



NAVAL
POSTGRADUATE
SCHOOL



Stratified, Cluster & Complex Sampling

Professor Ron Fricker
Naval Postgraduate School
Monterey, California

Reading:
Scheaffer et al. chapters 5 & 8



Goals for this Lecture

- Define stratified & cluster sampling
- Discuss reasons for stratified or cluster sampling
- Derive estimators for the population
 - Mean
 - Total
 - Percentage
- Discuss examples of complex sampling designs
- Explain the Kish grid



Stratified Sampling

- **Stratified sampling** divides the sampling frame up into strata from which separate probability samples are drawn
- Examples
 - GWI survey, needed to obtain information from members of each military service
 - For external validity, WMD survey had to sample large urban areas
 - In a particular country, survey in support of IW requires information on both Muslim and Christian communities



Reasons for Stratified Sampling

- More precision
 - Assuming strata are relatively homogeneous, can reduce the variance in the sample statistic(s)
 - So, get better information for same sample size
- Estimates (of a particular precision) needed for subgroups
 - May be necessary to meet survey's objectives
 - May have to oversample to get sufficient observations from small strata
- Cost
 - May be able to reduce cost per observation if population stratified into convenient groupings

An Example



Record	Name	Group	
1	Bradburn, N.	High	
2	Cochran, W.	Highest	
3	Deming, W.	High	
4	Fuller, W.	Medium	
5	Habermann, H.	Medium	
6	Hansen, M.	Low	
7	Hunt, J.	Highest	
8	Hyde, H.	High	
9	Kalton, G.	Medium	→ Kalton, G.
10	Kish, L.	Low	
11	Madow, W.	Highest	
12	Mandela, N.	Highest	
13	Norwood, J.	Medium	→ Norwood, J.
14	Rubin, D.	Low	→ Rubin, D.
15	Sheatsley, P.	Low	
16	Steinberg, J.	Low	
17	Sudman, S.	High	
18	Wallman, K.	High	→ Wallman, K.
19	Wolfe, T.	Highest	
20	Woolsley, T.	Medium	

One SRS of Size 4

Figure 4.5 Frame population of 20 persons sorted alphabetically, with SRS sample realization of size $n = 4$.

Record	Name	Group	
2	Cochran, W.	Highest	
7	Hunt, J.	Highest	
11	Madow, W.	Highest	
12	Mandela, N.	Highest	
19	Wolfe, T.	Highest	→ Wolfe, T.
1	Bradburn, N.	High	→ Bradburn, N.
3	Deming, W.	High	
8	Hyde, H.	High	
17	Sudman, S.	High	
18	Wallman, K.	High	
4	Fuller, W.	Medium	→ Fuller, W.
5	Habermann, H.	Medium	
9	Kalton, G.	Medium	
13	Norwood, J.	Medium	
20	Woolsley, T.	Medium	
6	Hansen, M.	Low	
10	Kish, L.	Low	
14	Rubin, D.	Low	→ Rubin, D.
15	Sheatsley, P.	Low	
16	Steinberg, J.	Low	

One Stratified Random Sample of Total Size 4

Figure 4.6 Frame population of 20 persons sorted by group, with stratified element sample of size $n_h = 1$ from each stratum.

- When selecting a stratified random sample, must clearly specify the strata
 - Non-overlapping categories into which each sampling unit must be classified
 - Sampling units can only be in one strata
 - Strata based on information about whole population
- Can have more than one type of classification
 - E.g., officer/enlisted and male/female
 - Then each sampling unit (person in this case) must be classified into one of four stratum



- Additional notation for stratified sampling
 - L is the number of strata
 - N_i is the number of sampling units in stratum i
 - n_i is the sample size in stratum i
 - N is the total number of sampling units in the population: $N = N_1 + N_2 + \dots + N_L$
- In this lecture, we use SRS within each strata
 - That does not necessarily mean each strata has the same selection probabilities
 - And often more complex sampling done within the strata

Mean Estimation

- Within a strata, probability a unit is sampled is

$$n_i/N_i$$

- Via the SRS lecture, we know that an unbiased estimate of the mean of strata i is

$$\bar{y}_i = \frac{1}{N_i} \sum_{j=1}^{n_i} \frac{1}{\pi_i} y_{ij} = \frac{1}{N_i} \sum_{j=1}^{n_i} \frac{1}{n_i / N_i} y_{ij} = \frac{1}{N_i} \sum_{j=1}^{n_i} \frac{N_i}{n_i} y_{ij} = \frac{1}{n_i} \sum_{j=1}^{n_i} y_{ij}$$

- So the estimated total for strata i is $\hat{\tau}_i = N_i \bar{y}_i$
- And the population mean estimate is just the sum of all the estimated totals divided by the population size

$$\bar{y}_{st} = \frac{\hat{\tau}}{N} = \frac{1}{N} \sum_{i=1}^L \hat{\tau}_i = \frac{1}{N} \sum_{i=1}^L N_i \bar{y}_i$$

