# 36-303: Sampling, Surveys and Society

Review for Final Midterm, Apr 12 Brian W. Junker 132E Baker Hall brian@stat.cmu.edu

### Handouts & Online Stuff

### These Notes

Formula Sheet(s) for Final Midterm

- Posted in Week12 on class website
- Do not bring to exam; I will provide fresh copies Tuesday!
- HW05 Solutions: I will post them tomorrow (Fri) on class website (bug me if not!)
- HW05 Graded Papers: I will return these <u>after</u> the upcoming exam.

# Outline

### Review for Final Midterm Exam

- Tues Apr 12, 2011
- Closed book, closed notes
- Formula sheets (old one plus new one) provided
- Calculator recommended (please don't forget!!)
- Cumulative, but concentrating on
  - Groves Ch's 4, 6, 10
  - Class notes, readings from Weeks 7-12
  - HW 05
- This exam very similar in format to last one

# Review

- Good sampling and data collection
- Nonresponse
- Stratified Sampling
- Cluster Sampling
- Post-Survey Processing
- Imputation
- Post-stratification
  - Weights
  - Variance Estimation (Taylor and Jackknife)

### Good Sampling and Data Collection (1)

- Adjusting sample size for anticipated response rate
  - Email: 20% is typical
  - □ Phone: E.g. 2007 Pew Religious Survey had 25%
  - □ Face to Face: Over 70%; we saw 73%
- Collect demographic variables so you can post-stratify (to check, and if necessary, reweight sample to be "representative")
- SRS from C-book, list of faculty emails, etc.
  - Other methods if no frame, or SRS from frame is hard.

### Good Sampling and Data Collection (2)

#### Contacting respondents

- Once the sample (e.g. SRS) and mode of data collection is chosen (e.g. surveymonkey) is chosen, stick to it
  - You can break the SRS into "waves" and contact people in each wave separately; then if response rate is better than expected, later waves do not have to be contacted.
- But you can try to contact respondents in any reasonable way: email, phone, Facebook, etc., to improve response rates
- Followup with nonrespondents directly rather than send out general 2<sup>nd</sup> and 3<sup>rd</sup> notices to everyone in sample
- Late responders can be thought of as being like neverresponders.
- Distinguish refusers vs procrastinators: <u>"No" means "no"!</u>
- Personal, polite contacts work best

# Nonresponse (1)

- Types of non-response
  - Unit non-response
  - Item non-response

#### Reasons

- MCAR missing completely at random (ignorable msgnss)
- □ MAR missing at random
- MNAR missing not at random (informative missingness)
- What to do about it
  - Ignore it (MCAR!)
  - Prevent it
  - Impute missing responses (MCAR, MAR; hard for MNAR!)

# Nonresponse (2)

### Preventing Missingness

- Survey content
- Time of survey
- Interviewer skills
- Data collection method
- Questionnaire design
- Burden on respondent
- Survey Introduction
- Incentives
- Followup
- Imputing missing responses
  - More below on post-processing survey data...

# Stratified Sampling (1)

#### H strata

- $\square$  N<sub>h</sub> = population size in each stratum
- $n_h = \text{sample size in each stratum}$
- $f_h = n_h/N_h$  = sampling fraction, each stratum

$$\square W_{h} = N_{h}/N$$

Mean

$$\overline{y}_{st} = \sum_{h=1}^{H} W_h \overline{y}_h$$
 unbiased estim. of  $\overline{y}_{pop} = \sum_{i=1}^{N} y_i = \sum_{h=1}^{H} W_h \overline{y}_{h,pop}$ 

Variance

$$Var(\overline{y}_{st}) = \sum_{h=1}^{H} W_h^2 (1 - f_h) \frac{s_h^2}{n_h}$$

Η

 $N = \sum_{h=1} N_h$ 

H

 $n = \sum n_h$ 

h=1

# Stratified Sampling (2)

The <u>design effect</u> is a measure of how much better or worse <u>Stratified</u> is than <u>one SRS</u>:

$$DEFF = \frac{Var(\overline{y}_{st})}{Var(\overline{y}_{srs})} = \frac{\sum_{h=1}^{H} W_h^2 (1 - f_h) \frac{s_h^2}{n_h}}{(1 - f) \frac{s^2}{n}}$$

