

Sampling I

LPO 9951 / Fall 2015

PURPOSE In most stats classes, all samples are assumed to be simple random samples from the population, with each unit having exactly the same probability of being selected. In practice, this is extremely rare. Samples are usually designed with unequal probabilities of selection across different groups. Because survey methodology is complex, this lecture cannot pretend to be comprehensive. Instead, it is meant to expose you to various sampling designs often found in education research as well as the formulas for computing means and variances of some of the simpler designs.

Simple random sampling (SRS)

Formulas

Where y_i is value of y for the i th unit:

sample mean

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i \quad (1)$$

sample variance

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$$

standard error of sample mean

$$\bar{y}_{se} = \sqrt{\frac{s^2}{n}} \quad (3)$$

Compute population means and variances

First, let's compare our hand computations of the population mean and variance with Stata's output using `summarize` and `mean`. To do this, we'll read in some fake SAT-like score data from a population of 1.5 million test takers.

```
. use ${datadir}fakesat, clear

. // calculate population mean, variance, sd, and sem by hand
. egen scoretot = total(score)           // total of all scores

. scalar popmean = scoretot / _N         // population mean score

. gen sqdiff = (score - popmean)^2       // (Xi - Xbar)^2

. egen sst = total(sqdiff)              // total of squared differences

. scalar popvar = sst / (_N - 1)         // population variance (not super pop)
```

```

. scalar popsd = sqrt(popvar)           // population standard deviation
. scalar popsem = popsd / sqrt(_N)      // standard error of population mean

. // compare to stata internal calculation
. scalar list popmean popsd popsem
    popmean = 499.87537
    popsd = 99.859033
    popsem = .08153456

. summarize score

```

Variable	Obs	Mean	Std. Dev.	Min	Max
score	1,500,000	499.8754	99.85903	200	800

```

. mean score

```

Mean estimation Number of obs = 1,500,000

	Mean	Std. Err.	[95% Conf. Interval]
score	499.8754	.0815346	499.7155 500.0352

NB: Even though we're talking about these fake test-takers as the population, they can still be thought of as a sample from a superpopulation of potential test takers. This is why I still use $n - 1$ correction, which is what Stata also uses. With such large N , however, the adjustment is functionally moot.

Compute SRS mean and variances

Now we'll take a simple random sample (SRS) of 10% of our test takers and compute our statistics.

```

. keep score

. // randomly sample 10%, all with equal probability of selection
. sample 10
(1,350,000 observations deleted)

. // calculate sample mean, variance, sd, and sem by hand
. egen scoretot = total(score)           // total of all scores

. scalar sampmean = scoretot / _N        // sample mean score

. gen sqdiff = (score - sampmean)^2      // (Xi - Xbar)^2

. egen sst = total(sqdiff)               // total of squared differences

. scalar sampvar = sst / (_N - 1)        // sample variance

. scalar sampsd = sqrt(sampvar)          // sample standard deviation

```

```
. scalar sampsem = sampsd / sqrt(_N)          // standard error of sample mean

. // compare to stata internal calculation of sample
. scalar list sampmean sampsd sampsem
    sampmean = 499.98176
    sampsd = 99.790632
    sampsem = .2576583
```

```
. summarize score
```

Variable	Obs	Mean	Std. Dev.	Min	Max
score	150,000	499.9817	99.79063	200	800

```
. mean score
```

Mean estimation Number of obs = 150,000

	Mean	Std. Err.	[95% Conf. Interval]
score	499.9817	.2576583	499.4767 500.4867

As you can see, our sample mean, variance, and standard error of the mean are about the same as the population values. \bar{y}_{se} is a little higher, which is to be expected since we are basing our estimate off fewer observations. And in both cases, our hand calculations are the same as those given by Stata. That is always a good sign!

