

AS: Sampling methods for social surveys

Gaia Bertarelli, Sant'Anna School of Advanced Studies (Italy)

January 2023

The presentation at a glance

Simple Probability Samples

Estimators

Simple Random Sampling

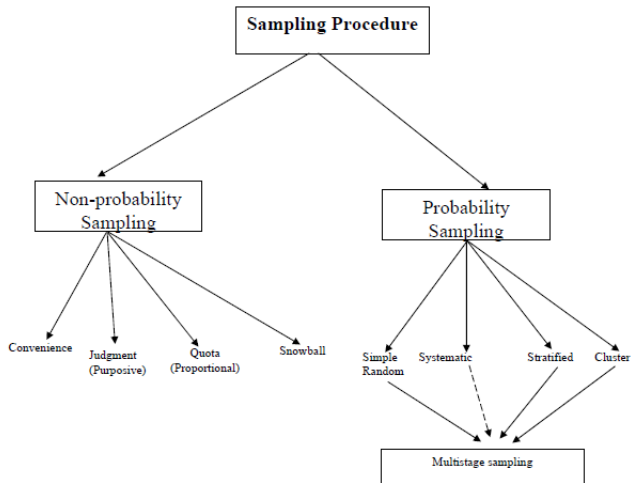
Stratified Sampling

Cluster sampling

Non-probability sampling techniques

Simple Probability Samples

Types of Samples



Simple Probability Sample

- In a **probability sample** each unit in the population has a **known probability of selection**, and a random number table or a randomization mechanism is used to choose the specific units to be included in the sample.
- If a probability sampling design is implemented well, an investigator can use a relatively small sample to make inferences about an arbitrarily large population.

Basic types of Probability Samples - Simple Random Sampling

- Simple Random Sampling (SRS) is the simplest form of probability sample.
 - An SRS of size n is taken when every possible subset of n units in the population has the same probability of being part of the sample.
 - SRSs are at the basis of more complex designs.
 - In taking a random sample the investigator is mixing every member of the population before selecting n units.
 - The investigator does not need to examine every member of the population.

Stratified random sample

- In the **Stratified random sample** the population is divided into subgroups called **strata**.
- Then a **SRS** is selected **from each stratum**.
- The SRSs in the strata are selected independently.
 - The strata are often subgroups of interest to the investigator (to the survey).
 - Elements in the same stratum often tend to be **more similar** than randomly selected elements from the whole population.
 - As a consequence **stratification increases precision**.

Cluster sample

- In a **Cluster sample** observation units in the population are aggregated into larger sampling units called **clusters**.
- Cluster sampling is used when natural groups are present in a population.
- The investigator takes a **SRS of the clusters (e.g. schools)** and **then subsample all or some members of the clusters**. They then randomly select among these clusters to form a sample.
- *Pay attention:* In stratified sampling, individuals are randomly selected from all strata to make up the sample. On the other hand cluster sampling, the sample is formed when all individuals are taken from randomly selected clusters.

- In stratified sampling, there is homogeneity within the group, while in the case of cluster sampling the homogeneity is found between groups.
- Heterogeneity occurs between groups in stratified sampling. In contrast, group members are heterogeneous in cluster sampling.
- When the sampling method adopted by the researcher is stratified, then the categories are imposed by him. Otherwise, categories are groups that already exist in cluster sampling.
- *Stratified sampling aims to improve accuracy and representation. Unlike cluster sampling, the goal of which is to improve cost effectiveness and operational efficiency.*

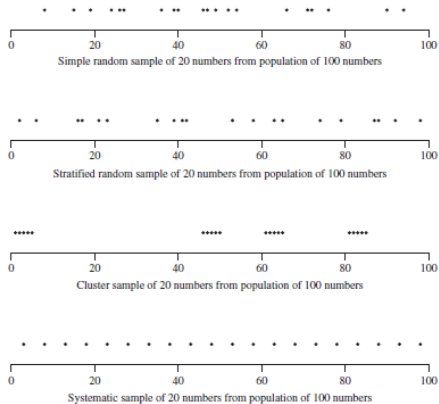
Systematic sample

- In a **Systematic sample** a starting point is chosen from a list of population members using a random number.
- That unit, and every k -th unit thereafter, is chosen in the sample.
- A systematic sample thus consists of units that are equally spaced in the list.

Figure 1: Source: S. Lohr (2019). Sampling Design and Analysis. CRC press. (pg. 27)

FIGURE 2.1

Examples of a simple random sample, stratified random sample, cluster sample, and systematic sample of 20 integers from the population $\{1, 2, \dots, 100\}$.



