

# STA304/1003 H1 F - Summer 2014: Surveys, Sampling, and Observational Data

## Lecture 10: Examples of Sampling Designs

Ramya Thinniyam

June 17, 2014

## Example 1: Lipstick Preference

A marketing research firm (MRF) estimates the proportion of customers preferring a certain brand of lipstick by “randomly” selecting 100 women who came by their booth in a shopping mall. Of the 100 sampled, 65 women stated a preference for brand A.

a) Identify the following for this study:

- ▶ Target Population : all women/customers who use lipstick
- ▶ Sampled Population : all women who came by MRF booth at the shopping mall.
- ▶ Variable : whether or not brand A is preferred or band preference
- ▶ Sampling Unit : a woman
- ▶ Observation Unit : a woman

b) What type of sampling design is it?

“Convenience sample” If women were truly randomly selected, could be considered an SRS with  $n=100$ ,  $N$  large, (but sampling frame is not available a prior !)

c). proportion of women who prefer

brand A = p

$$\hat{p} = \frac{65}{100} = 0.65$$

$$se(\hat{p}) = \sqrt{\frac{p(1-p)}{n}} = \sqrt{(1-\frac{n}{N})\frac{\hat{p}(1-\hat{p})}{n-1}}$$
$$= \sqrt{\frac{\hat{p}(1-\hat{p})}{n-1}}$$
$$= \sqrt{\frac{0.65(0.35)}{99}}$$
$$= 0.0479$$

Bound on error of est. based on 95% confidence is

$$1.96 se(\hat{p}) = 1.96(0.0479) = 0.0940$$

$$N \text{ large } fpc(1-\frac{n}{N}) \rightarrow 1$$
$$1 - \frac{n}{N} \approx 1$$

## Lipstick Preference - Questions (cont'd)

- c) Assuming all is correct with this sampling method, estimate the true proportion of women preferring brand A, and place a bound on the error of estimation.
- d) If a more accurate estimate has to be found, propose the minimal sample size required to estimate the true proportion within 3% of the true value. A reasonable guess is that the true proportion is within  $\pm 10\%$  of the value obtained in c).
- e) Did MRF select a simple random sample from the target population? Explain. Even if the sample may not be an SRS, could it still be reasonably representative of the target population?
- f) How would you help MRF to improve their sampling strategy? Revise the sampling design.

d).  $e=0.03$

$$0.55 \leq p \leq 0.75$$

$$N \text{ large } 1 - f \approx 1$$

Since no fpc required sample size

$$= n_0 \text{ maximize } S^2 = p(1-p) \text{ is maximized at } p=0.55$$

under constraint  $0.55 \leq p \leq 0.75$

$$n_0 = \left( \frac{Z + S^*}{e} \right)^2 = \frac{(1.96)(0.55)(0.45)}{(0.03)^2}$$
$$= 1056.44$$
$$= 1057$$

e). undercoverage : Not all women go to the mall.

overcoverage : Not all women go to both uses lipstick.

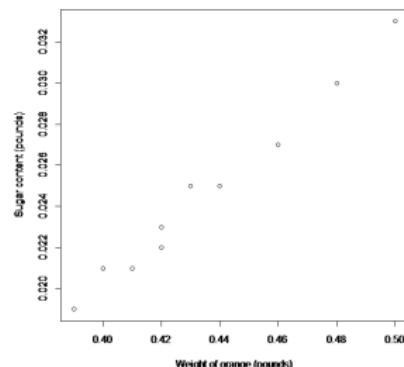
f). random select many mall and set many booth ask women (not only who comes to booth) anywhere in the mall

## Example 2: Sugar Content in Oranges

A company wishes to make inferences about the sugar content of its oranges based on a truckload of oranges that has just arrived. The total number of oranges is unknown and tedious to count. However, the total weight of all the oranges is determined to be 1800 pounds, by first weighing the truck loaded then unloaded. A random sample of 10 oranges is taken and for each orange the weight (in pounds) and sugar content (in pounds) is measured. Below is a summary of the sample data:

	Sugar Content (S)	Weight (W)	$(S - r * W)^2$
Total:	0.246	$y_i$	$4.35 x_i$

$$(res_i)^2 = (y_i - \hat{r}^* x_i)^2$$



$$n=10$$

N unknow

$$\hat{S}_i = r w_i$$

↑ slope estimate

$$S_i = R w_i + \varepsilon_i \quad r = \hat{R}$$

$$\bar{y} = \frac{0.246}{10} = 0.0246$$

$$\bar{w} = 0.435$$

$$S_r^2 = \frac{0.000053}{9}$$

$$1800 = T_w = T_x$$

a). SRS formulas cannot be used to estimate total sugar content since  $N$ , the total # of oranges is unknown.

b). Ratio Estimation:

① Sugar content and weight highly (+ve) correlated by common sense & scatterplot.

