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

Lecture 4 - Part I: Simple Random Sampling (cont'd)

Ramya Thinniyam

May 27, 2014

### Sample Size Estimation

Determine sample size needed for survey based on your expectations:

1. Specify tolerable error, *e* (how close you want estimate to be to parameter):

P(
$$|\bar{y} - \bar{y}_{U}| \le e$$
) = 1 -  $\alpha$  can be other parameter as well, e called margin of error.

- 2. Determine N
- 3. Find an equation to relate tolerable error and sample size:

$$e = z_{\alpha/2} \sqrt{\left(1 - \frac{n}{N}\right) \frac{S}{\sqrt{N}}}$$

4. Estimate unknown quantities:

Be conservative: when in doubt, overestimate the variances — wider CI

- ▶ Find s from past research
- For proportion, use  $S^2 = \frac{1}{4}$  ie.  $p = \frac{1}{2}$  for large populations (if no other information is given)
- For proportion,  $Var(y_i) = p(1-p) = p p^2$  is a parabola so you can obtain its maximum if there are constraints on p
- 5. Solve for n
- 6. If calculated *n* is larger than you can afford, change expectations for the survey and retry

#### FORMULA FOR SAMPLE SIZE:

$$n = \frac{z_{\alpha/2}^2 S^2}{e^2 + \frac{z_{\alpha/2}^2 S^2}{N}} = \frac{n_0}{1 + \frac{n_0}{N}} \text{ ; where } n_0 = \left(\frac{z_{\alpha/2} S}{e}\right)^2$$

# **Examples: Sample Size Calculations**

Use benchmarks of 95% confidence when not specified.

1. 'mydata' Example: estimate mean within 0.5 of its true

value. Previous study yields a sample variance of 10. Estimate 
$$y_u$$
 with 0.5 of the true value.  $e=0.5$ ,  $\alpha=0.05$ ,

2. Estimate the proportion of students who own a computer in a population with 1000 students, with a margin of error of

0.03. We know at least 80% in the population own computers.

N=1000

Proportion 
$$P \ge 0.8$$
 $e = 0.03$ 
 $d = 0.05$ ,  $Z_{\frac{1}{2}} = 1.96$ 

The population own for  $S^{\frac{1}{2}} = (1-p)p$  under constraint  $P \ge 0.8$ 
 $S = 0.8$ 

The population own for  $S^{\frac{1}{2}} = (1-p)p$  under constraint  $P \ge 0.8$ 
 $S = 0.8$ 

The population own for  $S^{\frac{1}{2}} = (1-p)p$  under constraint  $P \ge 0.8$ 
 $S = 0.8$ 

The population own for  $S^{\frac{1}{2}} = (1-p)p$  under constraint  $P \ge 0.8$ 
 $S = 0.8$ 

The population own for  $S = (1-p)p$  under constraint  $P \ge 0.8$ 
 $S = 0.8$ 

The population own for  $S = (1-p)p$  under constraint  $P \ge 0.8$ 
 $S = 0.8$ 

The population own for  $S = (1-p)p$  under constraint  $P \ge 0.8$ 
 $S = 0.8$ 

The population own for  $S = (1-p)p$  under constraint  $P \ge 0.8$ 
 $S = 0.8$ 

The population own for  $S = (1-p)p$  under constraint  $P \ge 0.8$ 
 $S = 0.8$ 

The population own for  $S = (1-p)p$  under constraint  $P \ge 0.8$ 
 $S = 0.8$ 

The population own for  $S = (1-p)p$  under constraint  $P \ge 0.8$ 
 $S = 0.8$ 

The population own for  $S = (1-p)p$  under constraint  $P \ge 0.8$ 
 $S = 0.8$ 

The population own for  $S = (1-p)p$  under constraint  $P \ge 0.8$ 

The population own for  $S = (1-p)p$  under constraint  $S = (1-p)p$ 

The population own for  $S = (1-p)p$ 

The population of  $S = (1-p)p$ 

The

 $\int_{0}^{\infty} \left( \frac{Z_{\frac{1}{2}}S^{2}}{e} \right)^{2} = \frac{1.96 \cdot \sqrt{10}}{0.5}^{2}$  = |53.664|  $\int_{0}^{\infty} \frac{1 + \sqrt{10}}{1 + \sqrt{10}} = 60.5772$   $\int_{0}^{\infty} \frac{1 + \sqrt{10}}{1 + \sqrt{10}} = 60.5772$   $\int_{0}^{\infty} \frac{1 + \sqrt{10}}{1 + \sqrt{10}} = 60.5772$   $\int_{0}^{\infty} \frac{1 + \sqrt{10}}{1 + \sqrt{10}} = 60.5772$ 