Mean Estimation Summary

- Estimator for the mean: $\bar{y}_{st} = \frac{1}{N} \sum_{i=1}^L N_i \bar{y}_i$
- Variance of \bar{y}_{st} :

$$\begin{aligned}\widehat{\text{Var}}(\bar{y}_{st}) &= \text{Var}\left(\frac{1}{N} \sum_{i=1}^L N_i \bar{y}_i\right) \\ &= \frac{1}{N^2} \sum_{i=1}^L N_i^2 \text{Var}(\bar{y}_i) \\ &= \frac{1}{N^2} \sum_{i=1}^L N_i^2 \left(1 - \frac{n_i}{N_i}\right) \left(\frac{s_i^2}{n_i}\right)\end{aligned}$$

- Bound on the error of estimation: $2\sqrt{\widehat{\text{Var}}(\bar{y}_{st})}$

Example: TV Viewing Time (1)

- Survey to estimate average TV viewing time
- Population consists of two urban areas and rural residents
- Reasons to stratify:
 - Similarity of viewing habits by towns and rural regions
 - Cost: perhaps cheaper to survey each region separately (?)
- Sample allocated proportional to strata size

TABLE 5.1

Television-viewing time, in hours per week

Town A	Town B	Rural
35	27	8
43	15	14
36	4	12
39	41	15
28	49	30
28	25	32
29	10	21
25	30	20
38		34
27		7
26		11
32		24
29		
40		
35		
41		
37		
31		
45		
34		

Example: TV Viewing Time (2)

FIGURE 5.1
Box plots of television-viewing time

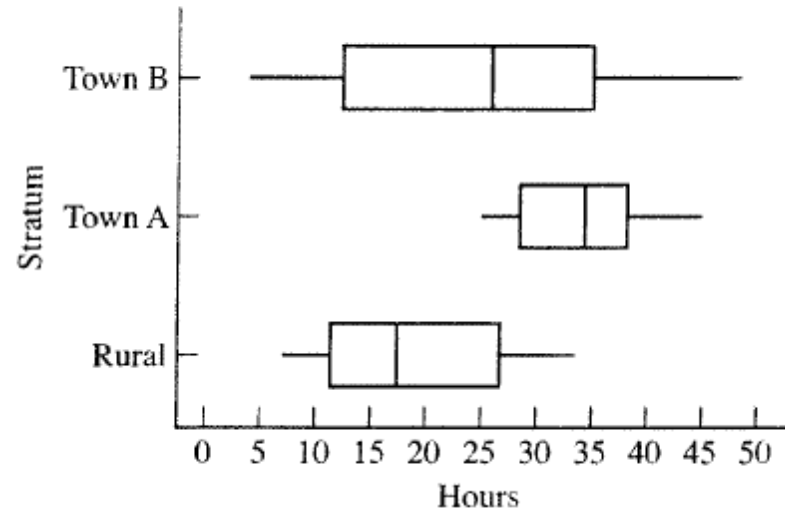


TABLE 5.2
Summary of the data from Table 5.1

	<i>N</i>	<i>n</i>	Mean	Median	SD
Town A	155	20	33.90	34.50	5.95
Town B	62	8	25.12	26.00	15.25
Rural	93	12	19.00	17.50	9.36

Example: TV Viewing Time (3)

A numerical descriptive summary of the data is shown in Table 5.2.

(a) From Table 5.2 and Eq. (5.1),

$$\begin{aligned}\bar{y}_{st} &= \frac{1}{N} [N_1 \bar{y}_1 + N_2 \bar{y}_2 + \cdots + N_L \bar{y}_L] \\ &= \frac{1}{310} [(155)(33.900) + (62)(25.125) + (93)(19.00)] \\ &= 27.7\end{aligned}$$

is the best estimate of the average number of hours per week that all households in the county spend watching television. Also,

$$\begin{aligned}\hat{V}(\bar{y}_{st}) &= \frac{1}{N^2} \sum N_i^2 \left(1 - \frac{n_i}{N_i}\right) \left(\frac{s_i^2}{n_i}\right) \\ &= \frac{1}{(310)^2} \left[\frac{(155)^2 (0.871) (5.95)^2}{20} + \frac{(62)^2 (0.871) (15.25)^2}{8} \right. \\ &\quad \left. + \frac{(93)^2 (0.871) (9.36)^2}{12} \right] \\ &= 1.97\end{aligned}$$

The estimate of the population mean with an approximate 2-SD bound on the error of estimation is given by

$$\bar{y}_{st} \pm 2\sqrt{\hat{V}(\bar{y}_{st})} \quad \text{or} \quad 27.675 \pm 2\sqrt{1.97} \quad \text{or} \quad 27.7 \pm 2.8$$

Example: TV Viewing Time (4)

- Correct analysis:

$$\bar{y}_{st} = 27.7, \hat{\sigma}_{\bar{y}_{st}} = 1.4, 95\% \text{ CI} = (24.9, 30.5)$$

- What if you just treated this as a SRS?