SRS with a finite population correction (FPC)

Description and formula

Consider the normal estimate of the standard error of the mean, shown in equation 3 above. In cases where the proportion of the population that is sampled is quite large, this will in fact be an overestimate of the standard error of the mean. This is because in classical statistical theory, the sample is conceived as being from an infinitely large population. The finite population correction is a way of adjusting for the fact that the sample actually may be more representative than standard approaches would suggest. The finite population correction (FPC) is calculated as:

$$fpc = \sqrt{\frac{N-n}{N-1}} \quad (4)$$

where N is the population size and n is the sample size. As you can see, as n grows small relative to N , the FPC will approach 1 and the correction will be very small. As n becomes a larger fraction of N , the opposite is true. The FPC is rarely used in practice, but it should be used whenever the population size is known. Calculating \bar{y}_{se} using the FPC is done as follows:

$$\bar{y}_{se} = \sqrt{\frac{s^2}{n}} \times (fpc) = \sqrt{\frac{s^2}{n}} \sqrt{\frac{N-n}{N-1}} \quad (5)$$

Example

We'll again use some fake data to test our formulas. This time we have test score data for 50 students from a single large class. Let's say, for some mysterious reason, we only have access to information from 30 students. Maybe we did an exit poll of grades after class and assume that the 30 responses represent a random sample (highly unlikely, but we'll go with it for now). This number of students represents a sizeable portion of the population so we should adjust our estimate of the error the average score to take that into account.

```
. use ${datadir}singleclasstest, clear

. // check population stats; store population number
. sum score
```

Variable	Obs	Mean	Std. Dev.	Min	Max
score	50	80.02087	5.986087	64.91638	93.89526

```
. scalar N = _N

. // sample large part of population
. sample 30, count
(20 observations deleted)

. scalar n = _N

. // calculate sample stats without fpc correction
. egen scoretot = total(score)           // total of all scores

. scalar xbar = scoretot / n             // sample mean score

. gen sqdiff = (score - xbar)^2           // (Xi - Xbar)^2

. egen sst = total(sqdiff)               // total of squared differences

. scalar var_x = sst / n                 // sample variance

. scalar sd_x = sqrt(var_x)              // sample standard deviation

. scalar sampsem = sd_x / sqrt(n)        // standard error of sample mean

. scalar list xbar var_x sd_x sampsem
      xbar =  78.874406
      var_x = 26.844847
      sd_x =  5.1812013
      sampsem = .9459536

. // calculate fpc
. scalar fpc = sqrt((N - n) / (N - 1))

. // correct sem with fpc
. scalar sampsemfpc = sampsem * fpc

. scalar list sampsem sampsemfpc
      sampsem = .9459536
```

```

sampsemfpc = .60434759

. // 95% CI for the mean
. di xbar - invnormal(.975) * sampsem
77.020371

. di xbar + invnormal(.975) * sampsem
80.728441

. // 95% CI for fpc-adjusted mean
. di xbar - invnormal(.975) * sampsemfpc
77.689906

. di xbar + invnormal(.975) * sampsemfpc
80.058905

```

Comparing the 95% confidence intervals, we can see they are a little tighter when the FPC is used. This is a reflection of our knowledge that our sample represents a sizeable portion of the population and therefore is a better estimate than the standard formula will compute.

SRS and frequency weights

Sometimes data are reported in what's known as a frequency-weighted design. In such a setup, observations that take on the same values are reported only once, with a weight that is equal to how many times this particular set of observations was reported. This was a very common way of formatting data when computer memory was expensive, but less common now. You may still run across it from time to time so it's good to know about.

To demonstrate, we'll again use the fake SAT data. This time, however, the data are in a frequency format. Using the weight option for the `mean` command, we can set `[fw = freq]` and get the same estimates for mean score as we did with the full dataset.

```

. use ${datadir}fakesat_freq, clear

. // list first few observations
. list if _n < 11

```

```

+-----+
| score  freq |
+-----+
1. |   200   2412 |
2. |   210    908 |
3. |   220   1213 |
4. |   230   1539 |
5. |   240   2045 |
+-----+
6. |   250   2725 |
7. |   260   3313 |
8. |   270   4305 |
9. |   280   5393 |
10. |   290   6621 |
+-----+

```

```
. // compare simple mean with frequency-weighted mean
. mean score
```

```
Mean estimation                Number of obs   =           61
```

```
-----+-----
          |      Mean   Std. Err.   [95% Conf. Interval]
-----+-----
    score |      500    22.7303    454.5326    545.4674
-----+-----
```

```
. mean score [fw = freq]
```

```
Mean estimation                Number of obs   =   1,500,000
```

```
-----+-----
          |      Mean   Std. Err.   [95% Conf. Interval]
-----+-----
    score |  499.8754   .0815346    499.7155    500.0352
-----+-----
```

SRS with inverse probability weights

For most survey-based social science data, it is unlikely that all population members have the same probability of being sampled. This is a problem when the probability of selection is correlated with the quantities we hope to estimate. Without accounting for the probability of selection, our estimates will be biased, perhaps severely.