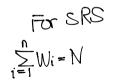
# Sampling Weights

$$\pi_i = P(unit \ i \ is \ in \ sample) : Inclusion Probability (not restricted to SRS)$$

For any sampling design, the sampling weight is the reciprocal of the inclusion probability:

$$w_i = \frac{1}{\pi_i}$$

- Interpreted as the number of population units represented by i different design has different weights
- For SRS:  $w_i = \frac{N}{n}$ : each unit represents itself +  $\frac{N}{n}$  1 unsampled units in the population
- ► For SRS: all weights are the same (each unit in the sample represents the same number of units in the population). called a self-weighting sample.



# Large Sample CIs for a Location - Finite Populations

\* Use Finite Population Correction (fpc) in variance estimates -  $\left| \left( 1 - \frac{n}{N} \right) \right|$  \*

#### For Mean:

$$ar{y} \pm z_{lpha/2} \sqrt{\left(1 - rac{n}{N}
ight) rac{S^2}{n}} \quad ext{ OR } \quad ar{y} \pm z_{lpha/2} \sqrt{\left(1 - rac{n}{N}
ight) rac{S^2}{n}}$$

#### For Total:

$$\hat{t} \pm z_{\alpha/2} \sqrt{\left(1 - \frac{n}{N}\right) \frac{N^2 S^2}{n}}$$
 OR  $\hat{t} \pm z_{\alpha/2} \sqrt{\left(1 - \frac{n}{N}\right) \frac{N^2 S^2}{n}}$ 

For Proportion:

$$\hat{p} \pm z_{\alpha/2} \sqrt{\left(1 - \frac{n}{N}\right) \frac{\hat{p}(1 - \hat{p})}{n - 1}}$$

with replacement.

For infinite populations or SRSWR, we don't have fpc. b/c in this case  $[-\frac{n}{N} \rightarrow 1]$  so all Stuffs above.

# Example: Local News Coverage

Let p=proportion of

 $A = \frac{\sum y_i}{y_i} = y_i = \frac{640}{1600} = 0.4$ 

news wronger.

A major metropolitan newspaper selected a SRS of 1,600 readers from their list Ui=local contenage; = \( \) , if i answered yes \( \) |= local contenage; = \( \) , o.w. \( \) |ength(local contenage) = length of vector = n of 100,000 subscribers. They asked whether the paper should increase its coverage of local news (1='yes', 0='no'). length (localcoverage) M = 100000[1] 1600 > mean(localcoverage) n = 16000.4 um (local coverage)  $\overline{y} = 0.4$ Sum (local coverage) =  $\sum_{i=1}^{n} y_i = \#$  of sampled subscribers who said yes. > sum(localcoverage) subscribers who wants more local [1] a) Find a 99% CI for the proportion of readers who would like more coverage of  $\hat{p} \pm Z_{02}\sqrt{(-\frac{n}{N})\hat{\rho}(-\frac{\hat{p}}{\hat{p}})} = 0.4 \pm 2.58\sqrt{(-\frac{1600}{60000})} \xrightarrow{0.4 \times 0.6} \Rightarrow (0.3686, 0.4314)$ local news.

- b) Find a 99% CI for the percent of readers who would like more coverage of 99% CI for 100p: (36.86,4314) local news.
- c) Find a 99% CI for the proportion of readers who would do not want more (1-1).4314,1-0.3686)=(0.5686,0.6314) local news coverage.

6/12

# Comparing Two Means - Infinite Populations

 $100(1-\alpha)\%$  Approximate CI for difference of two means:

$$(\hat{\mu}_1 - \hat{\mu}_2) \pm z_{\alpha/2} \sqrt{V(\hat{\mu}_1) + V(\hat{\mu}_2) - 2Cov(\hat{\mu}_1, \hat{\mu}_2)}$$

- ► If zero is in the interval, then there is no statistically significant difference between the two population means
- ▶  $100(1 \alpha)\%$  of all samples generate an interval that captures the true difference between the two means

Example: Lifetime of Lightbulbs — independent Infinite Pop. N. N. Infinite

assume 2 pops

Random sample of 140 traditional light bulb lifetimes:  $\bar{V}_1 = 1348.2, s_1 = 22.65$ ,  $n_i = 140$ 

Random sample of 80 new technology light bulb lifetimes:  $\bar{y}_2 = 1387.7, s_2 = 23.06, \eta_2 = 80$ 

Find a 95% CI for the difference between the two mean lifetimes. Can the manufacturer of the new technology light bulbs claim their mean lifetime is better than that for the traditional ones?

traditional ones?