② Make sense to fit regression through origin.

③  $N$  is unknown but could be estimated/not even needed for ratio est.

$$c). \hat{N} = \frac{T_w}{w} = \frac{1800}{0.435} = 4137.9310$$

so approx 4138 oranges

$$d). \text{Estimate } R = \frac{\bar{s}_u}{w_u} = \frac{T_s}{T_w}$$

$$r = \frac{\bar{s}}{w} = \frac{\hat{T}_s}{\hat{T}_w} = \frac{0.296}{4.35} = 0.0566 \text{ pounds}$$

of sugar per pound of orange

## Sugar Content in Oranges - Questions:

- Explain why Simple Random Sampling formulas cannot be used to estimate the total sugar content in the truckload.
- What type of estimation would be reasonable to use for this scenario? Justify.
- Estimate the total number of oranges in the truckload.
- Give a point estimate for the mean amount of sugar content per pound of orange.
- The company has now been informed that the truckload contained 3000 oranges. Find a 95% CI for the total sugar content in the truckload.
- In the previous part, suppose you were to instead use SRS to find the 95% CI. Which CI (the one from this part or the previous part) would you expect to be wider. Justify.

e). now  $N=3000$  known  
 $\hat{T}_{sr} = \Gamma T_w = \frac{0.0246}{0.435} (1800)$   
 $= 101.88 ?$

$$\text{se}(\hat{T}_{sr}) = \sqrt{(1-\frac{n}{N}) \left( \frac{T_w}{w} \right)^2 \frac{S_{sr}^2}{n}}$$
$$= \sqrt{(1-\frac{10}{3000}) \left( \frac{1800}{0.435} \right)^2 \frac{0.000053}{9}} = 3.1701$$

95% CI for  $T_s$  (95.58, 108.18) pounds  
( $\frac{10}{3000}$ )

f). wider : SRS less precise

## Example 3: Escalators in Subways

A city transportation system includes 30 subway stations, each containing 6 escalators, and is interested in the number of days that the escalators were down for repair in the past year. 3 subway stations were randomly selected and the maintenance records for all 6 escalators in each are examined. Below are the results:

Station	Number of Days Escalator is Down	Average	Sample Variance
1	4, 3, 7, 2, 11, 0	4.5	15.5
2	11, 4, 3, 1, 0, 2	3.5	15.5
3	0, 3, 6, 4, 3, 2	3	4

Analysis of Variance Table

Response: days

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
station	2	7	3.500	0.3	0.7452
Residuals	15	175	11.667		

a) one-stage cluster sample psu=station psu=escalator?

$N=30$ ,  $n=3$ ,  $M_i=m=6$  (equal cluster sizes)

$$M = \sum_{i=1}^N M_i = Nm = 30 \times 6 = 180$$

$$\bar{M} = 6$$

$$b) \hat{Y}_{\text{umb}} = \frac{1}{M} \sum_{i \in S} M_i \bar{y}_i = \frac{1}{180} \cdot \frac{3}{3} \cdot 6(4.5 + 3.5 + 3)$$

equal cluster sizes

$$Se(\hat{Y}_{\text{umb}}) = \sqrt{\frac{1}{nM^2} \left(1 - \frac{1}{N}\right) \sum_{i \in S} (M_i \bar{y}_i - M \hat{Y}_{\text{umb}})^2}$$

$$= \sqrt{\frac{1}{3(6)^2} \left(1 - \frac{3}{30}\right) \sum_{i \in S} (\bar{y}_i - \hat{Y}_{\text{umb}})^2}$$

$$= \frac{1}{3(6)} \left(1 - \frac{3}{30}\right) (3.5) \xrightarrow{\text{MS Station}}$$

$$= 0.1708$$

CI skipped

$$c) \hat{T}_{\text{umb}} = M \times \hat{Y}_{\text{umb}} = 180 \times 3.6667 = 660.006 \text{ total down days}$$

d) using SRS: simply ignore the station,

$$Se(\bar{y})_{\text{SRS}} = \sqrt{(1 - \frac{n}{N}) \frac{s^2}{nm}} = \sqrt{(1 - \frac{1}{10}) \frac{SS_{TO}}{17(18)}}$$

$$= \sqrt{(1 - \frac{1}{10}) \frac{182}{17(18)}}$$

$$= 0.7316 \quad SS_{TO} = SS_{\text{Station}} + SS_{\text{Res}} = 7 + 175 = 182$$

Cluster better.