Going back to our fake SAT data, let's assume that our sample data come from voluntary responses and that test takers are more likely to report their scores if those scores are high (a not unreasonable situation). For purposes of the example, let's assume that we know the probability that a test taker will select to report his or her results (something that we are generally *very* unlikely to know). We generate our 1% sample this time based on the probability of reporting and check the unadjusted sample mean.

```
. use ${datadir}fakesat, clear

. // assume probability of reporting score is corrected with score
. gen preport = score / 1000 + .1 * (score / 10000)^2 + rnormal(0, .025)

. // sample based on probability of reporting
. gsample 1 [w = preport], percent
(81 observations created)
(1,485,081 observations deleted)

. // mean of sample
. mean score
```

```
Mean estimation                Number of obs   =       15,000
```

```
-----+-----
          |      Mean   Std. Err.   [95% Conf. Interval]
-----+-----
    score |  521.0153   .8015293    519.4442    522.5864
-----+-----
```

As expected the mean score is much higher than the population average (which should be around 500). Since we are omniscient researchers, we can generate inverse probability weights using the probability of reporting. Thinking it through, these weights will downweight those with a high likelihood of reporting and upweight those with a low likelihood, hopefully improving our estimate of the population mean score in the process.

```
. gen pweight = 1 / preport
```

```
. // check probability-weighted mean
. mean score [pweight = pweight]
```

Mean estimation Number of obs = 15,000

	Mean	Std. Err.	[95% Conf. Interval]	
score	500.7936	.8872305	499.0545	502.5327

Stratified random sampling with probability proportional to size

Description

Stratified sampling is a widely used and broadly applicable way of designing a sample. In stratified sampling, a set of strata are selected from the population, then samples are taken from within each strata. An example would be taking a sample of students from within elementary, junior, and high schools, with level of school as the strata. The idea is that strata are different in some fundamental way from each other but internally similar. Strata should effectively partition the population space, that is, not overlap and fully account for the population when put together.

Formulas

The notation for this type of sampling design is as follows:

stratum mean

$$\bar{y}_h = \frac{1}{N_h} \sum_{j=1}^{N_h} y_{hj} \quad (6)$$

stratum variance

$$s_h^2 = \frac{1}{N_h - 1} \sum_{j=1}^{N_h} (y_{hj} - \bar{y}_h)^2 \quad (7)$$

population mean

$$\bar{y} = \frac{1}{N} \sum_{h=1}^L N_h \bar{y}_h \quad (8)$$

population mean variance

$$s^2 = \sum_{h=1}^L \left(\frac{N_h}{N} \right)^2 \left(\frac{N_h - n_h}{N_h - 1} \right) \left(\frac{s_h^2}{n_h} \right) \quad (9)$$

where

- N is the population total
- y_{hj} is observation j within stratum h
- \bar{y}_h is the mean within stratum h
- N_h is the total number within stratum h
- n_h is the number sampled within stratum h
- s_h^2 is the variance within stratum h

Example

This time we'll using fake data on a high school with grades 9-12. A student in this school is considered to be *at risk* if his or her test score falls below a certain cut off, commensurate with the student's grade. Looking at the administrative data we can see the proportion at risk, within each grade and across the school.

```
. // read in fake highschool data; store full student population
. insheet using ${datadir}fakehs.csv, clear
(8 vars, 2,003 obs)
```

```
. scalar stupop = _N
```

```
. // proportion of students at risk
. mean atrisk                                // overall
```

```
Mean estimation                                Number of obs    =        2,003
```

```
-----+-----
          |      Mean   Std. Err.   [95% Conf. Interval]
-----+-----
    atrisk |    .337993   .0105719   .3172599   .3587261
-----+-----
```

```
. mean atrisk, over(grade)                    // within each grade
```

```
Mean estimation                                Number of obs    =        2,003
```

```
      9: grade = 9
     10: grade = 10
     11: grade = 11
     12: grade = 12
```

```
-----+-----
      Over |      Mean   Std. Err.   [95% Conf. Interval]
-----+-----
```


atrisk					
9		.3252336	.0202723	.2854766	.3649907
10		.3089431	.0208524	.2680485	.3498377
11		.3677932	.0215218	.3255857	.4100008
12		.3509514	.021968	.3078688	.394034