- Incorrect analysis with fpc:

$$\bar{y}_{SRS} = 27.7, \hat{\sigma}_{\bar{y}} = 10.6, 95\% \text{ CI} = (6.5, 48.9)$$

- Incorrect analysis without fpc:

$$\bar{y}_{SRS'} = 27.7, \hat{\sigma}_{\bar{y}} = 11.3, 95\% \text{ CI} = (5.1, 50.3)$$

Total Estimation Summary



- Estimator for the total: $\hat{\tau}_{st} = N\bar{y}_{st} = \sum_{i=1}^L N_i \bar{y}_i$
- Variance of $\hat{\tau}_{st}$:

$$\begin{aligned}\widehat{\text{Var}}(\hat{\tau}_{st}) &= \text{Var}\left(\sum_{i=1}^L N_i \bar{y}_i\right) \\ &= \sum_{i=1}^L N_i^2 \text{Var}(\bar{y}_i) \\ &= \sum_{i=1}^L N_i^2 \left(1 - \frac{n_i}{N_i}\right) \left(\frac{s_i^2}{n_i}\right)\end{aligned}$$

- Bound on the error of estimation: $2\sqrt{\widehat{\text{Var}}(\hat{\tau}_{st})}$

Proportion Estimation Summary

- Proportion estimator: $\hat{p}_{st} = \frac{1}{N} \sum_{i=1}^L N_i \hat{p}_i$
- Variance of \bar{y}_{st} :

$$\begin{aligned}\widehat{\text{Var}}(\hat{p}_{st}) &= \text{Var}\left(\frac{1}{N} \sum_{i=1}^L N_i \hat{p}_i\right) \\ &= \frac{1}{N^2} \sum_{i=1}^L N_i^2 \text{Var}(\hat{p}_i) \\ &= \frac{1}{N^2} \sum_{i=1}^L N_i^2 \left(1 - \frac{n_i}{N_i}\right) \left(\frac{\hat{p}_i(1 - \hat{p}_i)}{n_i}\right)\end{aligned}$$

- Bound on the error of estimation: $2\sqrt{\widehat{\text{Var}}(\hat{p}_{st})}$

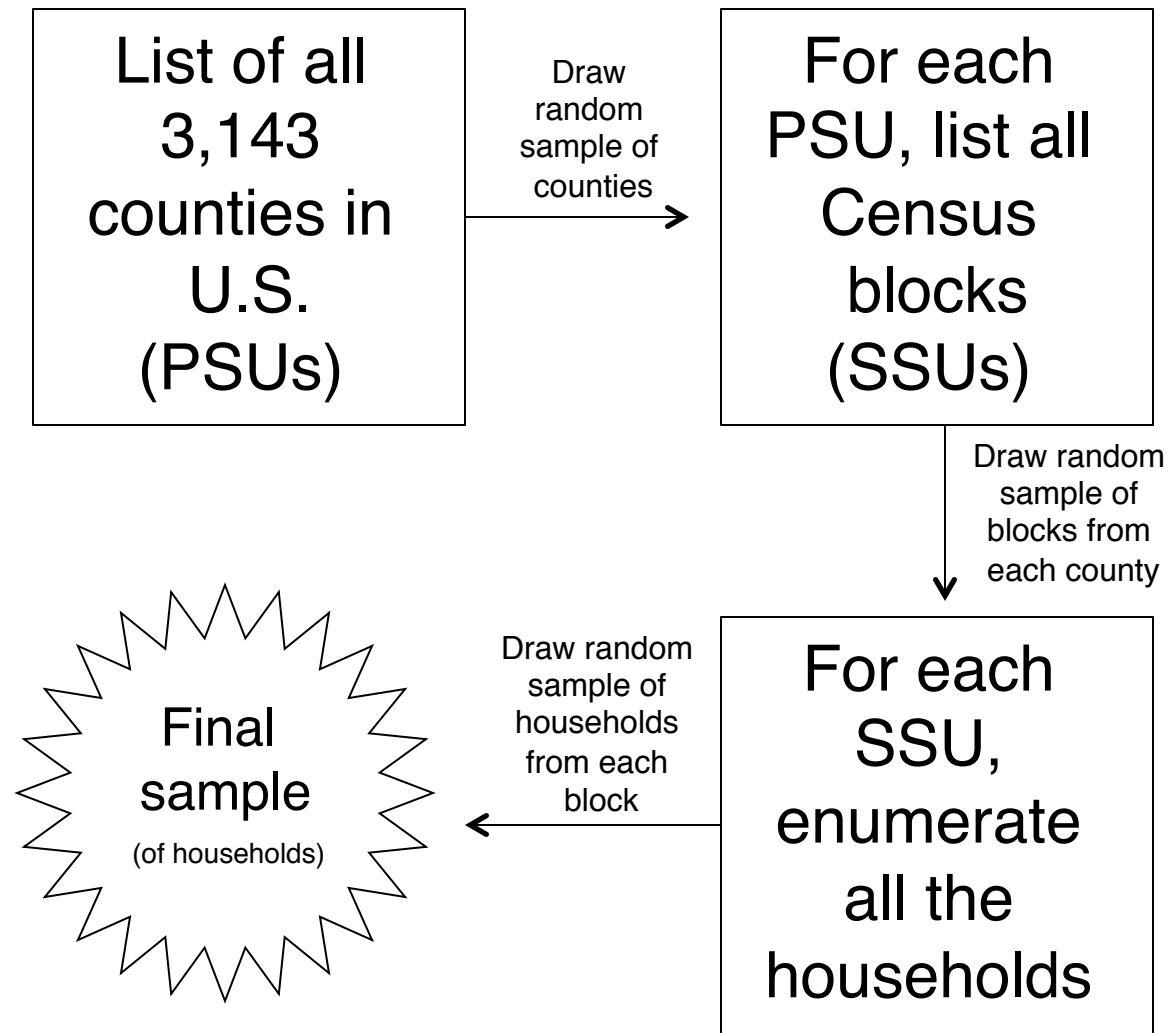
- **Poststratification** used when strata variable(s) not observed until after survey completed
 - That is, stratify on data collected in the survey
- Useful when sample demographics turn out not to match the population demographics
 - As with stratification, still need to know the distribution of the strata variable in the population
 - But don't need to know the values for each element in the sampling frame
- See Scheaffer et al. for more detail