## Escalators in Subways - Questions:

- What type of sampling design is used here. Identify the characteristics and parameters.
- Estimate the mean downtime per escalator with a 95% CI.
- Estimate the total number of days down for all escalators in the subway.
- Compare the efficiency of this estimator with that of the SRS in this case. Decide which sampling design is better for this situation.
- Estimate the Intracluster Correlation Coefficient (ICC). Is the value of ICC in accordance to what you found in the previous part? Explain. If the SRS is more efficient than this sampling design, suggest another sampling design that would give better precision.

✓ by (d) expect ICC negative.

$$(e) SS_{\text{Res}} = \widehat{SSW} = \widehat{MSW} \times 30 \times (6-1) = 11.667 \times 150 = 1750.05$$

$$SS_{\text{Station}} = \widehat{SSB} = \widehat{MSB} \times (30-1) = 3.5 \times 29 = 101.5$$

$$ICC = 1 - \frac{m}{m-1} \frac{\widehat{SSB}}{\widehat{SS}_{TO}} = 1 - \frac{6}{5} \frac{1750.05}{1750.05 + 101.5} = -0.1342 < 0$$

as expected since var of cluster < var of SRS  $\Rightarrow$

Within each cluster it's heterogeneous, between clusters homogeneous

## Example 4: Refereed Publications by Department

A university has 807 faculty members of which 102 work in Biological Sciences, 310 in Physical Sciences, 217 in Social Sciences, and 178 in Humanities. The university is interested in estimating the number of refereed publications of its faculty members. Some faculty were randomly selected from each department for a total of 50 faculty members in the sample. The number of refereed publications of each selected member was carefully investigated and then recorded. Below is the frequency table for the number of refereed publications in each department:

Number of Refereed Publications	Faculty Members			
	Biological	Physical	Social	Humanities
0	1	10	9	8
1	2	2	0	2
2	0	0	1	0
3	1	1	0	1
4	0	2	2	0
5	2	1	0	0
6	0	1	1	0
7	1	0	0	0
8	0	2	0	0
Total	7	19	13	11

a). STRS by department

merits: ...

b). summary stats table:

	Bio	Phy	Soc	Hum
$\bar{Y}_i$ : Average	3.14	2.10	1.23	0.45
$S_i^2$ : S.var	6.81	8.21	4.36	0.87
$n_i$ : s.size	7	19	13	11

$L=4$  strata

$$N_1 = 102$$

$$N_2 = 310$$

$$N_3 = 217$$

$$N_4 = 178$$

$$\hat{T}_{\text{Bio}} = N_1 \cdot \bar{Y}_1 = 102 \cdot 3.14 = 320.57$$

$$\begin{aligned} \hat{V}\text{ar}(\hat{T}_{\text{Bio}}) &= N_1^2 \left(1 - \frac{1}{N_1}\right) \frac{S_i^2}{n_i} = (102)^2 \left(1 - \frac{1}{102}\right) \frac{6.81}{7} \\ &\approx 9426.327 \end{aligned}$$

Do for each department.

Combine 4, get total.

$$\hat{T}_{\text{str}} = 1321.189 \quad \sqrt{(\hat{T}_{\text{str}})} = 6560.78$$

## Refereed Publications by Department - Questions:

- What type of sampling design was used? Identify its parameters. Discuss the merits of using this type of sampling design and its applications.
- Estimate the total number of refereed publications by faculty members in this university, with a standard error.
- Give a 95% CI for the proportion of faculty members in this university with no refereed publications.
- Give a 95% CI for the percentage of faculty members in this university that have at least one refereed publication.
- Aside from the estimation done above, what other inferences (using any statistical methods) can be made. Name the methods and specify what types of inferences / questions of interest could be answered for each.

	$N_i$	$n_i$	$P_i$	$\frac{N_i}{n_i} \hat{P}_i$	$\frac{(1-\hat{P}_i)(N_i)^2}{n_i(n_i-1)}$	$\frac{(1-\hat{P}_i)(N_i)^2}{N_i(n_i-1)} \hat{P}_i(1-\hat{P}_i)$	$\frac{n_i}{n_i-1}$
Bio	102	7	1/7	0.018		0.0003	
Phy	310	9	10/9	0.202		0.0019	
Soc	217	13	9/13	0.186		0.0012	
Hum	178	11	8/11	0.160		0.0009	
Total	807	50		0.567		0.0043	

$$\hat{P}_{\text{str}} = 0.567$$

$$Se(P_{\text{str}}) = \sqrt{0.0043} = 0.0656$$

$$\begin{aligned} 95\% \text{ CI for } p: & 0.567 \pm 1.96(0.0656) \\ & = (0.4384, 0.6956) \end{aligned}$$

$$\begin{aligned} d). 95\% \text{ CI for } (1-p): \\ & (1 - 0.6956, 1 - 0.4384) \\ & = (30.44\%, 69.56\%) \end{aligned}$$

- do other sampling based on gender/ages/...other factors ...
- linear regression, do prediction.
- check if department & # of pbs are independent.