

# Optimal sample size allocation for cluster sampling

Week 8 (5.4)

Stat 260, St. Clair

1 / 19

## Determining sample sizes for a cluster sample

**Problem:** You have a quantitative variable  $y$  and you want to estimate its population mean/total.

**Question 1:** How many SSU (elements) to sample?

**Question 2:** How many PSU (clusters) to sample?

**Optional:** How to do this in an optimal way?

2 / 19

## Optimal Allocation

This allocation is **optimal** because it either

- **minimizes costs** for a fixed SE/margin of error, *or*
- **minimizes SE/margin of error** for a fixed survey cost.

An optimal solution is "easy" to derive assuming equal cluster sizes:

- $M_i = M$ : cluster sizes are equal
- $m_i = m$ : cluster sample sizes are equal

3 / 19

## Optimal Allocation

**Mathematical Problem:**

Let  $c_1$  be the cost per PSU (cluster) and  $c_2$  be the cost per SSU (element). With  $c_0$  fixed costs, the total survey costs are

$$C(m, n) = c_0 + c_1 n + c_2(mn)$$

Variance is also a function of  $m$  and  $n$  and ANOVA MS.

$$V(\hat{y}_{unb}; m, n) = \left(1 - \frac{n}{N}\right) \frac{MSB}{nM} + \left(1 - \frac{m}{M}\right) \frac{MSW}{nm}$$

4 / 19

## Optimal Allocation: 1. SSU sample size

**Solution:** Use Lagrange Multiplier method to minimize one function (C or V) subject to the constraints of the other function.

The optimal SSU (element) sample size is

$$m_{opt} = \sqrt{\frac{c_1 M(N-1)(1-R_a^2)}{c_2(NM-1)R_a^2}} \approx \sqrt{\frac{c_1(1-R_a^2)}{c_2 R_a^2}} \text{ when } N \gg M$$

where  $R_a^2 = 1 - \frac{MSW}{S^2}$

5 / 19

## Optimal Allocation: 2. PSU sample size:

### (a) achieving a margin of error

**Problem:** How many PSU to sample to estimate  $\bar{y}_{\mathcal{U}}$  with  $(1 - \alpha)100\%$  confidence and a margin of error  $e = z_{\alpha/2} SE(\hat{\bar{y}}_{unb})$ ?

**Solution:** Get  $m_{opt}$ , if you ignore the FPC then

$$n_{opt} = \frac{\nu z_{\alpha/2}^2}{e^2} \text{ where } \nu = \frac{MSB}{M} + \left(1 - \frac{m_{opt}}{M}\right) \frac{MSW}{m_{opt}}$$

7 / 19

## Optimal Allocation

- We know  $m_{opt}$
- final sample size is then determined by  $n$ :

$$n \times m_{opt} = \text{number of observation units sampled}$$

**Question 2:** Determine  $n$  subject to a constraint:

- fixed SE/margin of error, *or*
- fixed survey cost.

6 / 19

## Optimal Allocation: 2. PSU sample size:

### (a) achieving a margin of error

- If  $N$  is smaller, don't ignore FPC and use:

$$n_{opt} = \frac{\nu z_{\alpha/2}^2}{e^2 + \frac{z^2 MSB}{NM}}$$

- To estimate  $t$  with  $e_t$  margin of error, set  $e = e_t / (NM)$ .

8 / 19

## Optimal Allocation: 2. PSU sample size:

### (b) Do not go over budget

**Problem:** How many PSU to sample if your budget is  $C$  dollars (or man hours, etc...)?

**Solution:** Get  $m_{opt}$ , then

$$n_{opt} = \frac{C - c_0}{c_1 + c_2 m_{opt}}$$

9 / 19

## Example: Dorms

- New GPA study: want to estimate average dorm GPA with a 95% ME of 0.2
  - $N = 100$  rooms with  $M = 4$  students per room
- Previous study: **One-stage example**
  - $msw = 0.18504$ ,  $msb = 0.56392$  and  $\hat{S}^2 = 0.279$
  - $\hat{R}_a^2 \approx 0.337$
- Costs?
- $c_1 = 20$  minutes to travel between rooms and
- $c_2 = 10$  minute to talk to each student.

11 / 19

## Optimal Allocation

- Note: The **cost** and **ME** solutions for  $n$  work for *any* values of  $m$  given a desired cost or ME.
- You need a guess at MSB and MSW
  - $MSB = S_t^2 / M$  (how to cluster totals vary?)
  - $MSW = \sum_i^N S_i^2 / N$  (within cluster variation?)
- HW #7: How to compute MSB and MSW from guesses of  $R_a^2$  and  $S^2$

10 / 19

## Example: Dorms

What is the optimal number of student to sample per room?

12 / 19

## Example: Dorms

How many rooms to sample to get a ME of  $e = 0.2$  for estimating mean GPA?