- We need to be able to list the N units in the finite population.
- The finite population (**Universe**) of N units is denoted by the index set $U = \{1, 2, \dots, N\}$
- Out of this population we can choose different samples, which are subsets of U
- Suppose $U = \{1, 2, 3, 4\}$ we can select 6 samples from this finite population $S_1 = \{1, 2\}$, $S_2 = \{1, 3\}$, $S_3 = \{1, 4\}$, $S_4 = \{2, 3\}$, $S_5 = \{2, 4\}$, $S_6 = \{3, 4\}$.

- Each possible sample S from the population has a known probability $P(S)$ of being selected (and the possible probabilities sum to 1).
- In a probability sample, since each possible sample has a known probability of being the chosen sample, each unit in the population has a known probability of appearing in our selected sample. It is called **Inclusion probability**

$$\pi_i = P(\text{unit } i \text{ in sample})$$

- The **Sampling Weight**, for any sampling design, is the reciprocal of the inclusion probability

$$w_i = \frac{1}{\pi_i}$$

The sampling weight of unit i is the number of population units represented by unit i .

Estimators

Estimators

- An **Estimator** is a Statistic (a random variable (rv)) whose calculated value is used to estimate a population parameter θ .
- An **Estimate** A particular realization of an estimator $\hat{\theta}$.
- **Type of Estimators:**
 1. **point estimate:** single number that can be regarded as the most plausible value of θ .
 2. **interval estimate:** a range of numbers, called a confidence interval indicating, can be regarded as likely containing the true value of θ
- In the Frequentist world view parameters are fixed, statistics are rv and vary from sample to sample (i.e., have an associated sampling distribution).
- In theory, there are many potential estimators for a population parameter.

Estimators (ii)

- Good Estimators Are:

1. **Consistent**: as the sample size increases $\hat{\theta}$ gets closer to θ .
2. **Unbiased**: $E(\hat{\theta}) = \theta$ and we call **Bias** the quantity:

$$Bias(\hat{\theta}) = E(\hat{\theta}) - \theta$$

3. **Precise**: Sampling distribution of $\hat{\theta}$ should have small variance

$$V(\hat{\theta}) = E[(\hat{\theta} - E(\hat{\theta}))^2]$$

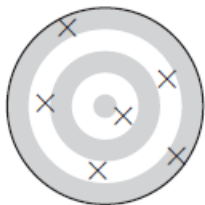
4. **Accurate**: Mean Squared Error is small

$$MSE(\hat{\theta}) = V(\hat{\theta}) + (Bias(\hat{\theta}))^2$$

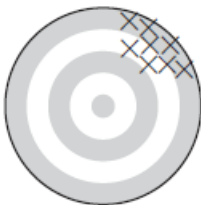
(we have to work with MSE because sometimes we work with biased estimators)

Properties of Estimators

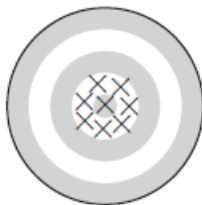
Figure 2: A: unbiased; B: precised but not unbiased; C: accurate



Archer A



Archer B



Archer C

Estimators (iii)

- The finite population $U = \{1, 2, \dots, N\}$ has as measured values $\{y_1, y_2, \dots, y_N\}$.
- It is possible to select a sample S of n units from U using the probabilities of selection who defined a sampling design.
- y_i fixed but unknown unless that unit i appears in the selected sample S .
- Without any other statistical assumptions, the only information we have about the set of y_i 's in the population in the set $\{y_i : i \in S\}$.

Population Quantity

- Population Total:

$$t = \sum_{i=1}^N y_i$$

- Population Mean:

$$\bar{y}_U = \frac{t}{N}$$

- Population Variance:

$$S^2 = \frac{1}{N-1} \sum_{i=1}^N (y_i - \bar{y}_U)^2$$

- Population Proportion:

$$p = \frac{\text{Number of units with the characteristic of interest in the population}}{N}$$

Simple Random Sampling

Simple Random Sampling

- **Two** types of Simple random sampling:
 1. Simple random sampling with replacement (**SRSWR**): the same unit may be included more than one in the sample.
 2. Simple random sampling without replacement (**SRS**): all units in the sample are distinct.
- In SRSWR a sample of size n from a population of N units can be seen as n independent sample of size 1. One unit is randomly selected from the population to be ALWAYS the first sampled unit with probability $\frac{1}{N}$.