What is Cluster Sampling?



- **Cluster sampling**: a probability sample in which each sampling unit is a collection, or cluster, of elements
 - Elements for survey occur in groups (clusters)
 - So, sampling unit is the cluster, not the element
 - Aka single stage cluster sampling
 - When sampling clusters by region, called area sampling
- There are more complicated types of cluster sampling such as **two-stage cluster sampling**
 - First select primary sampling units (PSUs) by probability sampling
 - Then within each selected PSU, sample secondary sampling units (SSUs) via probability sampling
 - Survey all units in each selected SSU

Illustration



Advantages and Disadvantages



- Advantages:
 - For some populations, cannot construct explicit sampling frame
 - But can construct a frame by cluster (area, organizational, etc.) then sample within
 - For some efforts, too expensive to conduct a SRS
 - E.g., drawing a SRS from the US population for an in-person interview
- Disadvantage: Cluster sampling generally provides less precision than SRS or stratified sampling

When To Use Cluster Sampling



- Use cluster sampling only when economically justified
 - I.e., when cost savings overcome (or require) loss in precision
- Most likely to occur when
 - Constructing a complete list-based sampling frame is difficult, expensive, or impossible
 - The population is located in natural clusters (schools, city blocks, etc.)

Example: Small Village

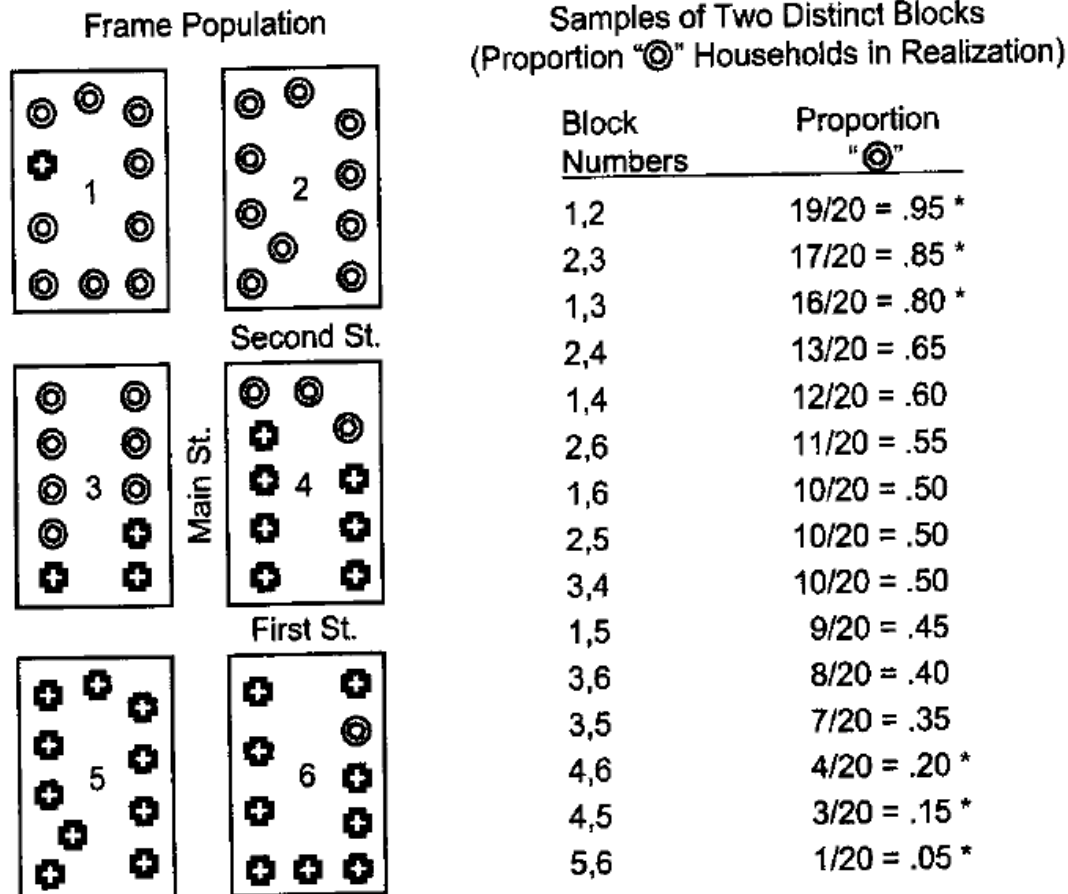


Figure 4.4 A bird's eye view of a population of 30 “+” and 30 “•” households clustered into six city blocks, from which two blocks are selected.



Effect of Cluster Sampling

- First, note that the true proportion is 0.5
 - Further, sampling 2 blocks out of 6 produces an unbiased estimate
