

## Final Sampling Project Report

*Team Heeringa:*  
Chia Wen Cheng  
Makenna Harrison  
Noah R Marcotte  
Wenqing Qian  
Zhuoyu Wang

SURVMETH 625: Applied Sampling

19 April 2023

University of Michigan, Ann Arbor

## 1. Population and Frame Characteristic

For this study, our goal is to monitor rates of drug use and cigarette smoking among teenagers in the State of Michigan. Our client is the Michigan Department of Education. To produce estimates of these behaviors, we will draw a sample of students to participate in our research. Here, our population is students currently enrolled in grades 7 through 12 in schools across the 83 counties in Michigan. We will implement a two-stage sampling procedure. First, we will sample schools, then we will sample students within each selected school. Therefore, our sampling frame has two components. For the first stage, we will use a list of public and private schools in the state. This list was obtained from Michigan's Department of Education in 2022. Then, for the second stage frame we will use a list of currently enrolled students for each school. Those schools selected will serve as our primary sampling units (PSU). In addition, we will use previous year's classroom enrollment as measures of size for each school.

Additionally, decimal numbers in this report are rounded to two digits for presentation purposes, but when performing calculations with those numbers, we use their full, unrounded fractions to ensure accuracy.

## 2. Overall Design

Our client is interested in the following three key variables related to teenage smoking and drug use, "ever smoked one cigarette", "ever smoked marijuana", and "age when first approached to smoke cigarettes or marijuana". The Department of Education would like to generate estimates of means and proportions for each variable, having a coefficient of variation ( $cv$ ) of no more than 0.05. We calculate these estimates along with their respective element variance, desired sampling variance (level of precision), and estimated necessary minimum sample size to maintain a simple random sample (SRS). These estimated values are provided in *Table 2.1*. Additionally, the calculations to achieve these values are as follows:

- (1) Estimates of the element variances ( $s^2$ ) for each variable sample proportion ( $p$ );
  - For variables "ever smoked one cigarette" and "ever smoked marijuana", we use the formula:  $s^2 = p(1 - p)$ , (*ignoring the finite population correction*)
- (2) Estimates of the desired sampling variance for each variable mean( $\bar{y}$ )/proportion were derived from the following inequality by solving for the standard error ( $se$ ) of the estimate and squaring it,
  - $cv = \frac{se(\bar{y})}{\bar{y}} \leq 0.05 \rightarrow se(\bar{y}) \leq 0.05 * \bar{y}$  (the formula is the same for the proportion)
- (3) Now that we have the sampling variance for all three variables, we are able to calculate the sample size( $n$ ) needed for SRS by the following formulas.
  - For variables "ever smoked one cigarette" and "ever smoked marijuana", we use the formula:  $Var(p) = \frac{1}{n-1} p(1 - p) \rightarrow n = \frac{p(1-p)}{Var(p)} + 1$
  - For variable "age when first approach to smoke cigarettes or marijuana", we use the formula:  $Var(\bar{y}) = \frac{1}{n} s^2 \rightarrow n = \frac{1}{Var(\bar{y})} s^2$

**Table 2.1: Sampling variance and sample size needed in an SRS design (cv = 0.05)**

Variable	Sample Mean or Proportion	Element Variance	Desired Sampling Variance	Sample Size Needed for SRS	Design Effect	roh
"ever smoked one cigarette"	0.25	0.19	0.0002	1201	2.5	0.03
"ever smoked marijuana"	0.15	0.13	0.00006	2268	2	0.02
"age when first approach to smoke cigarettes or marijuana"	12	1	0.36	3	1.7	0.01

Furthermore, a similar study was previously conducted using a sample size of  $n = 7,500$  students between the ages of 13 and 19, selected from a total of  $a = 150$  clusters. From this study we obtained the design effect for each variable mean/proportion, and using this information, computed synthetic estimates of roh (*Table 2.1*). The calculation of synthetic rate of homogeneity (roh) is as follows ( $b$  is the number of elements selected within each cluster):

$$b = \frac{n}{a} = \frac{7500}{150} = 50 \quad roh = \frac{deff - 1}{b - 1}$$

The total budget for data collection for this project is \$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. Using this information, we calculate the optimum sub-sample size ( $b_{opt}$ ), optimum cluster size ( $a_{opt}$ ), optimum overall sample size ( $n_{opt}$ ), and estimated cost for each of our key variables as summarized in *Table 2.2*. Since roh is portable and we had previously calculated it for each of the variables, we were able to use it in calculating our optimum sample size. The calculations are as follows, where  $c_a$  indicates cost per primary stage cluster and  $c_b$  equals cost per completed questionnaire within a cluster:

$$\begin{aligned} C &= 500000 \\ c_b &= 50 \\ c_a &= 3000 \\ b_{opt} &= \sqrt{\frac{c_a}{c_b} \frac{1 - roh}{roh}} \\ a_{opt} &= \frac{C}{c_a + b_{opt} * c_b} \\ deff &= 1 + (b_{opt} - 1) * roh \end{aligned}$$

**Table 2.2: Estimates of optimum sub-sample and expected cost and design effect**

Variable	$b_{opt}$	$a_{opt}$	$n_{opt}$	Estimated Cost	Expected Design Effect
"ever smoked one cigarette"	44	96	4224	\$499,200	2.32
"ever smoked marijuana"	54	88	4752	\$501,600	2.08
"age when first approached to smoke cigarettes or marijuana"	64	81	5184	\$502,200	1.9

We decided to use  $b_{opt} = 64$  ( $b^* = 64$ ) for our overall optimum subsample size in our design under our fixed cost constraint of \$500,000. While slightly over budget, we chose this design because it gives us the largest overall sample size. Assuming equal sized clusters, using this roh in our calculations yielded the largest overall sample size ( $n_{opt} = 5184$ ) and the smallest design effect ( $deff = 1.9$ ). For our estimates of  $b_{opt}$  and  $a_{opt}$ , we rounded to the nearest integer. If we had truncated our values instead, we would have come under budget for all three designs. However, considering we are not too far over budget, we feel our choice to round to the nearest integer is justified.