- A SRS of size n is selected so that **every possible subset** of n distinct units in the population has the same probability of being selected as the sample.
- There are $\binom{N}{n}$ possible samples.
- The probability of selecting any individual sample S of n units is

$$P(S) = \frac{1}{\binom{N}{n}} = \frac{n!(N-n)!}{N!}$$

- Remember:
 - $k! = k(k-1)(k-2)\dots 1$
 - $0! = 1$
 - $\binom{N}{n} = \frac{N!}{n!(N-n)!}$

Sampling weights is SRS

- In SRS each unit has inclusion probability $\pi_i = \frac{n}{N}$
- As a consequence each unit has the same weight

$$w_i = \frac{N}{n}$$

.

- Every unit in the sample represents itself plus $\frac{N}{n} - 1$ unsampled units in the population.

Estimators SRS

| Population Quantity | Estimator | Standard Error of Estimator |
|--|---|---|
| Population total, $t = \sum_{i=1}^N y_i$ | $\hat{t} = \sum_{i \in \mathcal{S}} w_i y_i = N\bar{y}$ | $N\sqrt{\left(1 - \frac{n}{N}\right)\frac{s^2}{n}}$ |
| Population mean, $\bar{y}_U = \frac{t}{N}$ | $\frac{\hat{t}}{N} = \frac{\sum_{i \in \mathcal{S}} w_i y_i}{\sum_{i \in \mathcal{S}} w_i} = \bar{y}$ | $\sqrt{\left(1 - \frac{n}{N}\right)\frac{s^2}{n}}$ |
| Population proportion, p | \hat{p} | $\sqrt{\left(1 - \frac{n}{N}\right)\frac{\hat{p}(1 - \hat{p})}{n - 1}}$ |

Finite Population Correction

- $\left(1 - \frac{n}{N}\right)$ it is called **Finite Population Correction**(fpc): this correction is made because, if we have a small population, the greater our sampling fraction $\frac{n}{N}$ is, the more information we have about the population and thus the smaller is the variance.
- For most samples coming from large populations fpc is generally equal to 1 (when n is very small relative to N , the finite population correction is almost equal to one).
- Variance is estimated from the sample, but the fpc it is used to assess the error in estimation, which is due to the fact that not all data from the finite population are observed.
- Practically it is used when you sample without replacement from more than 5% of a finite population.

Finite Population Correction (ii)

- Some formulas used to compute standard errors are based on the idea that (1) samples are selected from an infinite population or (2) samples are selected with replacement.
- This does not present much of a problem when the sample size (n) is small relative to the population size (N); that is, when the sample is less than 5% of the population.
- If you are using a SRSWR you do not have to use the fpc.

Confidence Intervals

- When reporting the results of a survey it is necessary to give an idea of how accurate the estimates are.
- **Confidence Intervals** (CI) are used to indicate the accuracy of an estimate.
- We do not know the values of the statistics from all possible samples so we can not compute the exact confidence interval
→ Asymptotic results. Our population is supposed to be part of a **superpopulation**.
- For big enough sample the CI is given by

$$\text{estimate} \pm z_{\alpha/2} SE(\text{estimate})$$

- The margin of error of an estimate is the half-width of a confidence interval, i.e. $z_{\alpha/2} SE(\text{estimate})$

Sample Size Estimation

- **Specify the tolerable error:** the precision needed can be defined as

$$P(|\bar{y} - \bar{y}_U| \leq e) = 1 - \alpha$$

- e is the **margin of error**.
- e and α are fixed by the researcher.
- If the relative precision must be fixed and it is preferred to check the CV with respect to the absolute error

$$P\left(\left|\frac{\bar{y} - \bar{y}_U}{\bar{y}_U}\right| \leq r\right) = 1 - \alpha$$

Sample size estimation (ii)

- SRSWR

$$n_0 = \left(\frac{z_{\alpha/2} S}{e} \right)^2$$

- SRS (absolute error)

$$n = \frac{n_0}{1 + \frac{n_0}{N}} = \frac{z_{\alpha/2}^2 S^2}{e^2 + \frac{z_{\alpha/2}^2 S^2}{N}}$$

- SRS (relative precision)

$$n = \frac{n_0}{1 + \frac{n_0}{N}} = \frac{z_{\alpha/2}^2 S^2}{r \bar{y}_U^2 + \frac{z_{\alpha/2}^2 S^2}{N}}$$

Stratified Sampling

Stratified Sampling

- If the variable of interest takes on different mean values in different subpopulations, we may be able to obtain more precise estimates of population quantities by taking a **stratified** random sample.
- We divide the population into H subpopulations. Each subpopulation is called a Stratum. Strata do not overlap.
- We draw an independent probability sample from each stratum, **then pool the information to obtain overall population estimates.**

Why Stratified Sampling?

- