36-303: Sampling, Surveys and Society

Stratified Samples and Sample Size Calculations
Brian Junker
132E Baker Hall
orian@stat.cmu.edu

06 March 2012

Outline

- Team Projects This Week
- Midterm
- Stratification
 - What is it; Notation
 - Weights and Proportionate Sampling
 - Variances and Design Effect
 - Examples

Handouts

- These Lecture Notes
- Handout on Stratified Sampling
- Handout on Sampling Details
 - Selecting an SRS from C-Book
 - Contacting respondents
 - Nonresponse followup on surveymonkey.com
- Reading:
 - Stratified Sampling: Groves Sect 4.5,
 - Nonresponse: Groves Ch 6

06 March 2012

Team Projects This Week

- II.5 Due Thursday (Blackboard)
 - □ Include a paragraph or so on your research question
 - Pretest a version of your questionnaire (or observation protocol) on a group of possible respondents/units.
 - Write 2-4 paragraphs: how many respondents/units you used in the pretest; how similar they were to units in the population; and the changes you made in the survey based on this pretest.
 - Include both old and revised questionnaire/protocol
- IRB Form Due Thursday (Blackboard)
 - If you are surveying people, and you have not turned one in to me yet, you need to do this by Thursday.
 - Form at http://www.stat.cmu.edu/~brian/303. You don't need to include any attachments
- Peer Evaluations Due Thursday (email to me)
 - Each person should email me forms for all other members of their team, in one email with subject "36303 Peer Evaluations"
 - Form at http://www.stat.cmu.edu/~brian/303

3

1

06 March 2012

2

Team Projects After Spring Break

- II.6 Project Plan (Tue Mar 20, Blackboard)
 - Final, full project proposal (items A-M on the "designing a sample survey" handout, except don't include the IRB form [item I]).
 - This should be easy: copy and update the latest thing you have done for each of the items A-M up to now into a single electronic file to submit on blackboard. for each team.
 - From this proposal, anyone outside our class should be able to read and understand completely what you are proposing to do

Get started Collecting Data!

06 March 2012

Stratified Sampling

- Strata are just subgroups of the target population that have some feature in common (gender, major, region, income, ...)
- Why stratify?
 - We need to make a separate inference for each stratum (e.g. we want to estimate men's and women's incomes separately)
 - Different sampling schemes would be used in each stratum (PA voters in PA, vs PA voters in Afganistan)
 - Population is geographically diverse (Minnesota, Illinois, Ohio, Pennsylvania)
 - Reduce variance of estimates (and reduce <u>sample size</u>) by exploiting similarity among members of the same stratum

Midterm Exam ...

Seemed to go well...

>	<pre>summary(exam1)</pre>				
	Min. 1st Qu.	Median	Mean 3rd Qu.	Max.	NA ' s
	64.00 76.00	85.00	82.41 89.00	99.00	3.00
>	<pre>stem(exam1)</pre>				
	6 4				
	6 8999		C (11)		
	7 2				
	7 567799				
	8 04		B (11)		
	8 566778899				
	9 224		A (6)		
	9 559				

06 March 2012

What is Stratification?

Record	Name	Group		Record	Name	Group	
1	Bradburn, N.	High		2	Cochran, W.	Highest	One Stratified Random
2	Cochran, W.	Highest		7	Hunt, J.	Highest	Sample of Total Size 4
3	Deming, W.	High		11	Madow, W.	Highest	v.
4	Fuller, W.	Medium	One SRS of Size 4	12	Mandela, N.	Highest	
5	Habermann, H.	Medium		19	Wolfe, T.	Highest	Wolfe, T.
6	Hansen, M.	Low		1	Bradburn, N.	High	Bradburn, N.
7	Hunt, J.	Highest		3	Deming, W.	High	
8	Hyde, H.	High		8	Hyde, H.	High	
9	Kalton, G.	Medium	Kalton, G.	17	Sudman, S.	High	
10	Kish, L.	Low		18	Wallman, K.	High	
11	Madow, W.	Highest		4	Fuller, W.	Medium	
12	Mandela, N.	Highest		5	Habermann, H.	Medium	
13	Norwood, J.	Medium	Norwood, J.	9	Kalton, G.	Medium	
14	Rubin, D.	Low		13	Norwood, J.	Medium	
15	Sheatsley, P.	Low		20	Woolsley, T.	Medium	
16	Steinberg, J.	Low		6	Hansen, M.	Low	
17	Sudman, S.	High		10	Kish, L.	Low	
18	Wallman, K.	High		14	Rubin, D.	Low	Rubin, D.
19	Wolfe, T.	Highest		15	Sheatsley, P.	Low	
20	Woolsley T	Medium		16	Steinberg, J.	Low	