This design has a smaller number of clusters and a larger number of subsamples than the other designs. While this design may have a larger optimum sample size, it has a higher sampling variance. However, we thought the increase in sampling variance was not large enough to constitute the use of another design. Using other roh values resulted in smaller sample sizes and larger design effects. With this design, we take a subsample of 64 students from 81 different schools. Using the estimated size of the target population provided by our client ( $N = 830138$ ), our sampling fraction ( $f$ ) for this design is 0.00624.

Table 2.3 shows the design effect ( $deff_{(2)}$ ), effective sample size ( $n_{eff}$ ), and expected coefficient of variation ( $cv(\bar{y})$ ) under our design, assuming the portability of estimates of roh and the element variance. All key variable estimates meet the client's request to have an expected coefficient of variation 0.05. Calculations are as follows, where  $b_{(2),opt} = 64$ ,  $n_{opt} = 5184$ ,  $p$  equals the proportion for "ever smoked one cigarette" or "ever smoked marijuana" and  $\bar{y}$  is the mean for "age when first approached to smoke cigarettes or marijuana".

$$deff_{(2)} = 1 + (b_{(2),opt} - 1) * roh; var_{(2),SRS}(p) = \frac{p(1-p)}{n_2} \text{ or } var_{(2),SRS}(\bar{y}) = \frac{se(\bar{y})}{n_2}$$

$$var(p) = deff_2 * var_{(2),SRS}(p) \text{ or } var(\bar{y}) = deff_2 * var_{(2),SRS}(\bar{y})$$

$$n_{eff} = \frac{n_{opt}}{deff_{(2)}}; cv(\bar{y}) = \frac{se(\bar{y})}{\bar{y}}, \text{ where } \bar{y} \text{ is the proportion or mean of the variable and } se(\bar{y}) \text{ is the square root of the design sampling variance.}$$

**Table 2.3: Expected design effects given  $b_{opt} = 64$**

Variable	Design Effect		SRS Variance		Design Sampling Variance		Effective Sample Size	Expected Coefficient of Variation
	Total	20%	Total	20%	Total	20%	Total Sample	Total Sample
<i>“ever smoked one cigarette”</i>	2.93	1.37	4e-5	2e-4	1e-4	2e-4	1770	0.04
<i>“ever smoked marijuana”</i>	2.29	1.24	2e-5	1e-4	5e-5	2e-4	2268	0.05
<i>“age when first approach to smoke cigarettes or marijuana”</i>	1.9	1.17	2e-4	1e-3	4e-4	7e-4	2728	0.002

Additionally, if we were to calculate a hypothetical subclass of 20% that is evenly distributed across strata/schools, we would take 20% from our  $b$  and  $n$ ,  $b_{opt} = 64 * .20 = 13$ ; and  $n_{opt} = 5184 * .20 = 1037$ . Here, we calculate the design effect, SRS variance, design sampling variance, and expected coefficient of variation for this 20% subclass using the same formulas and sample proportions and means given to us for the total sample. Now that we have the estimated design sampling variance for the total sample and the 20% subclass, we can calculate confidence intervals for our variable proportion/mean estimates. The standard error is calculated as  $se = \sqrt{var_2(p)} = \sqrt{deff_2 * var_{(2),SRS}(p)}$  and the confidence interval is calculated as  $p \pm se * t_{0.975,80}$ . Where  $p$  is the sample mean/proportion,  $se$  is the standard error, and  $t_{0.975,80} = 1.99$  (the  $t$  score at the 95% confidence level with  $a-1 = 80$  degrees of freedom). Table 2.4 contains these estimates and their confidence intervals:

**Table 2.4: Sample Mean/Proportion with Confidence Intervals**

	Proportion/Mean (SE)		95% Confidence Interval	
	Total Sample	20% Subclass	Total Sample	20% Subclass
<i>“ever smoked one cigarette”</i>	0.25 (0.01)	0.25 (0.02)	[0.23,0.27]	[0.22, 0.28]
<i>“ever smoked marijuana”</i>	0.15 (0.01)	0.15 (0.01)	[0.14,0.16]	[0.13, 0.17]
<i>“age when first approached to smoke cigarettes or marijuana”</i>	11 (0.02)	11 (0.03)	[11.96, 12.04]	[11.93, 12.07]

### 3. Allocations

In this section, we will allocate how many schools to sample from each stratum and how many students to subsample from each selected school. In order to maintain *EPSEM* (equal probability of selection) across all strata, we will use the stratified *Probability Proportion to estimate Size* (PPeS) design to sample clusters and subsample students.

Our sampling fraction (not inflated) is calculated as:  $f = \frac{a_{opt} \times b_{opt}}{N} = \frac{5148}{830,183} = 0.00624$

As previously stated, our optimum number of clusters ( $a_{opt}$ ) is 81. We allocate the number of clusters to each stratum proportionate to the number of students in each stratum ( $W_h$ ). We do this by multiplying  $a_{opt}$  by the proportion of students in each stratum ( $W_h$ ). Thus, we have the number of expected schools,  $a_h$ , for each stratum. The overall student count from the previous year was treated as the measure of size ( $MOS$ ) for the PPeS sampling.

$$\text{According to PPeS: } f_h = f_1 \times f_2 = \left( \frac{a_h MOS_{ha}}{\sum_{a=1}^{A_h} MOS_{ha}} \right) \times \left( \frac{b_h^*}{MOS_{ha}} \right) \Rightarrow b_h^* = f_h \times \frac{\sum_{a=1}^{A_h} MOS_{ha}}{a_h}$$

$f_1$  is the first stage sampling rate and  $f_2$  is the second stage sampling rate.  $A_h$  is the number of schools for each stratum.  $a_h$  is the indicator for the school in stratum  $h$ .  $MOS_{ha}$  is the measures of size of school  $a$  in stratum  $h$ . We calculate the target subsample size for each stratum ( $b_h^*$ ) using the above formula. These calculated values are presented below (Table 3.1):

**Table 3.1: Uninflated PPeS Allocation**

$h$	$n_h$	$a_h$	$b_h^*$
1	74	1	74.13
2	57	1	56.97
3	60	1	59.89
4	51	1	50.53
5	132	2	65.59

$h$	$n_h$	$a_h$	$b_h^*$
6	189	3	62.70
7	1152	18	64.46
8	1152	18	64.44
9	2304	36	63.94
<b>Total</b>	<b>5171</b>	<b>81</b>	

However, since we expect a 30% response rates among schools, and a 70% response rate among students, we will need to generate inflated estimates for  $a_{opt}$  and  $b_{opt}$  (see Table 3.2):

**Table 3.2: Inflated Cluster and Subsample Sizes**

$b_{infl}$	$a_{infl}$	$n_{infl}$	Sampling fraction
92	270	24686	0.0297

Furthermore, we will need to apply this same inflation technique to the first and second stage sampling across the 9 strata. We derive the inflated number of clusters ( $a_{h, infl}$ ) by dividing each  $a_h$  by 0.3 and rounding to the nearest integer. This same procedure will need to be applied to subsample sizes in each stratum ( $b_h^*$ ), where we will divide by 0.7 to get  $b_{h, infl}^*$  (see *Table 3.3*).

To select clusters from strata, we stratify based on two auxiliary variables we believe to be related to students' smoking behavior. First, we achieve implicit stratification within the regions by geographic location. Here, we ordered the 83 counties into a list moving across the map (see *Appendix II*) in an s-shaped to ensure geographically contiguous units were adjacent to each other. Then, within each county, we ordered all the schools by type (public ahead of private).

Next, we calculate zone sizes for each stratum. We first define the sampling interval for each stratum ( $k_h$ ) as  $MOS_{ha}/a_{h, infl}$  (see *Table 3.3*). For each stratum, if the value of  $k_h$  is a non-integer, we will multiply it by either 10 or 100 (depending on the number of decimals) to create a new interval ( $k_{h, new}$ ). Next, we selected a random start between 1 and  $k_{h, new}$  using `sample()` function in R. Then, we add  $k_{h, new}$  repeatedly until we obtained a total of  $a_{h, infl}$  numbers. We then divide each selected case by the same number it was multiplied by earlier (1, 10, or 100), and truncate the decimals to get integers. We then used this sequence of integers to make our selection from the cumulative counts by taking schools where the cumulative number of students first exceeded  $k_{h, new}$ . *Table 3.3* enumerates the inflated values alongside the number of schools for each stratum ( $n_h$ ), the estimated number of students in each strata ( $MOS_{ha}$ ), the random start for PSU selection, and minimum measure of size.

**Table 3.3: Inflated PPeS Allocation**

$h$	$n_h$	$a_{h, infl}$	$MOS_{ha}$	$b_{h, infl}^*$	$k_h$	Random Start (seed: 625)	Minimum MOS <sup>1</sup>
1	74	1	3561	105.89	3561	2226	105
2	57	2	5474	81.39	2737	543	82
3	60	3	8631	85.55	2877	415	85
4	51	2	4855	72.19	2427.5	13268	73
5	132	6	18907	93.71	3151.17	71675	94
6	189	11	33133	89.57	3012.09	127418	89
7	1152	62	191992	92.08	3096.65	249087	93
8	1152	61	188830	92.05	3095.57	34754	93
9	2304	122	374755	91.34	3071.76	214121	91

<sup>1</sup> We set seed at 1000 and generated the following 9 random numbers: 0.34993746, 0.75546158, 0.31708425, 0.86581282, 0.76416075, 0.07288487, 0.49261774, 0.63516203, 0.08485787.

#### 4. An Illustration of Cluster Selections in Strata 8 & 9 and Unit Selections in Stratum 7

##### *First-stage Selection (Strata 8 & 9)*

We implemented first-stage selection for all nine strata using our inflated estimates ( $a_{h, infl}$ ). For the sake of illustration, we will only discuss in detail the sampling procedure and results for the last two strata. This same selection methodology should be applied to strata 1-7 (see *Table 3.3*).

The total number of schools in these last strata are 61 and 122 (for Regions 8 and 9, respectively). Within each strata we adopted a systematic sampling to select schools. Here, we calculated the sampling interval for both strata ( $k_h$ ) based on the sum of all high schools in the region ( $MOS_{ha}$ ) and the inflated number of schools selected from the region ( $a_{h, infl}$ ). There are 188830 students in Region 8 and 374755 students in Region 9. Thus, for Region 8, the sampling interval ( $k_g$ ) is calculated to be 3095.57, and Region 9's sampling interval ( $k_g$ ) is 3071.76. For Stratum 8, we selected a random start between 1 and 309,557 ( $k_{8, new}$ ): 34,754. We added  $k_{8, new}$  60 times to our random start to obtain 61 schools for selection. Each time, we truncated the last two digits from whole integers which corresponded to the selected schools. We applied this same sampling procedure to Stratum 9. Here, our random start was selected from between 1 and 307176 ( $k_{9, new}$ ): 214121. We then repeatedly added  $k_{9, new}$  to the random start 121 times, removed the last two digits, to get an index for each of the 122 schools we select.

There were no oversized units in either of the strata. There were also no selected units in Region 8 that failed to meet the minimum sufficient size. However, 7 units adjacent to the selected school were of insufficient size. Therefore, we implemented a "linking after selection" procedure to combine these undersized units with neighboring schools to yield large enough units for selection.

For instance, *Skyline High School* is one of the selected units, whose size is 420. However, its next unit *Early College Alliance - Milan* only has 26 students. The two schools following *Early College Alliance - Milan* — *Saline Alternative High School* (53) and *Washtenaw County Drop Back In Academy* (38) — are also not large enough until we reached the third one after it. Then we moved backwards, linking *Saline Alternative High School*, *Washtenaw County Drop Back In Academy*, and *Early College Alliance - Milan* to create a linked unit of size 117, and our selected unit *Skyline High School* ended up not being linked to any other units.

There are 4 schools not meeting the minimum sufficient size in Stratum 9: *Cloverdale Dvlpmntl Training Ctr.* (49), *ST PAUL LUTHERAN SCHOOL* (39), *Commonwealth Community Development Academy* (66), *Hope of Detroit Academy* (45). Taking *Cloverdale Dvlpmntl Training Ctr.* as an example, the next school on the list (*Birmingham Covington School*) has already reached the minimum size, so we linked *Cloverdale Dvlpmntl Training Ctr.* with the school in front of it (*Clifford H. Smart Middle School*) and obtained a linked unit sized 734. This size is large enough, so we simply substituted the undersized school with the linked one. In addition, another 14 schools also need to be linked because those next to them on the list are too small in size. A similar procedure is followed when linking these undersized units. The complete information about unit linking in both strata is displayed in *Appendix III*. The full list of



first-stage units we sampled from Region 8 and 9 is included in the same appendix, with linked units highlighted in light green shades.

### *Second-stage selection (Stratum 7)*

The example school given is located in Region 7. In this stratum, the average cluster size is 64.46 before being inflated for nonresponse and 92.08 after inflation. Since its measure of size is 242, the within-PSU sampling rate ought to be  $\frac{b_h^*}{MOS_{h\alpha}} = \frac{92.08}{242} = 0.38$ . However, in reality there are only 219 students in the school. Thus, the actual subsample size should be  $x_{h\alpha} = \frac{b_h^*}{MOS_{h\alpha}} B_{h\alpha} = \frac{92.08}{242} \times 219 = 83.3327$  (where  $B_{h\alpha}$  is the actual size of the school ( $\alpha$ ) in stratum ( $h$ )). Here, the probability of selecting 83 students from this school is 66.73% and the probability of selecting 84 students is 33.27%. According to the random number we drew between 0 and 1, we finally sampled 83 students from Applesamp Middle School to maintain the within-PSU sampling rate and *epsem* across all students.

We still performed systematic sampling to select students, following a similar protocol as in first stage selection. Since the sampling interval is  $k = B_{h\alpha} / x_{h\alpha} = 2.63855$ , we specified a random starting point between 1 and 263 ( $k_{new}$ ) as 178. After repeatedly adding  $k_{new}$  and truncating to form integers, we generated 83 student IDs. The list of students selected is attached in *Appendix III*.

## **5. Estimation Procedures**

We now have the numbers of the expected schools for each region,  $a_h$ , as stratum, to plan for the nonresponse error. We then collapsed strata 1 & 2 and strata 3 & 4 because they only contained 1 PSU, and recoded the stratum in the ascending order of the regions. We generated seven pseudo-stratum for the next steps for sampling error computing unit (SECU) and variance estimation.

Within each of the seven pseudo-stratum, we form 2 SECUs. Thus, the first stratum has 1 expected school in each SECU; the second has 1; the third has 1; the fourth has 1 or 2; the fifth stratum has 9; the sixth has 9; the last stratum has 18. We will choose clusters randomly when forming SECUs.

Since we have two SECUs within each collapsed stratum, we are able to apply a paired difference method by treating SECUs within each collapsed stratum as pairs.

**Table 5.1: Sampling error codes**

<i>h</i>	<i>a<sub>h</sub></i>	collapsed stratum codes	number of SECUs	expected size per SECU
1	1	1	2	1
2	1			
3	1	2	2	1
4	1			
5	2	3	2	1
6	3	4	2	1 or 2
7	18	5	2	9
8	18	6	2	9
9	36	7	2	18

The estimations of the standard errors and the 95% confidence intervals of the mean/proportion for both the total sample and the 20% subclass sample are mentioned in the “overall design” section. You may refer to tables 2.4 to 2.5 for the calculations. To clarify our calculation process, again, we hold the assumption that the 20% subclass is evenly distributed across strata/school.

Taylor Series Linearization (TSL) can be used to approximate the variance of mean and the standard error. TSL is a method that obtains a first-order linear approximation for the ratio estimator, and then uses the variance estimate for this approximation to estimate the variance of the estimate itself. We select TSL because the three expected coefficients of variation for each variable of interest are less than 0.1, in which the second term in the TSL remains small enough for ignorance and the approximation works well under this circumstance. In addition, the nonresponse rate is fixed for the two sampling stages in our case, which does not raise concerns in adequacy of our final sample sizes. Last but not least, we use the stratified PPeS design, and thus do not need to weight sample totals.

To calculate the variance estimate of the total sample, we will use the following formula.