- Now note that 6 of the 15 possible samples had proportions > 0.8 or < 0.2
 - That is, 40 percent of the samples are “far off”
- With a SRS*, the probability of getting 16 or more of the same type of household is about 0.3 percent...

Cluster Sampling Notation

- Lots of notation for cluster sampling:
 - M is the number of clusters in the population
 - m is the number of clusters selected via SRS
 - n_i is the number of elements in cluster i
 - $N = \sum_{i=1}^M n_i$, the total number of population elements
 - $\bar{N} = N / M$, the average cluster size in the population
 - $\bar{n} = \frac{1}{m} \sum_{i=1}^m n_i$, the average cluster size in the sample
 - y_{ij} is the j^{th} observation in the i^{th} cluster
 - $t_i = \sum_{j=1}^{n_i} y_{ij}$, the total of the i^{th} sampled cluster

Mean Estimation Summary

- Estimator for the mean: $\bar{y}_{cl} = \frac{1}{\sum_{i=1}^m n_i} \sum_{i=1}^m t_i = \frac{1}{m\bar{n}} \sum_{i=1}^m t_i$
 - Variance of \bar{y}_{cl} : $\widehat{\text{Var}}(\bar{y}_{cl}) = \left(1 - \frac{m}{M}\right) \frac{s_r^2}{m\bar{N}^2}$
- where $s_r^2 = \frac{1}{m-1} \sum_{i=1}^m (t_i - \bar{y}_{cl} n_i)^2$
- ↖ If \bar{N} not known, estimate with \bar{n}

✓ Key idea: Only the clusters are random

Total Estimation Summary



- Estimator for the total: $\hat{\tau} = N \bar{y}_{cl} = \frac{N}{\sum_{i=1}^m n_i} \sum_{i=1}^m t_i = \frac{N}{m\bar{n}} \sum_{i=1}^m t_i$
- Variance of $\hat{\tau}$:

$$\begin{aligned}\widehat{\text{Var}}(\hat{\tau}) &= N^2 \widehat{\text{Var}}(\bar{y}_{cl}) \\ &= N^2 \left(1 - \frac{m}{M} \right) \frac{s_r^2}{m\bar{N}^2} \\ &= M^2 \left(1 - \frac{m}{M} \right) \frac{s_r^2}{m}\end{aligned}$$

where $s_r^2 = \frac{1}{m-1} \sum_{i=1}^m (t_i - \bar{y}_{cl} n_i)^2$

What if N is Not Known?