- Usually, DEFF < 1, i.e. stratified does better than one big SRS!
  - Usually best if:
    - Elements are more similar to each other within strata than between (e.g., substantively meaningful strata)
    - Proportionate sampling (*f<sub>h</sub>* same in every stratum)
  - Cochran (1961) suggests 2-6 strata usually give the best results; greater than 6 OK, but there are diminishing returns

# Stratified vs. Cluster Sampling





Take an SRS from every stratum:

Stratified Sampling



Variance of the estimate of  $\overline{y}_U$  depends on the variability of values within strata.

For greatest precision, individual elements within each stratum should have similar values, but stratum means should differ from each other as much as possible. Take an SRS of clusters; observe all elements within the clusters in the sample:



The cluster is the sampling unit; the more clusters we sample, the smaller the variance. The variance of the estimate of  $\overline{y}_U$  depends primarily on the variability *between* cluster means.

For greatest precision, individual elements within each cluster should be heterogeneous, and cluster means should be similar to one another.

# Cluster Sampling (1)

### One-stage clustering, equal cluster sizes:

For each cluster i in the SRS of clusters S, we can calculate the cluster mean

$$\overline{y}_i = \frac{1}{M} \sum_{j=1}^M y_{ij}$$

where M is the cluster size. Since S is an SRS of n clusters

$$\overline{y}_{cl} = \frac{1}{n} \sum_{i \in \mathcal{S}} \overline{y}_i$$

The standard error (SE) needed for constructing confidence intervals is the square root of

$$Var(\overline{y}_{cl}) \approx \left(1 - \frac{n}{N}\right) \frac{1}{n} s_{\overline{y}_i}^2 = \left(1 - \frac{n}{N}\right) \frac{1}{n} \left[\frac{1}{n-1} \sum_{i \in \mathcal{S}} (\overline{y}_i - \overline{y}_{cl})^2\right]$$

# Cluster Sampling (2)

As with stratified sampling we can calculate a *design effect* 

$$DEFF = \frac{Var(\bar{y}_{cl})}{Var(\bar{y}_{srs})} = \frac{Ms_{\bar{y}_i}^2}{s_{y_{ij}}^2} \approx 1 + (M-1)\rho ,$$

where  $\rho$  is the *intraclass correlation (ICC)*, to see what the effect on precision of clustering is.

- In stratified sampling we usually get DEFF < 1 if we design the strata to have very different means and little variation within stratum.
- In clustered sampling, we usually get DEFF > 1. We can make  $DEFF \approx 1$  by making the clusters have very similar means and lots of variation within cluster.

# Post-survey Processing



- <u>Top row:</u> Raw data collection process
  - The order of Coding, Data Entry and Editing will depend on the data collection design (FTF, phone, www, computer assisted, ...)
  - Computer-based surveys require you to design the Data Entry and Edit Checks when you build the form in surveymonkey.com, questionpro.com, etc.
- <u>Bottom row</u>: Calculations based on the data and/or design

# Imputation (1)

- <u>Weights</u> are a good solution for unit nonresponse (missed that whole person)
- <u>Imputation</u> is a good solution for item nonresponse (person never answered question #17).
- Basic ideas of imputation:
  - Build a model for <u>what sort of person wouldn't respond</u>, and use the model to fill in a value for this person
  - □ Find one or more other people like this person who <u>did</u> answer #17, and use their answers for this person

### Alternative to imputation: <u>Case-wise deletion</u>

- Delete this person from the survey so you don't have to deal with the nonresponse to question #17
- Pro's and con's of case-wise deletion??
- MCAR: Missing Completely at Random

# Imputation (2)

#### Mean Value Imputation

- If question #17 is a numerical item, take the average of everyone else's answer to #17, and fill that in for this person
- MCAR: Missing completely at random

#### Hot Deck Imputation

- Among all the other people who answered question #17, find the one person who matches this person on important variables: age, sex, occupation, answers to other questions, etc.
- Fill in that person's answer for this person's #17.
- MAR: Missing at Random (within covariates)

#### Regression Imputation

 Among all the people who answered question #17, fit a regression model (or logistic regression, or whatever) for response to question #17 as a function of other variables:

 $y_{17} = \beta_0 + \beta_1(age) + \beta_2(sex) + \beta_3(occupation) + \beta_4(answer to Q3) + ... + \epsilon$ 

- Use the fitted model to predict what this person would have answered to #17, and fill that value in
- MAR

### Post-Stratification (1)

- As part of survey data collection it is a good idea to get general demographic information (e.g. in our surveys: sex, age, class, major, hometown, etc.)
- After data collection we compare the proportions in each of these categories in our sample with the same proportions in the population
- If they agree, great. If they disagree, we may reweight the sample to make them agree

weight = (population proportion)/(sample proportion)

These categories are called "post-strata", and the weights are called "post-stratification weights"

# Post-Stratification (2)

- Post-stratification weights can fix
  - disproportionate sampling of post strata
  - disproportionate nonresponse across poststrata
- Only works if the sampling/nonresponse process is <u>ignorable</u> within post-strata
  - That is, nonresponse does not depend on the answer you would have gotten if the person had responded

#### If the sampling/nonresponse process is non-ignorable then these weights don't work; other weights have to be used

- The weights are only as good as your model for nonresponse
  - These weights are a very big deal in pre-election phone surveys for example (resp. rate as low as 5%, weights account for ignorable and nonignorable nonresponse)

### Example from HW05

Sex	College	$\mathrm{Hrs}/\mathrm{Wk}$	_	Sex	College	$\mathrm{Hrs}/\mathrm{Wk}$	
Μ	Eng	28	-	F	Eng	36	
$\mathbf{M}$	Eng	29		$\mathbf{F}$	Eng	33	
$\mathbf{M}$	Eng	23		Μ	$\operatorname{Lib}$	27	
$\mathbf{M}$	Eng	35		Μ	$\operatorname{Lib}$	28	
$\mathbf{M}$	Eng	29		$\mathbf{F}$	$\operatorname{Lib}$	29	
$\mathbf{M}$	Eng	30		$\mathbf{F}$	$\operatorname{Lib}$	30	
$\mathbf{M}$	Eng	34		$\mathbf{F}$	$\operatorname{Lib}$	28	
$\mathbf{M}$	Eng	31		$\mathbf{F}$	$\operatorname{Lib}$	28	
$\mathbf{F}$	Eng	30		$\mathbf{F}$	Lib	32	
$\mathbf{F}$	Eng	31		$\mathbf{F}$	Lib	30	

Sample Post-strata:				
$\mathbf{Sex}$	Eng	Lib		
Μ	8	2		
$\mathbf{F}$	4	6		

Population Post-strata:

$\mathbf{Sex}$	Eng	Lib
Μ	617	380
F	450	551

Post-strat. weights:

(617/1998)/(8/20) = 0.77	(380/1998)/(2/20) = 1.90
(450/1998)/(4/20) = 1.13	(551/1998)/(6/20) = 0.92

Sex	College	$\mathrm{Hrs}/\mathrm{Wk}$	Wgt		Sex	College	$\mathrm{Hrs}/\mathrm{Wk}$	Wgt
Μ	Eng	28	0.77	-	F	Eng	36	1.13
$\mathbf{M}$	Eng	29	0.77		$\mathbf{F}$	Eng	33	1.13
$\mathbf{M}$	Eng	23	0.77		$\mathbf{M}$	Lib	27	1.90
$\mathbf{M}$	Eng	35	0.77		$\mathbf{M}$	Lib	28	1.90
$\mathbf{M}$	Eng	29	0.77		$\mathbf{F}$	Lib	29	0.92
$\mathbf{M}$	Eng	30	0.77		$\mathbf{F}$	Lib	30	0.92
$\mathbf{M}$	Eng	34	0.77		$\mathbf{F}$	Lib	28	0.92
$\mathbf{M}$	Eng	31	0.77		$\mathbf{F}$	Lib	28	0.92
$\mathbf{F}$	Eng	30	1.13		$\mathbf{F}$	Lib	32	0.92
$\mathbf{F}$	Eng	31	1.13		$\mathbf{F}$	Lib	30	0.92

Unweighted mean:

$$\overline{y}_{srs} = \frac{1}{20} \sum_{i=1}^{20} y_i = 30.05$$

Weighted mean:

$$\overline{y}_w = \frac{\sum_i w_i y_i}{\sum_i w_i} = 29.91$$

### Taylor Series Variance Approximation

 $Var_{TS}(\overline{y}_{w}) \approx 0.46 =$  $\frac{1}{\left(\sum_{i} w_{i}\right)^{2}} \left| Var\left(\sum_{i} w_{i}y_{i}\right) - 2\overline{y}_{w}Cov\left(\sum_{i} w_{i}y_{i}, \sum_{i} w_{i}\right) + (\overline{y}_{w})^{2}Var\left(\sum_{i} w_{i}\right) \right|$ where  $\overline{y}_w = 29.91, \, \overline{w} = 1.00, \, \overline{wy} = 29.91$  and  $Var(\sum_{i=1}^{n} w_i) \approx n \cdot \frac{1}{n-1} \sum_{i=1}^{n} (w_i - \overline{w})^2 = 2.26$  $Var(\sum_{i=1}^{n} y_i w_i) \approx n \cdot \frac{1}{n-1} \sum_{i=1}^{n} (w_i y_i - \overline{wy})^2 = 1788.84$  $Cov(\sum_{i=1}^{n} y_i w_i, \sum_{i=1}^{n} w_i) \approx n \cdot \frac{1}{n-1} \sum_{i=1}^{n} (w_i y_i - \overline{wy})(w_i - \overline{w}) = 60.64$ 

### Jackknife Variance Approximation:

• Replicate 
$$\overline{y}_{w}^{(r)} = \frac{\sum_{i=1}^{n} w_{i}^{(r)} y_{i}^{(r)}}{\sum_{i=1}^{n} w_{i}^{(r)}}$$
 's:

29.99382 29.94970 30.21439 29.68501 29.94970 29.90558 29.72912 29.86147 30.09879 30.02371 29.64834 29.87356 30.00619 29.81600 29.93868 29.88352 29.99383 29.99383 29.77321 29.88352

#### Calculate

$$\overline{y}_{JK} = \frac{1}{n} \sum_{r=1}^{n} \overline{y}_{w}^{(r)} = 29.91$$

$$Var_{JK}(\overline{y}_{w}) = \frac{n-1}{n} \sum_{r=1}^{n} (\overline{y}_{w}^{(r)} - \overline{y}_{jk})^{2} = 0.34$$

 Now confidence intervals can be calculated in the usual way, e.g.

$$(\overline{y}_w - 2\sqrt{(1 - n/N)Var(\overline{y}_w)} \ , \ \ \overline{y}_w + 2\sqrt{(1 - n/N)Var(\overline{y}_w)})$$

for either the Taylor Series or Jackknife estimate of variance.

# Review

- Final Midterm Exam
  - Tues Apr 12, 2011, in class
  - Closed book, closed notes
  - □ Formula sheets (old one plus new one) provided
  - Calculator recommended
  - Cumulative, but concentrating on
    - Groves Ch's 4, 6, 10
    - Class notes, readings from Weeks 7-12
    - HW 05