#### Formula for Variance Estimation (Method: Taylor Series)

$$Var(r) \approx (1/x^2)[var(y) + r^2 var(x) - 2rcov(y, x)]$$

$$var(x) = \sum_{h=1}^{H=7} (x_{h1} - x_{h2})^2$$

$$var(y) = \sum_{h=1}^{H=7} (y_{h1} - y_{h2})^2$$

$$cov(y, x) = \sum_{h=1}^{H=7} (x_{h1} - x_{h2})(y_{h1} - y_{h2})$$

$x$  : the total sample size

$y$  : the sum of the variable interested(sample totals) of total sample

$$y = \sum_{h=1}^{H=7} \sum_{\alpha=1}^A \sum_{\beta=1}^B y_{h\alpha\beta}$$

$r$ : the mean of the variable interested of total sample

$$r = y/x$$

$h$ : the indicator of the collapsed stratum

$H$ : the number of the collapsed strata

$y_{h1}$ : the sample totals in the first pseudo PSU in the  $h_{th}$  collapsed stratum

$x_{h1}$ : the sample size in the first pseudo PSU of the  $h_{th}$  collapsed stratum

For the variance estimate of the 20% subclass, we calculate the estimation of the standard error of a mean/proportion still using the above formula. However, the some meanings of the symbols change:

- (1) the  $x$  refers to the sample size of the 20% subclass;
- (2)  $y$  refers to the sum of the variable interested(sample totals) for the 20% subclass;
- (3)  $r$  refers to the mean of the variable interested for the 20% subclass;
- (4)  $y_{h1}$  refers to sample totals in the first pseudo PSU in the  $h_{th}$  SECU for the 20% subclass;
- (5)  $x_{h1}$  refers to the sample size in the first pseudo PSU of the  $h_{th}$  SECU for the 20% subclass.

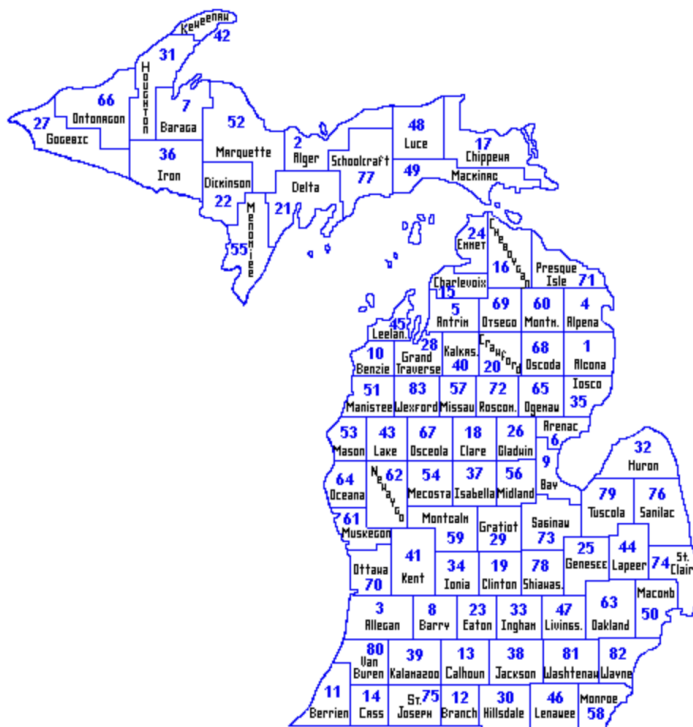
**Appendices:**  
**Appendix I: Notation Table**

Notation	Meaning	Notation	Meaning
$a$	number of clusters	$h$	number of stratum
$b$	the number of elements selected within each cluster	$a_h$	number of expected schools for each stratum
$C$	budget for data collection for this project	$a_{h, infl}$	inflated number of expected schools for each stratum
$c_a$	cost per primary stage cluster	$b_h^*$	target subsample size for each stratum
$c_b$	cost per completed questionnaire within a cluster	$b_{h, infl}^*$	the inflated subsample size of each element of each stratum
$cv$	coefficient of variation	$n_h$	sample size for each stratum
$var$	variance	$f_h$	sampling fraction of each stratum
$n$	sample size	$A_h$	number of schools for each stratum
$\bar{y}$	mean of population totals	$B_{h\alpha}$	actual size of the $\alpha$ _th school in $h$ _th stratum
$s^2$	estimates of the element variance	$MOS_{h\alpha}$	the measures of size of $\alpha$ _th school in $h$ _th stratum
$p$	sample proportion	$W_h$	the proportion of the number of students in stratum $h$ _th in total sample size
$deff$	design effect	$k$	sampling interval for second-stage sampling
$roh$	ratio of homogeneity	$r$	the mean of the variable interested from total sample

$f$	overall sampling fraction (without inflation)	$k_h$	sampling interval for each stratum
$t$	t-score	$k_{h, new}$	the new sampling interval of each stratum
SECU	sampling error computing unit	$x_{h\alpha}$	the actual subsample size
PSU	primary sampling units	$a_{opt}$	optimum cluster size
$f_1$	the first stage sampling rate	$b_{opt}$	optimum sub-sample size
$f_2$	the second stage sampling rate	$n_{opt}$	optimum overall sample size
		$n_{eff}$	effective sample size

## Appendix II: Michigan Counties Map used when sorting the school list

<https://www.michigan.gov/mdhhs/doing-business/licensing/child-welfare/institutions/online-lookups/michigan-counties-map>



**Appendix III: Linked schools in Strata 8 & 9**  
**Undersized and linked units in Stratum 8**

<b>School selected</b>	<b>Undersized units (Self/next/both)*</b>	<b>Size before linking</b>	<b>Schools linked</b>	<b>Total size after linking</b>
Skyline High School	Next	420	None	420
Holt Junior High School	Next	933	Stockbridge Alternative High School	953
Lake Fenton Middle School	Next	294	Special Services - North	307
LUTHERAN HIGH SCHOOL-NW	Next	309	OAKLAND STEINER SCHOOL ANN ARBOR CHRISTIAN SCHOOL	368
ST MICHAEL SCHOOL	Next	182	STS PETER & PAUL SCHOOL	199
LUTHERAN HIGH SCH-WESTLAND	Next	208	GO LIKE THE WIND SCHOOL BIBLE BAPTIST CHURCH SCHOOL SAINT MARTHA SCHOOL GTB LIFE LONG LEARNING CTR LOUISIANA HOMES SCHOOL	282
ST PATRICK SCHOOL	Next	96	ST PAUL SCHOOL	133

\* "Self" means the exact unit selected is undersized; "next" means the unit selected is not undersized but the next one on the list is; and "both" means both the selected unit and the one following are undersized.