- Alternate estimator for the total:

$$\hat{\tau} = M \times \text{average cluster total} = M \times \frac{1}{m} \sum_{i=1}^m t_i \equiv M \times \bar{y}_t$$

- Where now, variance of $\hat{\tau}$ is

$$\widehat{\text{Var}}(\hat{\tau}) = M^2 \left(1 - \frac{m}{M} \right) \frac{s_t^2}{m}$$

with $s_t^2 = \frac{1}{m-1} \sum_{i=1}^m (t_i - \bar{y}_t)^2$

- Again, note that the calculations are based on the idea that only the clusters are random

Example: NAEP



- Assume:
 - 40,000 4th grade classrooms in US
 - $n_i = 25$ students per classroom
- Sampling procedure:
 - Select m classrooms
 - Visit each classroom and collect data on all students
 - If $m = 8$, will have data on 200 students
- Note the differences from SRS
 - All groups of 200 students cannot be sampled
 - Students in each classroom more likely to be alike

Mean & Variance Computations

- Calculate the mean test score as

$$\bar{y}_{cl} = \frac{1}{m\bar{n}} \sum_{i=1}^m t_i = \frac{1}{8 \times 25} \sum_{i=1}^8 \sum_{j=1}^{25} y_{ij}$$

- The variance is

$$\widehat{\text{Var}}(\bar{y}_{cl}) = \left(1 - \frac{8}{40,000}\right) \frac{s_r^2}{8 \times 25^2} \approx \frac{s_r^2}{8 \times 25^2}$$

$$\text{where } s_r^2 = \left(\frac{1}{7}\right) \sum_{i=1}^8 (t_i - 25\bar{y}_{cl})^2$$

Sampling with Probability Proportional to Size (pps)



- Sometimes it makes sense to sample clusters in proportion to their size
 - Puts more emphasis on the larger clusters where most of the observations are
 - E.g., sampling clusters equally when there are a lot of small clusters and only a few very large ones could be very inefficient
- For cluster sampling, sampling with **probability proportional to size** (pps) means

$$\Pr(\text{select cluster } i) = n_i / N$$

Estimation Summary With PPS

- Plugging $\pi_i = n_i / N$ into the H-T estimator gives:

- Mean

- Estimator for the mean: $\hat{\mu}_{pps} = \frac{1}{m} \sum_{i=1}^m \bar{y}_i$

- Variance of $\hat{\mu}_{pps}$: $\widehat{\text{Var}}(\hat{\mu}_{pps}) = \frac{1}{m(m-1)} \sum_{i=1}^m (\bar{y}_i - \hat{\mu}_{pps})^2$

- Total

- Estimator for the total: $\hat{\tau}_{pps} = \frac{N}{m} \sum_{i=1}^m \bar{y}_i$

- Variance of $\hat{\tau}_{pps}$: $\widehat{\text{Var}}(\hat{\tau}_{pps}) = \frac{N^2}{m(m-1)} \sum_{i=1}^m (\bar{y}_i - \hat{\mu}_{pps})^2$

Complex Sampling for Real-World Surveying



- Usually, real world requirements and constraints result in complex sampling
 - Some combination of stratification and clustering along with unequal sampling probabilities
- For example, geographic clustering arises with face-to-face interviewer-based surveys
 - Often it's multi-stage clustering as well
- Stratification often also necessary to ensure desired representation in sample
- When combined, estimation gets much more complicated



NAEP Sampling Scheme

- First stage: 96 PSUs consisting of metropolitan statistical areas (MSAs), a single non-MSA county, or a group of contiguous non-MSA counties
 - About a third of the PSUs are sampled with certainty
 - Remainder are stratified and one selected from each stratum with probability proportional to size
- Second stage: selection public and nonpublic schools within the PSUs
 - For elementary, middle, and secondary samples, independent samples of schools are selected with probability proportional to measures of size
- Third and final stage: 25 to 30 eligible students are sampled systematically with probabilities designed to make the overall selection probabilities approximately constant
 - Except students from private schools and schools with high proportions of black or Hispanic students oversampled
- In 1996 nearly 150,000 students were tested from just over 2,000 participating schools

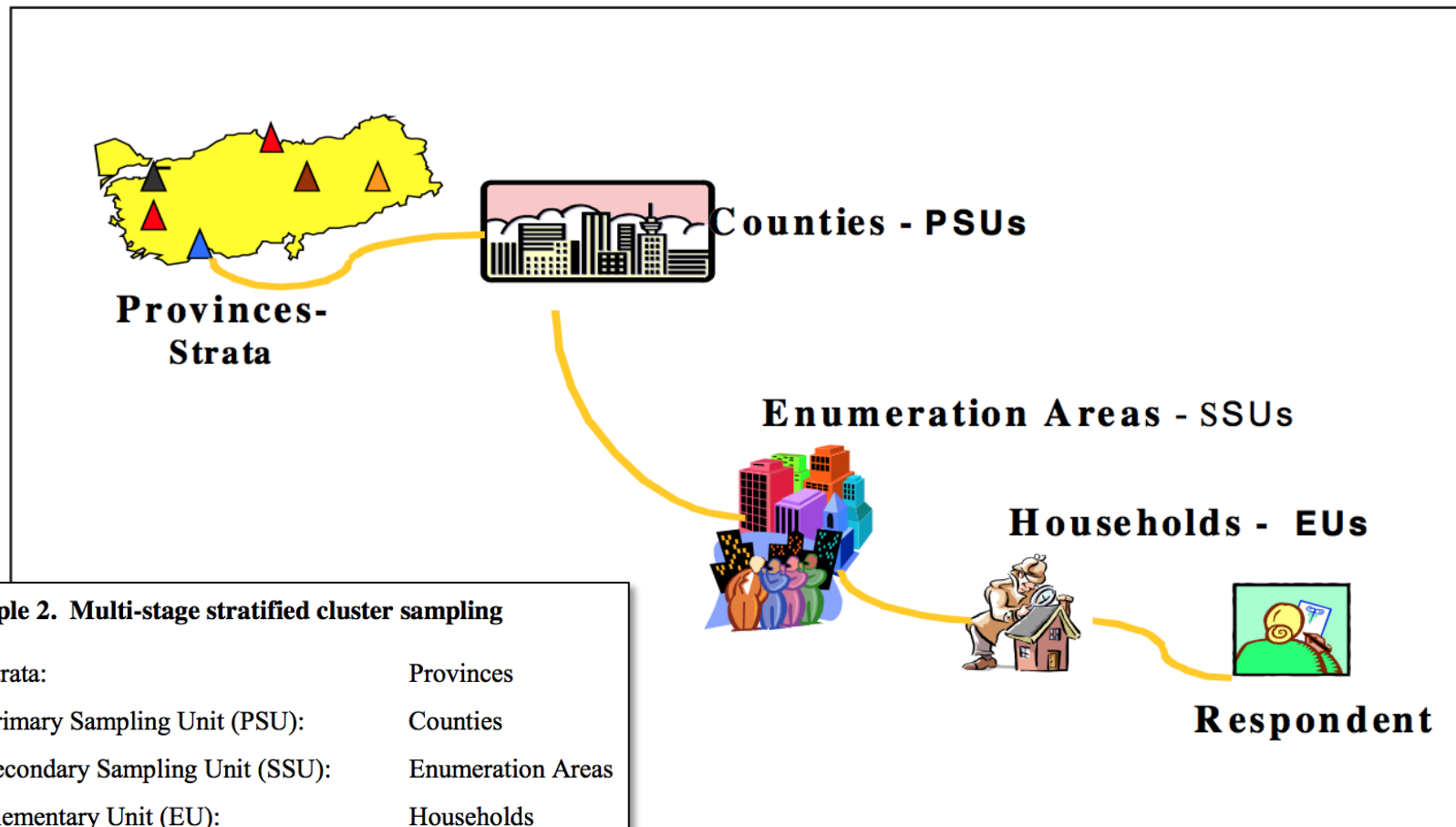
National Survey of Third World Country



- First step: Stratify sample by state/province proportional to population
 - Oversample any state with less than 100 or 200 interviews to allow for state-to-state comparisons
- Second step: Within state/province, stratify by urban and rural
 - Urban/rural stratification used to make sure that all localities are represented
 - As a general rule, locations of 10,000 or more classified urban, otherwise classified rural
- Third step: Select PSUs within state/provinces and by urban/rural location
- Fourth step: Select starting point within each PSU for each interviewer
 - Starting points defined as locations with sufficient public presence to be known by local residents, such as schools, markets, etc.