7

5

Some Basic Notation

H strata
□ N_h = population size in each stratum $_H$ $N = \sum_{h=1}^{N_h} N_h$
\square n _h = sample size in each stratum $n = \sum n_h$
\Box f _h = n _h /N _h = sampling fraction, each stratum
The population average
$\overline{y}_{pop} = \frac{1}{N} \sum_{i=1}^{N} y_{i} = \frac{1}{N} \sum_{i=1}^{H} \sum_{j=1}^{N_{h}} y_{hi} = \sum_{i=1}^{H} \frac{N_{h}}{N} \frac{1}{N_{h}} \sum_{j=1}^{N_{h}} y_{hi} = \sum_{i=1}^{H} \frac{N_{h}}{N} \overline{y}_{h,pop}$
In stratified sampling we mimic this
$\overline{y}_{st} = \frac{1}{n} \sum_{i=1}^{n} y_i = \sum_{h=1}^{H} \frac{N_h}{N} \overline{y}_h \text{ where } \overline{y}_h = \frac{1}{n_h} \sum_{i=1}^{n_h} y_{hi}$
06 March 2012 9

Weights, and Proportionate Sampling

Let
$$W_h = N_h/N$$
. Then
 $\overline{y}_{pop} = \sum_{h=1}^{H} W_h \overline{y}_{h,pop} \text{ and } \overline{y}_{st} = \sum_{h=1}^{H} W_h \overline{y}_h$

- In proportionate sampling we let f_h = n_h/N_h = f for all strata h. Then n_h/n = N_h/N (why??)
 - □ The sample is called "self-weighting"
 - Sample mean is "simple" for self-weighting

$$\overline{y}_{st} = \sum_{h=1}^{N} W_h \overline{y}_h = \sum_{h=1}^{H} \frac{N_h}{N} \overline{y}_h = \sum_{h=1}^{H} \frac{n_h}{n} \overline{y}_h =$$

$$\sum_{h=1}^{H} \frac{n_h}{n} \frac{1}{n_h} \sum_{i=1}^{n_h} y_{hi} = \frac{1}{n} \sum_{h=1}^{H} \sum_{i=1}^{n_h} y_{hi} = \frac{1}{n} \sum_{i=1}^{n} y_i \overline{y}_{srs}$$

06 March 2012

Sampling Variances (SRS w/o replacement in each stratum)

Within each stratum it's the same old answer

$$Var(\overline{y}_h) = (1 - f_h) \frac{s_h^2}{n_h}$$
 where $s_h^2 = \frac{1}{n_h - 1} \sum_{i=1}^{n_h} (y_{hi} - \overline{y}_h)^2$

• Then we combine across strata using weights $(W_h)^2$: $Var(\overline{y}_{st}) = Var\left(\sum_{h=1}^{H} W_h \overline{y}_h\right)$ $= \sum_{h=1}^{H} Var(W_h \overline{y}_h) = \sum_{h=1}^{H} W_h^2 Var(\overline{y}_h)$

 $= \sum_{h=1}^{H} W_h^2 (1 - f_h) \frac{s_h^2}{n_h}$

Design Effect

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

$$d^{2} = \frac{Var(\overline{y}_{st})}{Var(\overline{y}_{srs})} = \frac{\sum_{h=1}^{H} W_{h}^{2}(1-f_{h}) \frac{s_{h}^{2}}{n_{h}}}{(1-f) \frac{s^{2}}{n}}$$

- Usually, d² < 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
 - Cochran (1961) suggests 2-6 strata usually give the best results; greater than 6 OK, but there are diminishing returns

06 March 2012

10

Handout on Stratified Sampling

(Briefly) Handout on Sampling Details

13 06 March 2012
Review
Team Projects This Week
Midterm Exam
Stratification
What is it; Notation
 Weights and Proportionate Sampling
 Variances and Design Effect
 Handout on Stratified Sampling
Handout on Sampling Details
Reading:
 Stratified Sampling: Groves Sect 4.5,
Nonresponse: Groves Ch 6
06 March 2012 15