**Undersized and linked units in Stratum 9**

<b>School selected</b>	<b>Undersized units (Self/next/both)</b>	<b>Size before linking</b>	<b>Schools linked</b>	<b>Total size after linking</b>
Uby Community High School	Next	264	Huron Learning Center	298
Croswell-Lexington Middle School	Next	395	Sanilac ISD Special Education Services; Marlette Learning Center	439
Yale Junior High School	Next	333	Phoenix Alternative School Capac Adult and Community Education St. Clair County Learning Academy	418
Cloverdale Dvlpmntl Training Ctr.	Self	49	Clifford H. Smart Middle School	734
Bloomfield Hills Andover H.S.	Next	951	Crossroads for Youth	1027

Hart Middle School	Next	774	Huron Valley Adult Education	800
MANISTEE CATHOLIC CENT SCHS	Next	98	MENOMINEE CATHOLIC CENTRAL	121
ST ELIZABETH ANN SETON M S	Next	161	ST JOSEPH'S SCHOOL	166
Armada Middle School	Next	310	Career Preparation Center	382
ST PAUL LUTHERAN SCHOOL	Self	39	ST MARY'S SCHOOL ST MARY/MC CORMICK CATHOLIC ACADEMY CONCORDIA LUTHERAN NORTH	101
Annapolis High School	Next	824	Dickinson Center	849
Ferguson Academy for Young Women	Next	286	Barnard Center	300
Communication and Media Arts HS	Next	516	Downriver High School	563
Colin Powell Academy	Next	237	Summit Academy	288
Commonwealth Community Development Academy	Both	66	Edison Public School Academy	307
Hope of Detroit Academy	Both	45	George Crockett Consortium High School	94
ST PATRICK SCHOOL	Next	102	ST PAUL LUTHERAN SCHOOL	130
LUTHERAN HIGH SCHOOL-NORTH	Next	628	None	628

#### Appendix V: List of schools sampled in Strata 8 & 9 (after linking)

##### Stratum 8 (after linking):

Selection number	Stratum	School	Cumulative MOS	Within School Interval
1	8	Davis Middle School	387	3.5539
2	8	Addison High School	3721	5.4162
3	8	Madison High School	6690	5.8041
4	8	Deerfield Public Schools	9660	2.6538
5	8	Chelsea High School	13571	14.7431
6	8	Forsythe Middle School	15844	6.8749

7	8	Milan High School	18983	15.0380
8	8	Willow Run High School	22108	8.5976
9	8	Lincoln Middle School	25338	11.5151
10	8	Community High School	28503	7.1698
11	8	Ypsilanti High School	31485	18.7781
12	8	Skyline High School	34787	6.5180
13	8	East Jackson High School	37801	6.7508
14	8	Michigan Center Jr/Sr High School	41068	11.0185
15	8	TA Wilson School	43724	4.2212
16	8	Columbia Central High School	47249	8.6441
17	8	Charlotte Senior High School	50426	14.7741
18	8	Grand Ledge High School	53938	26.7704
19	8	Waverly Senior High School	56403	16.4657
20	8	C.W. Otto Middle School	59273	6.7353
21	8	Everett High School	63742	25.1874
22	8	Holt Senior High School	66495	21.5870
23	8	Mason High School	69138	16.7916
24	8	Webberville High School	71679	3.3366
25	8	Haslett High School	75483	13.9982
26	8	Holt Junior High School Stockbridge Alternative High School	78190	14.7897
27	8	Brighton High School	82329	34.8093
28	8	Hartland High School	85585	27.9964
29	8	Howell High School	88838	40.0392
30	8	Pinckney High School	90361	23.6355
31	8	Fowlerville Junior High School	93649	7.4647
32	8	Beecher High School	96681	7.4336
33	8	Clio Area High School	99522	16.2329
34	8	Andrew G. Schmidt Middle School	102552	8.4579



35	8	George R. Carter Middle School	105710	8.1630
36	8	Grand Blanc Community High School	109022	39.7598
37	8	LakeVille High School	112158	9.2028
38	8	LakeVille Middle School	115138	4.0815
39	8	Bendle Middle School	118100	2.9486
40	8	Fenton Senior High School	121429	18.6384
41	8	Northern High School	124344	14.1999
42	8	Bendle/Carman-Ainsworth Alternative Education	127451	7.5423
43	8	Lake Fenton Middle School Special Services - North	130618	4.7644
44	8	LUTHERAN HIGH SCHOOL-NW OAKLAND STEINER SCHOOL ANN ARBOR CHRISTIAN SCHOOL	133628	5.7110
45	8	Owosso Middle School	136911	8.2561
46	8	Corunna High School	140055	12.6791
47	8	St. Johns Middle School	143193	8.5200
48	8	ST MICHAEL SCHOOL STS PETER & PAUL SCHOOL	146025	3.0883
49	8	Alma Middle School	149182	4.9506
50	8	Birch Run High School	152263	8.9545
51	8	White Pine Middle School	155283	12.4928
52	8	Saginaw High School	158593	14.7586
53	8	Heritage High School	161739	27.6084
54	8	Academy for Technology and Enterprise	164709	4.9351
55	8	LUTHERAN HIGH SCH-WESTLAND GO LIKE THE WIND SCHOOL BIBLE BAPTIST CHURCH SCHOOL SAINT MARTHA SCHOOL GTB LIFE LONG LEARNING CTR	167738	4.3764

		LOUISIANA HOMES SCHOOL		
56	8	John Glenn High School	171328	14.3396
57	8	Cramer Junior High School	173717	5.2299
58	8	ST PATRICK SCHOOL ST PAUL SCHOOL	176898	2.0640
59	8	Jefferson Middle School	180168	8.8614
60	8	H.H. Dow High School	184157	22.2078
61	8	West Intermediate School	186292	8.1941