The continuous states that the traditional ones?

The continuous for independence of independence of traditional is significantly lower = 
$$(1.343.2-1387.7)\pm 1.96\sqrt{\frac{(2.65)^3}{140}+\frac{(23.06)^3}{140}}=(-45.7939,-33.2061)$$

The continuous formula is significantly lower traditional in the continuous looks are continuous looks.

# Comparing Two Independent Locations - Finite Populations

Approximate large sample  $100(1 - \alpha)\%$  CI for difference between

#### Two Means:

$$(\bar{y}_1 - \bar{y}_2) \pm z_{\alpha/2} \sqrt{\left(1 - \frac{n_1}{N_1}\right) \frac{s_1^2}{n_1} + \left(1 - \frac{n_2}{N_2}\right) \frac{s_2^2}{n_2}}$$

#### Two Totals:

$$(\hat{t}_1 - \hat{t}_2) \pm z_{\alpha/2} \sqrt{\left(1 - \frac{n_1}{N_1}\right) \frac{N_1^2 S_1^2}{n_1} + \left(1 - \frac{n_2}{N_2}\right) \frac{N_2^2 S_2^2}{n_2}}$$

#### Two Proportions:

$$(\hat{p}_1 - \hat{p}_2) \pm z_{\alpha/2} \sqrt{\left(1 - \frac{n_1}{N_1}\right) \frac{\hat{p}_1(1 - \hat{p}_1)}{n_1 - 1} + \left(1 - \frac{n_2}{N_2}\right) \frac{\hat{p}_2(1 - \hat{p}_2)}{n_2 - 1}}$$

 If zero is in the interval, then there is no statistically significant difference between the two population locations

2 proportions, both and finite

(1). F  $N_1=352$ ,  $n_1=38$ ,  $\sum_{j=1}^{n_1} y_j = 22$ Example: Two Proportions -  $N_1$  and  $N_2$  finite

(2). M  $N_2=315$ ,  $n_2=105$ ,  $\sum_{j=1}^{n_2} y_j = 21$ A company wishes to increase the sales of their new points. A company wishes to increase the sales of their new product by

95% (I for (p1-p2):

advertising/targeting the correct consumers. They take a simple random sample in the male and female population and ask consumers who bought their product, "Do you like this product?"

Female Population:  $N_1 = 352$ ,  $n_1 = 88$ , 22 answered 'yes'

Male Population:  $N_2 = 315$ ,  $n_2 = 105$ , 21 answered 'yes'

Use a 95% CI to answer the question of interest. What would you recommend the company to do?

you recommend the company to do?

$$(\hat{p},-\hat{p}_2) \pm \frac{\sum_{s \in N(1-\frac{n_1}{N_1})} \hat{p}_1(1-\hat{p}_1)}{n_1-1} + (1-\frac{n_2}{N_2}) \frac{\hat{p}_2(1-\hat{p}_2)}{n_2-1}}{n_2-1}$$

$$0.25 \pm 1.96\sqrt{(1-\frac{38}{352})} \frac{0.25 \cdot 0.75}{87} + (1-\frac{105}{352}) \frac{0.25 \cdot 0.8}{104} = (-0.0507, 0.1507)$$
Therefore, no Significant difference between the product.

# Example: Two Means- $N_1$ finite, $N_2$ infinite a general type of population

infinite

A school has 500 children of which sample 25 are sampled: the mean number of pets per child is 1.32, s=0.3. In the general population, 25 children are sampled: the mean number of pets per child is 1.08, s=0.5. We wish to compare the mean number of pets in the school population with the general population.

a) What assumptions must be made to use a CI?

N=500, n=25, y=1.32, s=0.3 N=00, n=25, y=1.08, s=0.5

b) Find a 95% CI and make conclusions.

95% CI for (M1-M2): (1.32-1.08) 
$$\pm 1.96\sqrt{1-\frac{25}{300}} = \frac{0.5^2}{25} = (0.0129, 0.4671)$$
  
O not inside, so the first pop has more pets. significant difference.

11/12