## Example: Dorms - check answer

- We used  $z = 1.96$  for 95% confidence, but we should be using a  $t$ -distribution with  $n - 1$  degrees of freedom for CI when  $n$  is "small"
- Check margin of error with our larger multiplier, suggests a larger  $n$

```
n <- 16
qt(.975, df= n-1)
## [1] 2.13145
se_squared <- (1-n/100)*0.56392/(n*4) + (1-2/4)*0.18504/(n*2)
qt(.975, df= n-1)*sqrt(se_squared) # 0.2 or less??
## [1] 0.2162418
```

13 / 19

14 / 19

## Example: Dorms - check answer

- Try  $n$  of 17, 18 and 19!

```
n <- c(17,18,19)
se_squared <- (1-n/100)*0.56392/(n*4) + (1-2/4)*0.18504/(n*2)
qt(.975, df= n-1)*sqrt(se_squared) # 0.2 or less??
## [1] 0.2077542 0.2000704 0.1930670
```

- Final answer:  $n = 19$  will give a ME of at most 0.2

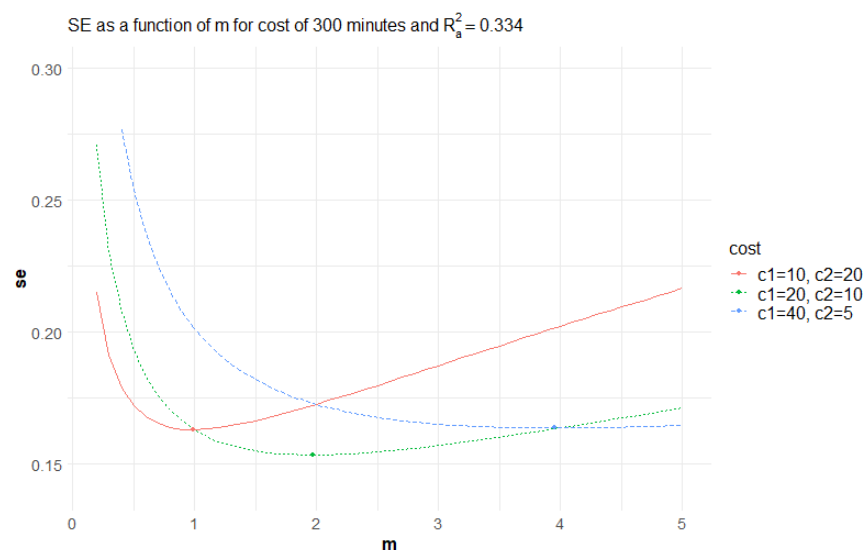
## Example: Dorms

How many rooms to sample if we have a fixed cost of 300 minutes?

15 / 19

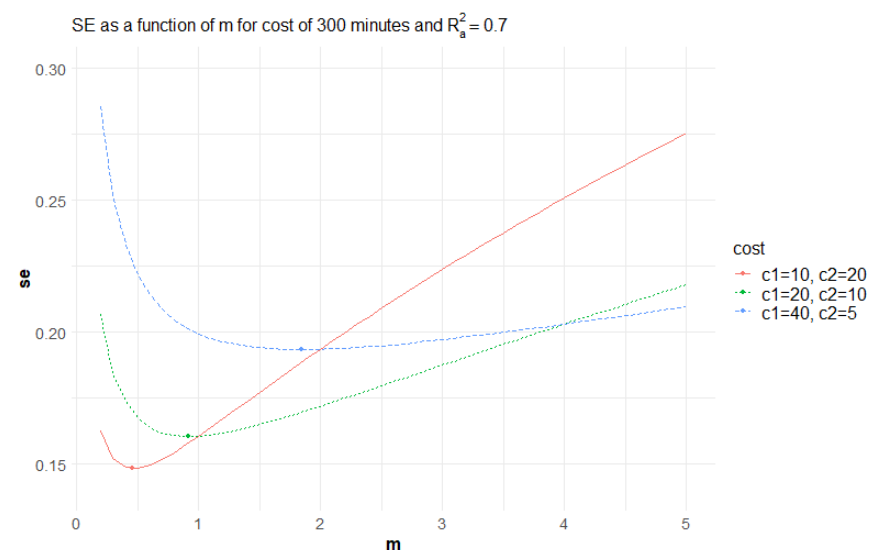
16 / 19

## Example: Dorms with $R_a^2 = 0.334$



17 / 19

## Example: Dorms with $R_a^2 = 0.7$



18 / 19

## Optimal Allocation: Unequal cluster sizes

If clusters are **not too variable with respect to size**, the (almost) optimal solution could use  $\bar{M}$  to get  $m_{opt}$

- use  $m_{opt}$  for all clusters or
- use an average of  $m_{opt}$ 
  - $m_i / M_i$  roughly constant

If clusters sizes are variable, don't use the optimal solution for equal sizes!

19 / 19