**Stratum 9 (after linking):**

Selection number	Stratum	School	Cumulative MOS	Within School Interval
1	9	Uby Community High School Huron Learning Center	2290	4.6605
2	9	Unionville-Sebewaing High School	5281	4.8951
3	9	Carsonville-Port Sanilac H.S.	8517	4.4885
4	9	Croswell-Lexington Middle School Sanilac ISD Special Education Services Marlette Learning Center	11417	6.1775
5	9	Marine City Middle School	14569	5.1766
6	9	Port Huron Northern High School	18560	23.9908
7	9	Port Huron High School	20671	26.1646
8	9	Yale Junior High School Phoenix Alternative School Capac Adult and Community Education St. Clair County Learning Academy	23898	6.5373
9	9	Imlay City High School	26809	10.4627
10	9	Lapeer East Senior High School	30956	22.5207
11	9	Anderson Middle School	32981	5.7553
12	9	Brandon High School	36631	19.2833

13	9	Clarkston High School	40611	29.1674
14	9	Clifford H. Smart Middle School Cloverdale Dvlpmntl Training Ctr.	42103	11.4793
15	9	Farmington High School	46150	20.9567
16	9	Hazel Park High School	48864	16.9218
17	9	Milford High School	51887	26.4462
18	9	Page Middle School	54457	6.4747
19	9	Lake Orion Community High School	59143	37.9411
20	9	MacArthur K-8 University Academy	60529	1.5796
21	9	North Farmington High School	64530	21.8325
22	9	Oak Park High School	67459	20.9098
23	9	Rochester High School	70511	27.6504
24	9	Southfield High School	73216	20.5814
25	9	Troy High School	75983	32.9678
26	9	Waterford Kettering High School	79790	24.6320
27	9	Wylie E. Groves High School	82477	21.0506
28	9	Bloomfield Hills Andover H.S.; Crossroads for Youth	85209	16.0616
29	9	Southfield-Lathrup High School	89499	23.8969
30	9	Waterford Mott High School	91244	27.2907
31	9	Walled Lake Western High School	94735	25.1794
32	9	Oxford Area Middle School	97467	10.8068
33	9	Alice M. Birney Middle School	100509	7.2879
34	9	Novi Middle School	103475	15.8427
35	9	Athens High School	107213	29.0736
36	9	Lakeland High School	110916	27.1030
37	9	Hart Middle School Huron Valley Adult Education	113320	12.5115
38	9	Walled Lake Community Education Center	116010	3.5189

39	9	Academy of Lathrup Village	118927	1.8142
40	9	Walled Lake Northern High School	123292	25.7580
41	9	Stoney Creek High School	125493	25.3358
42	9	Life Skills Center of Pontiac	128131	2.7212
43	9	MANISTEE CATHOLIC CENT SCHS MENOMINEE CATHOLIC CENTRAL	131201	1.5327
44	9	ST ELIZABETH ANN SETON M S ST JOSEPH'S SCHOOL	134303	2.5961
45	9	Armada Area High School	137777	9.8684
46	9	Warren Mott High School	140505	28.2916
47	9	Wyandot Middle School	143843	10.0874
48	9	E.F. Siefert Elementary School	146682	2.8307
49	9	Fitzgerald Senior High School	150617	16.3431
50	9	Kelly Middle School	153010	5.9117
51	9	Lakeview High School	155761	18.0478
52	9	Lincoln High School	159261	14.0441
53	9	Frank Jeannette Jr. High School	162199	14.2475
54	9	Romeo High School	165606	29.2143
55	9	Shelby Junior High School	168205	20.5032
56	9	Utica High School	171406	21.0349
57	9	Adlai Stevenson High School	174495	32.5142
58	9	Carter Middle School	177343	8.3358
59	9	Richards Middle School	180729	12.5584
60	9	Sterling Heights Senior H.S.	184619	22.1923
61	9	Henry Ford II High School	188178	31.1067
62	9	Flynn Middle School	189815	6.9282
63	9	Armada Middle School Career Preparation Center	192666	5.9742

64	9	North Lake High School	195701	2.2208
65	9	Dakota High School	199707	38.5823
66	9	L'Anse Creuse Middle School - East	201820	8.5078
67	9	ST MARY'S SCHOOL ST MARY/MC CORMICK CATHOLIC ACADEMY CONCORDIA LUTHERAN NORTH ST PAUL LUTHERAN SCHOOL	204882	1.5796
68	9	Stevenson High School	208456	32.9052
69	9	South Middle School	211289	5.9430
70	9	Bryant Middle School	214461	7.4600
71	9	Central High School	217506	14.2944
72	9	Cody High School	220611	16.8749
73	9	Davidson Middle School	224093	12.9337
74	9	Denby High School	226940	17.3440
75	9	Ecorse Community High School	229502	7.0221
76	9	Finney High School	233266	12.6522
77	9	Fordson High School	236466	37.2217
78	9	Garden City High School	239720	23.0681
79	9	Grosse Pointe South High School	242437	25.6642
80	9	Ford High School	245885	23.2557
81	9	Huron High School	248605	13.8721
82	9	John D. Pierce Middle School	251212	8.1012
83	9	John Glenn High School	254642	33.1398
84	9	Kosciuszko School	257122	6.6311
85	9	Discovery Middle School	260733	10.7130
86	9	Mumford High School	264154	32.3735
87	9	Northville High School	267072	34.3284
88	9	Osborn High School	270473	20.4250
89	9	Pershing High School	273280	19.2364