Of course, we don't really know these values. Why bother sampling if we did? But let's say that we do have a (correct) suspicion that the proportion of *at risk* students is different across grades. We therefore want to make sure that we first stratify on grade, then randomly sample, giving us an equal number of sampled students within each grade.

```
. global ss = 50                                // set within grade sample size

. sample $ss, count by(grade)                  // sample
(1,803 observations deleted)
```

To compute the mean within each stratum, we'll just take the simple mean of the observations within that stratum. We could weight each observation by the inverse of the selection, in this case, the number of students in the grade over the number selected and use equation (6), but the number of students in the grade would cancel out and we would be left with the simple sample average.

To compute the within stratum variance for each stratum, we'll use part of equation (9)

$$s_h^2 = \left(\frac{N_h - n_h}{N_h - 1} \right) \left(\frac{s_h^2}{n_h} \right)$$

In this simple stratification design, the population variance is just the weighted sum of the individual strata variances. So to get an estimate of the population variance, we'll weight each one by

$$\left(\frac{N_h}{N} \right)^2$$

and add them together.

To compute these values, we will use Stata's **preserve** and **restore** commands along with **collapse**, which allows us to compute the various means and standard deviations that we need.

```
. preserve

. collapse (mean) propatr = atrisk (sd) sdatr = atrisk (first) nstgrade, by(grade)

. scalar Ybar9 = propatr[1]                      // 9th grade average

. scalar Ybar10 = propatr[2]                     // 10th grade average

. scalar Ybar11 = propatr[3]                     // 11th grade average

. scalar Ybar12 = propatr[4]                     // 12th grade average

. // compute ((N_h - n_h) / N_h - 1) * (s_h^2 / n_h)
. gen varatr = ((nstgrade - $ss) / (nstgrade - 1)) * (sdatr^2 / $ss)
```

```

. // compute sem for each grade and store
. gen grade_sem = sqrt(varatr)

. scalar Ybar9_sem = grade_sem[1]           // 9th grade sem
. scalar Ybar10_sem = grade_sem[2]          // 10th grade sem
. scalar Ybar11_sem = grade_sem[3]          // 11th grade sem
. scalar Ybar12_sem = grade_sem[4]          // 12th grade sem

. gen weight = nstgrade / stupop           // (N_h / N)
. gen wproptr = weight * proptr            // weight strata proportions
. gen wvaratr = weight^2 * varatr           // weight strata variances
. collapse (sum) wproptr wvaratr            // sum within stata means and vars
. scalar Ybar = wproptr[1]                  // store estimate of pop. at risk
. scalar Ybar_sem = sqrt(wvaratr[1])        // compute root of above measure

. restore

```

Let's compare our estimates to the simple means computed by Stata.

```

. scalar list Ybar Ybar_sem
      Ybar = .33638542
      Ybar_sem = .03199795

. mean atrisk

```

Mean estimation Number of obs = 200

```

-----
      |           Mean   Std. Err.   [95% Conf. Interval]
-----+-----
      atrisk |           .335   .0334585   .2690213   .4009787
-----

```

```

. scalar list Ybar9 Ybar9_sem Ybar10 Ybar10_sem Ybar11 Ybar11_sem Ybar12 Ybar12_sem
      Ybar9 = .36000001
      Ybar9_sem = .06534968
      Ybar10 = .36000001
      Ybar10_sem = .06505993
      Ybar11 = .36000001
      Ybar11_sem = .06513891
      Ybar12 = .25999999
      Ybar12_sem = .05932035

. mean atrisk, over(grade)

```

Mean estimation Number of obs = 200

```

9: grade = 9
10: grade = 10
11: grade = 11
12: grade = 12

```

	Over	Mean	Std. Err.	[95% Conf. Interval]	
atrisk					
	9	.36	.0685714	.2247801	.4952199
	10	.36	.0685714	.2247801	.4952199
	11	.36	.0685714	.2247801	.4952199
	12	.26	.062662	.1364332	.3835668

As expected, the means are all almost exactly the same. Our various standard error estimates, however, are a little smaller. We didn't improve them much in this case, but we were able to make them smaller since we were able to incorporate our knowledge about the sample design.

