

# SURV625 Applied Sampling

2025-02-19

## SM 625: Week 4 Sampling Project Notes

For each of the three variables that will be the focus of the final course project, the Department of Education would like to generate estimates of means and proportions having a coefficient of variation of no more than 0.05. Using the numbers provided to you in the description of the final project, compute estimates of the element variances for each variable. Given these estimates, compute the desired level of precision (the desired sampling variance) for each estimate that corresponds to the desired coefficient of variation.

Now, given the desired levels of precision for each estimate, compute estimates of the necessary sample sizes for each of these three estimates (assuming simple random sampling), ignoring the finite population correction. These will be starting points for the eventual two-stage cluster sample design.

We first build a table to store our results for each week's assignments.

- We also add the expected averages for each outcome variable.

```
# build dataframe with inputs
MI_school_samples <- tibble(
  Outcome = c("smoked_cig", "smoked_mj", "age_approached_to_smoke"),
  type = c("prop", "prop", "mean"),
  desire_cv = rep(.05, 3),
  expect_mean = c(.25, .15, 12)
)
# calculate element
MI_school_samples |> kable()
```

Outcome	type	desire_cv	expect_mean
smoked_cig	prop	0.05	0.25
smoked_mj	prop	0.05	0.15
age_approached_to_smoke	mean	0.05	12.00

### Our process is to:

- 1st, calculate the estimated element variance.
  - For a proportion, to get the element variance we use  $\hat{p}(1 - \hat{p})$ .
  - For a mean, to get the element variance we simply just square the estimated standard deviation  $v(\bar{y}) = \sigma^2$ .
- 2nd, we calculate the estimated standard error as  $se(\hat{p}) = CV \times \hat{p}$ .
- 3rd, we compute the desired sampling variance as:  $var(\hat{p}) = se(\hat{p})^2$ , where  $se(\hat{p}) = \sqrt{var(\hat{p})}$

```
MI_school_samples <- MI_school_samples |>
mutate(
  # compute element variance
  s_sqrd = if_else(type=="prop", # for proportions
                  expect_mean * (1 - expect_mean),
                  if_else(type=="mean", # for means
                          1^2, NA)),
  # compute standard error
  se = desire_cv * expect_mean,
  # compute variance
  V = se^2
)

MI_school_samples |> select(-type) |> kable()
```

Outcome	desire_cv	expect_mean	s_sqrd	se	V
smoked_cig	0.05	0.25	0.1875	0.0125	0.0001563
smoked_mj	0.05	0.15	0.1275	0.0075	0.0000562
age_approached_to_smoke	0.05	12.00	1.0000	0.6000	0.3600000

We now estimate the desired sample sizes when we desire a  $CV = .05$  as  $n = \frac{s^2}{se^2}$

```
MI_school_samples <- MI_school_samples |>
  mutate(SRS_n = s_sqrd / V)

MI_school_samples |> select(1, SRS_n) |> kable()
```

Outcome	SRS_n
smoked_cig	1200.000000
smoked_mj	2266.666667
age_approached_to_smoke	2.777778

## SM 625: Week 5 Sampling Project Notes

For this week, we will consider the information available for stratified sampling of students. Eventually you are going to design a stratified cluster sample of students, where the clusters (or PSUs) are schools, but we aren't there yet.

Recall the regions of interest in the sampling project description:

```
school_frame <- read_xls(
  "~/repos/SURV625project/data/MI_school_frame_head_counts.xls")
```

Region	County_ID
1	07, 31, 66
2	22, 27, 36, 55
3	02, 21, 52
4	17, 48, 49, 77
5	01, 04, 06, 16, 20, 26, 35, 60, 65, 68, 69, 71, 72
6	05, 10, 15, 18, 24, 28, 40, 43, 45, 51, 53, 57, 67, 83
7	03, 08, 11, 12, 13, 14, 34, 39, 41, 54, 59, 61, 62, 64, 70, 75, 80
8	09, 19, 23, 25, 29, 30, 33, 37, 38, 46, 47, 56, 73, 78, 81
9	32, 44, 50, 58, 63, 74, 76, 79, 82

As “State officials are interested in providing, if at all possible, separate estimates for each of nine education regions in the state, where the regions are defined by groups of counties”, we will use these nine regions as strata.

Prepare a table that includes the:

- Overall population counts in each of these nine strata (the total count of students in the target population at each school is in the tot\_all column on the sampling frame).
- Given these counts, once you have the working overall sample size (unknown for now and will be decided by your team next week), what is the proportionate allocation plan of that sample of students across these nine strata?

```
# we will use Region, County_ID, and tot_all

# region counts
school_frame |>
  group_by(Region) |>
```

```

tally(tot_all) |>
mutate(prop_allocation = n/sum(n)) |>
rename(pop_count = n) |>
kable()

```

Region	pop_count	prop_allocation
1	3561	0.0042896
2	5474	0.0065941
3	8631	0.0103971
4	4855	0.0058484
5	18907	0.0227757
6	33133	0.0399126
7	191992	0.2312772
8	188830	0.2274682
9	374755	0.4514370

```

# what is the proportionate allocation plan of that sample
## of students across these nine strata?

```

## SM 625: Week 6 Sampling Project Notes

From a previous study, you obtain estimates of the following design effects for each of these three estimates:

- proportion ever smoked one cigarette = 2.5;
- proportion ever smoked marijuana = 2.0; and
- mean age when first asked to smoke = 1.7.

This previous study featured a sample of size  $n = 7,500$  students between the ages of 13 and 19, selected from a total of  $a = 150$  clusters. Using this information, compute a synthetic estimate of  $\rho_h$  for each of the three variables. These synthetic estimates of  $\rho_h$  will be used to consider alternative cluster sample designs as you continue with your project work. Finally, budget and cost information is now available. The total budget for data collection for this project will be \$500,000. The client and the data collection organization estimate that the data collection will cost \$3,000 per primary stage cluster (school), and \$50 per completed questionnaire within a cluster. We will use this cost information moving forward for optimal subsample size calculations.

We can estimate the sample ICC or  $\rho_h$  from the given design effect estimate as:

$$\hat{\rho}_h = \frac{deff - 1}{m - 1}$$

We now that the sample total is  $nm = 7500$  and the sample number of cluster is  $n = 150$ , which we can take the mean cluster size as  $m = nm/n = 7500/150 = 50$  and use it to calculate  $\rho_h$ .

```
nm <- 7500
n <- 150
m <- nm / n

MI_school_samples <- MI_school_samples |>
  # add deff and rho_h to our table
  mutate(deff = c(2.5, 2.0, 1.7),
         # compute rho_h
         rho_h = (deff - 1) / (m - 1))
```

```
)
```

```
MI_school_samples |> select(Outcome, deff, roh) |> kable()
```

Outcome	deff	roh
smoked_cig	2.5	0.0306122
smoked_mj	2.0	0.0204082
age_approached_to_smoke	1.7	0.0142857

## SM 625: Week 7 Sampling Project Notes

Recall that the client and the data collection organization estimated that the data collection would cost \$3,000 per primary stage cluster (school), and \$50 per completed questionnaire within a cluster. We will now use this information for optimum subsample size calculations. Recall that the total budget for data collection will be \$500,000.

Given this cost information and your estimates of  $\rho_{oh}$  for the three different variables of primary interest from last week, compute the optimum subsample size (and the corresponding optimal number of first stage clusters, given the total budget above) for each of the variables.

- How will you decide on a single overall optimum subsample size to use in your design?
- Think about a comparison of alternative cluster sample designs: under a fixed cost constraint, how would we decide which design would be best? What will be your overall sample size ( $n$ ) under this new optimum subsample size?

As you make progress in writing up what you have done so far, provide some discussion of the rationale for your choices in this regard.