90	9	Redford Union High School	276691	18.2511
91	9	Riverview Community High School	278612	13.9659
92	9	Romulus Senior High School	282282	19.4554
93	9	Allen Park Middle School	285067	9.7120
94	9	Southgate Anderson High School	288619	19.6274
95	9	Hoover Middle School	291135	9.9623
96	9	Wayne Memorial High School	295376	29.6679
97	9	Western International High School	297906	26.2272
98	9	Churchill High School	301364	32.7019
99	9	Annapolis High School Dickinson Center	303500	13.2778
100	9	Hally Magnet Middle School	306608	5.9899
101	9	Lowrey Middle School	309626	5.0984
102	9	Canton High School	313170	31.4977
103	9	Truman High School	316666	26.2741
104	9	Shumate Middle School	318574	9.1177
105	9	Ferguson Academy for Young Women Barnard Center	321890	4.6918
106	9	Communication and Media Arts HS Downriver High School	324744	8.8050
107	9	Colin Powell Academy Summit Academy	327865	4.5041
108	9	Henry Ford Academy	331152	7.5851
109	9	Edison Public School Academy Commonwealth Community Development Academy	333901	4.8013
110	9	Hillside Middle School	337456	8.0230
111	9	Plymouth High School	341345	31.7010

112	9	Hope of Detroit Academy George Crockett Consortium High School	343184	1.4701
113	9	Detroit Academy of Arts and Sciences High School	346499	11.6513
114	9	Detroit Premier Academy	349347	2.3303
115	9	Cody 9th Grade Academy	352667	6.0681
116	9	DIVINE CHILD HIGH SCHOOL	356250	13.8721
117	9	ST PATRICK SCHOOL ST PAUL LUTHERAN SCHOOL	358565	2.0331
118	9	LUTHERAN HIGH SCHOOL-NORTH	361549	9.8215
119	9	Bedford Senior High School	364444	26.7902
120	9	Jefferson Middle School	367871	5.1923
121	9	Bedford Junior High School	371256	13.0432
122	9	Dundee Middle School	373984	4.0662

## Appendix VI: List of Sampled Students from Applesamp Middle School

Grade	Surname	Given name(s)	ID
7	MAGEE	MONICA L	1
7	BELL	KENNETH MARK	4
7	SHY	RANDALL ERON	7
7	VINSON	STEVEN A	9
7	BURGARD	PETER JOSEPH	12
7	ABELES	JODI	14
7	WALL	SAMUEL THOMAS	17
7	BRENING	TERESA M	20
7	SEGEBARTH	DAVID	22
7	COHEN	SCOTT A	25
7	LUTZ	ADAM MARK	28
7	RAYNOR	GREGORY K	30
7	WILDSTEIN	EVAN	33
7	BARRY	MICHAEL J	35
7	DYKHOUSE	JAY D	38
7	HATCHER	MERRICK D	41
7	SERGEANT	STACY	43
7	MARIN	ERIK T	46
7	FOSTER	LORI A	49
7	HEADAPOHL	MARC DAVID	51
7	RENKER	COREY H	54
7	CURRY	TRACY YVETTE	57
7	BLACK	STEPHEN P	59
7	FRANCIS	DALE A	62
7	FEIGENBAUM	FRANK	64
7	CLARK	CHARLES M	67
7	LANGDON	JOHN DOUGLAS	70
7	TEMPLE	LOUISE ELIZABETH	72
7	WILLIAMS	LINDA K	75
7	ASHBROOK	ANDREW S	78
7	BRENNAN	T CASSEY	80
7	BEAUCAIRE	JERI	83
7	TANNER	DEBORAH SHIELDS	85
7	HINSLEY	PATRICIA NOEL	88
7	LAUNICZAK	MARA K	91
7	SCAMRNELLA	SANDRA ANNA MARIA	93
8	CONSTINE	RUSSELL P	96
8	ROGERS	ANTOINETTE DENISE	99



8	MOSS	ROBERT EARL	101
8	DYER	JENNIFER REBECCA	104
8	CUTHRELL	JASON ELLIS	106
8	NAUSS	MATTHEW P	109
8	COITEN	RICHARD S	112
8	KOLKMAN	ANN MARIE	114
8	KROSKY	DANIEL J	117
8	MILLS	RANDY SCOTT	120
8	CARSON	ERIK JOHN	122
8	KRAPFL	HEIDI R	125
8	BAVITZ	CHRISTOPHER THEODORE	128
8	GORNY	KAREN E	130
8	OCKER	MARK GRIFFITH	133
8	ROZENBLIT	IGOR A	135
8	MULLAN	BRIAN P	138
8	KANDIK	JULIE MARIE	141
8	FENSON	REBECCA R	143
8	JEZIC	NINA L	146
8	LEIBOVITZ	JAY BYRON	149
8	ALMLI	MARTA B	151
8	WEAVER	TRACEY DENISE	154
8	SPINDELMAN	MARC	156
8	SIMKIN	MARGARET RENEE	159
8	LOEVY	DEBRA L	162
8	HASKIN	JENNIFER L	164
8	DAHL	PHILLIP Y	167
8	EHRENBERG	JASON HERBERT	170
8	FROST	THORNTON HARRISON	172
8	ROMERO	ROSS I	175
8	BAILEY	THURSTON C	177
8	MENDEL	DAVID STUART	180
8	WILHELM	ANDREW J	183
8	CASCADE	JOSHUA CHARLES	185
8	BLANK	SALLY BAKER	188
8	GATTO	JULIA LYNN	191
8	MANSDORF	NICOLE LISA	193
8	EDIDIN	ERIC J	196
8	MARGOLIN	JONATHAN B	199
8	MORTON	MARCY DAWN	201
8	FELDSTEIN	RACHEL ELIZABETH	204
8	LEE	KAREN ELYSE	206

8	SILVER	ALLYSON	209
8	KUHN	LAURIE JOANNE	212
8	RISTIC	SYLIA C	214
8	CHILEWICH	MATTHEW ERIC	217