Stratified cluster sampling with probability proportional to size

Description

Cluster sampling involves taking a sample where a group of clusters (each of which contains multiple units) is designated, then a random sample of these clusters are drawn. Within each cluster, all units may be included or, in large-scale surveys, a second sample may be drawn. In either case, this last unit is typically the unit of analysis. For example, if we decided to conduct a study by taking a random sample of classrooms within a school, then taking a sample of students within those classrooms, this would be a cluster sampling design. In this example, the classrooms would be the primary sampling unit (*psu*) and the students would be the secondary sampling unit (*ssu*).

Formulas

population size

$$M = \sum_{h=1}^L \sum_{i=1}^{N_h} M_{hi} \quad (10)$$

population total

$$Y = \sum_{h=1}^L \sum_{i=1}^{N_h} \sum_{j=1}^{M_{hi}} Y_{hij} \quad (11)$$

sample total

$$m = \sum_{h=1}^L \sum_{i=1}^{n_h} m_{hi} \quad (12)$$

estimated population total

$$\hat{Y} = \sum_{h=1}^L \sum_{i=1}^{n_h} \sum_{j=1}^{m_{hi}} w_{hij} Y_{hij} \quad (13)$$

where

$$w_{hij} = \frac{N_h}{n_h}$$

estimated population size

$$\hat{M} = \sum_{h=1}^L \sum_{i=1}^{n_h} \sum_{j=1}^{m_{hi}} w_{hij} \quad (14)$$

estimated population total variance

$$\hat{V}(\hat{Y}) = \sum_{h=1}^L (1 - f_h) \frac{n_h}{n_h - 1} \sum_{i=1}^{n_h} (y_{hi} - \bar{y}_h)^2 \quad (15)$$

where

$$y_{hi} = \sum_{j=1}^{m_{hi}} w_{hij} y_{hij}$$

$$\bar{y}_h = \frac{1}{n_h} \sum_{i=1}^{n_h} y_{hi}$$

and

$$f_h = \frac{n_h}{N_h}$$

Example

Using the fake highschool data, let's try to get an estimate of test scores. First, let's take a look at the population values (which, again, we normally don't know):

```
. insheet using ${datadir}fakehs.csv, clear
(8 vars, 2,003 obs)

. // get population estimates
. mean testscore                                // overall

Mean estimation                                Number of obs   =        2,003

-----+-----
          |      Mean   Std. Err.   [95% Conf. Interval]
-----+-----
 testscore |    509.0789   1.210999    506.7039    511.4538
-----+-----

. mean testscore, over(grade)                    // within grade

Mean estimation                                Number of obs   =        2,003
```

```

9: grade = 9
10: grade = 10
11: grade = 11
12: grade = 12

```

	Over	Mean	Std. Err.	[95% Conf. Interval]	
testscore					
	9	480.3551	2.132872	476.1723	484.538
	10	503.9045	2.26261	499.4672	508.3418
	11	515.5706	2.193136	511.2695	519.8716
	12	540.0465	2.318895	535.4988	544.5942