Next, given this optimum subsample size and treating the values of  $\rho_{oh}$  as portable, compute the new expected DEFF for each estimate given the new design (this can be specific to each variable / estimate, given the different optimum subsample sizes). In addition, compute a new expected SRS variance for each variable under the new design, using the new “optimum” overall sample size (remember that you can treat the element variances for each variable estimated last week as portable). Finally, compute the new expected sampling variance for each estimate under this new cluster sample design. Are you still meeting the client’s precision requirements?

The client has also provided other new information: the estimated size of the target population is  $N = 830,138$ . Given this population size and your overall sample size ( $n$ ) under the new optimum subsample size computed above, what is your overall working sampling fraction ( $f$ )? Does it seem like finite population corrections will be necessary in your sampling variances if you choose to perform SRSWOR at some point?

The tables that you are developing and the text that accompanies them should carefully reflect the answers to all of the questions above.

We now have budget constraints and denote the cost per cluster as  $c_n = \$3,000$  and cost per element as  $c_m = \$50$ , with a total budget constraint of  $C = \$500,000$ . Since we know there are  $n = 150$  clusters and a total sample size of 7,500 students.

To compute the optimum  $m$  size we use the following equation:



$$m_{opt} = \sqrt{\frac{c_n}{c_m} \frac{1 - roh}{roh}}$$

Finally, we compute the sampling cost as  $n \times c_n + n \times m \times c_m$  which we defined these terms above.

```
c_n = 3000 # cost per cluster
c_m = 50 # cost per element within cluster
C = 500000 # total budget

MI_school_samples <- MI_school_samples |>
  mutate(
    # compute optimum m size
    m_opt = sqrt( (c_n / c_m) * ( (1-roh)/roh) ))

MI_school_samples |> select(Outcome, roh, m_opt) |>
  # we can print projected total cost for n=50
  #mutate(projected_cost = (c_n * n) + (c_m * n * m_opt),
    #projected_cost = scales::dollar(projected_cost)) |>
  kable()
```

Outcome	roh	m_opt
smoked_cig	0.0306122	43.58899
smoked_mj	0.0204082	53.66563
age_approached_to_smoke	0.0142857	64.34283

We now use the optimum m size and our budget constraints to calculate the an optimal number of  $n$  clusters for each outcome variable. For each outcome we use the optimum  $m$  to find the optimum  $n$  within our budget constraints as  $\$500,000 = n(c_n + m \times c_m)$ .

- We also compute the total sample size for these three options, and from here we can compute the total cost as  $n \times c_n + n \times m \times c_m$

```
MI_school_samples <- MI_school_samples |>
  mutate(
    # compute optimum n
    n_opt = C / (c_n + m_opt * c_m),
    # compute total SSU
```

```

    total_nm = m_opt * n_opt
  )

MI_school_samples |>
  select(Outcome, roh, m_opt, n_opt, total_nm) |>
  # we can print projected total cost for n=50
  mutate(projected_cost = (c_n * n_opt) + (c_m * n_opt * m_opt),
          projected_cost = scales::dollar(projected_cost)) |>
  kable()

```

Outcome	roh	m_opt	n_opt	total_nm	projected_cost
smoked_cig	0.0306122	43.58899	96.53536	4207.879	\$500,000
smoked_mj	0.0204082	53.66563	87.97734	4721.360	\$500,000
age_approached_to_smoke	0.0142857	64.34283	80.42281	5174.631	\$500,000

Note. if we round up for the optimum  $m$  and  $n$  we would slightly be above the total budget.

To be continued . . . . Answer . . .

How will you decide on a single overall optimum subsample size to use in your design?

- We will need to use the optimum  $n$  and  $m$  of each row and project them to the other rows to calculate the design effect:  $\text{deff} = (\text{SAM\_variance} / \text{SRS\_variance})$ . This should give us an idea of the mean deff for three total projections.