The World Health Survey Illustration

Figure 1. Multi-stage cluster sampling



Example 2. Multi-stage stratified cluster sampling

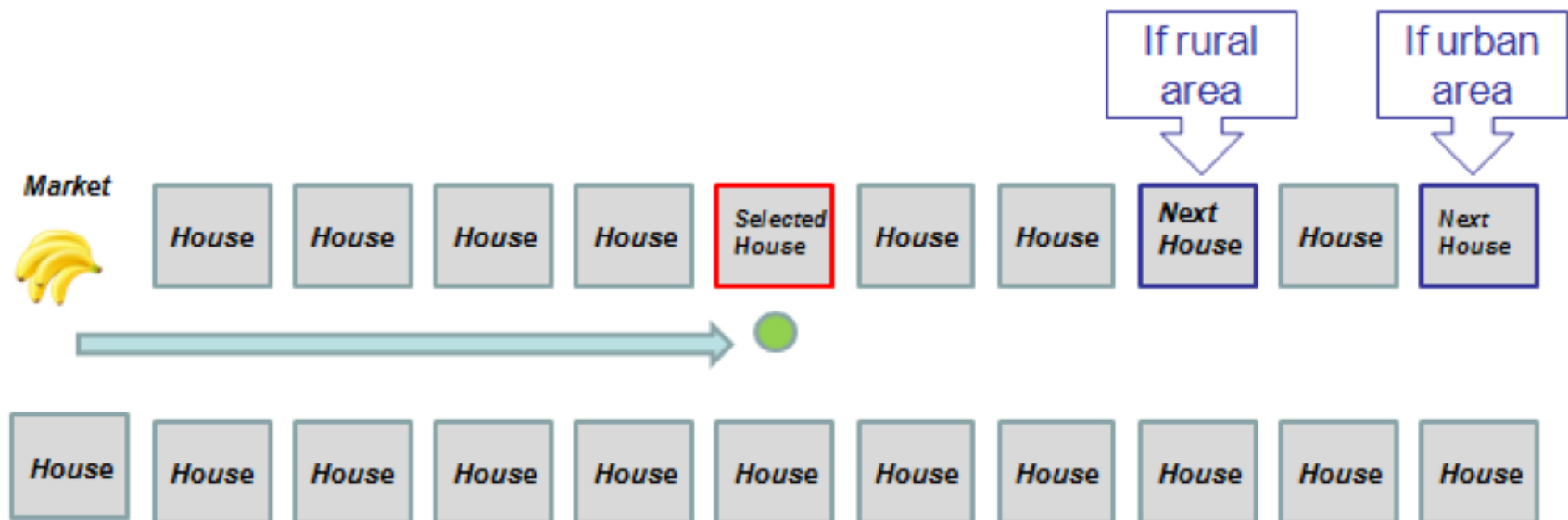
- Strata: Provinces
- Primary Sampling Unit (PSU): Counties
- Secondary Sampling Unit (SSU): Enumeration Areas
- Elementary Unit (EU): Households
- Final Unit: Persons

3/26/13

House Selection Via Systematic Sampling



1. From the starting point, the interviewers walk in different directions (north, south, east, and west).
2. To select the first dwelling, use the day code. To determine the day code, add the digits of day until only one number is remaining (EX: if November 13, $1+3=4$, so day code is 4).
3. Skip the number of dwellings of the day code and select the next dwelling, counting from the left (EX: in this case, skip 4 dwellings, and select the 5th dwelling.)
4. After the first dwelling is selected, each additional dwelling is determined by the skip pattern (every 3rd dwelling in rural areas or every 5th dwelling in urban areas).



Selection of Household in Multi-dwelling Structure



Multi-Household Selection Grid

The grid below is only to be used when one enters a one-level multi-household dwelling.

For an apartment building with multiple floors, one continues the skip interval, counting from the left on the bottom floor and working up.

Interviewer: Mark the number of households in dwelling in the left-hand column. Then calculate the day code, find that number in the top column and follow this column down until it intersects with the row corresponding to the marked number on the left. The number in that box will be the household you select. (Note: If there are more than 7 households, continue the counting by going back to "1" in the left-hand column—i.e., if 10 households, that would be #3 in the left-hand column.)

	Day Code Numbers								
Number of households in dwelling	1	2	3	4	5	6	7	8	9
1	1	1	1	1	1	1	1	1	1
2	1	2	2	2	1	1	1	2	1
3	1	2	2	3	1	1	3	1	2
4	2	3	3	3	1	2	4	3	4
5	2	4	4	5	5	3	1	5	3
6	5	5	5	2	3	4	1	3	2
7	1	2	2	1	2	4	2	6	7

Respondent Selection in Each House



- To select the person to interview within a household:
 - List all adult males and females aged 18 years and above in the household on a **Kish grid**
 - A Kish grid is essentially a table of randomly generated numbers
 - It's a pre-assigned table of random numbers to find the person to be interviewed
 - Alternative is the next-birthday method
 - One respondent is selected using the grid
 - Once the responded is selected, the interview is conducted with only that respondent

Kish Grid (aka Kish Tables) Example

Sequentially work down the list

Household	Kish Table
1	A
2	A
3	B1
4	B2
5	C
6	C
7	D
8	D
9	E1
10	E2
11	F
12	F
13	A
14	A
15	B1
16	B2
17	C
18	C
19	D
20	D
21	E1
22	E2
23	F
24	F
25	A
26	A
27	B1
Etc.....	

Selection table D	
If the number of adults in household is:	Select adult numbered:
1	1
2	2
3	2
4	3
5	4
6 or more	4

Overall Selection Probabilities

Adult numbered	If the number of adults in household is:					
	1	2	3	4	5	6 or more
1	1	1/2	1/3	1/4	1/6	1/6
2		1/2	1/3	1/4	1/6	1/6
3			1/3	1/4	1/4	1/6
4				1/4	1/6	1/6
5					1/4	1/6
6						1/6
7 or more						0

What We Have Covered



- Defined stratified & cluster sampling
- Discussed reasons for stratified and cluster sampling
- Derived estimators for the population
 - Mean
 - Total
 - Percentage
- Discussed examples of complex sampling designs
- Explained the Kish grid