Estimated means This time, rather than simply sampling students within each grade, let's sample entire classes. These will be our *PSUs* with students being the *SSUs*. After taking only 10 classes in each grade, we'll compute the mean score within each grade and overall, taking into account the survey design.

```

. global cut = 10                                // number of classes to keep in each grade

. keep if classid <= $cut                        // keep only sampled classes
(1,131 observations deleted)

. scalar m = _N                                  // number of students in sample

. // get estimated population (should be close to 2003)
. gen weight = nclgrade / $cut                  // 1 / (n_h / N_h) or just (N_h / n_h)

. qui sum weight                                // quietly -summarize-

. scalar Mhat = r(sum)                          // store sum of weights

. di Mhat                                        // estimated population
1963

. // get population score estimate, by class and overall
. preserve

. gen wscore = testscore * weight                // (w_hij * y_hij)

. // same as double sum: students to class, class to grade
. collapse (sum) wscore (first) nstgrade nclgrade, by(grade)

. gen Ybar_grade = wscore / nstgrade             // Ytotal_h / stupop_h = Ybar_h

. scalar Ybar9 = Ybar_grade[1]                  // 9th grade average score

. scalar Ybar10 = Ybar_grade[2]                 // 10th grade average score

. scalar Ybar11 = Ybar_grade[3]                 // 11th grade average score

. scalar Ybar12 = Ybar_grade[4]                 // 12th grade average score

```

```

. qui sum wscore                                // quietly -summarize-

. scalar Ybar_school = r(sum) / Mhat            // Ytotal / stupop = Ybar

. restore

. scalar list Ybar9 Ybar10 Ybar11 Ybar12 Ybar_school
      Ybar9 = 475.31216
      Ybar10 = 489.36139
      Ybar11 = 521.64276
      Ybar12 = 515.93848
Ybar_school = 510.17983

```

Estimated standard error of the mean We can see that our estimate of the population, \hat{M} , is close to the true value, but not exact. Our within grade estimates aren't that great overall, but the schoolwide estimate is pretty close.

First let's look at an unweighted estimate of the schoolwide sample mean and its standard error.

```

. mean testscore

Mean estimation              Number of obs   =           872

-----
      |      Mean   Std. Err.   [95% Conf. Interval]
-----+-----
 testscore |   511.1055   1.852679   507.4693   514.7417
-----

```

Next, let's try compute an estimate of the variance. Note that the equations above speak to the variance of the total. We don't want that. We want the variance of the mean score. Here's what we will do: compute the variance of the total, divide it by the square of the estimated number of *SSUs* to standardize it, and then divide it again by the number of sampled *SSUs* to get the standard error of the estimated mean score. This isn't quite right, as we don't really take into account the clustering of students when we divide, but it will get us a reasonable approximation without recourse to more complicated methods.

```

. gen wscore = testscore * weight                // (w_hij * y_hij)

. gen fpc_r = $cut / nclgrade                    // fpc rate by strata (grade)

. collapse (sum) wscore (first) nstgrade fpc_r, by(grade classid) // sum w/n class

. preserve

. collapse (mean) stmscore = wscore, by(grade) // mean weighted score w/n grades

. tempfile stratmeans                            // init temporary file

. save `stratmeans'                             // save temporary file
file /var/folders/p7/5ylvzpld32x89bnkw42jj7dw0000gr/T//St10270.000004 saved

. restore

```

```

. merge m:1 grade using `stratmeans', nogen    // merge grade means into file

      Result                                # of obs.
-----
not matched                                0
matched                                   40
-----

. // (1 - f_h) * (n_h / (n_h - 1)) * (y_hi - ybar_h)^2
. gen adjsqdiff = (1 - fpc_r) * ($cut / ($cut - 1)) * (wscore - stmscore)^2

. collapse (sum) adjsqdiff                // double sum: w/n strata, overall

. // var(Ybar) = var(total) / stupop^2; Ybar_sem = sd(var(Ybar)) / sqrt(sampstu)
. scalar Ybar_school_sem = sqrt(adjsqdiff / Mhat^2) / sqrt(m)

. scalar list Ybar_school Ybar_school_sem
Ybar_school = 510.17983
Ybar_school_sem = .33002674

```

While our estimate of the schoolwide mean is closer to the true mean than the naive estimation, our standard error is much improved. Great! We should be cautious, however, of the standard error. Wait, what? This is due to the fact that the standard error of complex survey designs cannot be directly computed, only estimated. Our estimate might be too generous. It likely is. There are better, albeit more complicated ways to compute the estimate we want.

Good news

Now that we've gone through this process, the good news is that with most national educational surveys, you won't have to compute weights or figure out means and variances by hand. Instead, the data files will give you the weights you need. Stata also has prepackaged routines to help you in this process. The most important one is `svyset` and its suite of commands. We will discuss these in the next lecture.

Init: 23 August 2015; Updated: 24 August 2015