estimate of $\frac{U}{N}(\overline{Y}_C \overline{Y}_U)$ Then we can estimate $\overline{Y}_C - \overline{Y}$ and fix the bias!

20 January 2011

Sampling Variability

- $\overline{\overline{y}}_s = \frac{1}{n_s} \sum_{i=1}^{n_s} y_{si}$ is the mean of the sample
- $\overline{Y}_C = \frac{1}{C} \sum_C Y_i^C$ is the mean of the frame

The Standard Error for estimating \overline{Y}_C with \overline{y}_s is

$$SE = \sqrt{rac{1}{S}\sum_{s=1}^{S}(\overline{y}_s - \overline{Y}_C)^2}$$

in case of $simple\ random\ sampling\ (next\ week!)\ we know that$

$$SE = SD/\sqrt{n_s} = \frac{\sqrt{\frac{1}{n_s - 1} \sum_{i=1}^{n_s} (y_{si} - \overline{y}_s)^2}}{\sqrt{n_s}}$$

13

14

Adjustment Error

- This usually comes in the forms of weights.
- If the proportion of units in the sample is systematically different from the population, we may weight each unit:

$$\overline{y}_w = rac{\sum_{i=1}^{n_s} w_i y_i}{\sum_{i=1}^{n_s} w_i}$$

- The main issues are (again) bias and variability of this estimate $\ \overline{y}_w - \overline{Y}$

More on the Project Outline Handout

- Some Examples of Proposals (I.1, due Jan 25)
- Shall we look at the whole project outline as well?
 - This is your chance to ask questions about any parts of the handout that you read, and are concerned about.

20 January 2011 17	20 January 2011 18
Review	Review Measurement Representation
 Quality in Surveys 	Construct μ_i Validity Validity Coverage Error
 More on the Project Outline Handout 	Measurement V
 Reading: Up to today: responsible for Groves Ch's 1, 2, 3 Next week: Groves Ch 5 Groves, Ch 11 (sections 1-6) Groves Ch 4 (sections 1-3; we will do more later) in that order 	Measurement Error Processing Error Edited Response $y_{\bar{\ell}}$ Edited Response $y_{\bar{\ell}}$ Adjustment Error Postsurvey Adjustments \bar{y}_{rror}
 Lecture notes online at http://www.stat.cmu.edu/~brian/303 	Survey Statistic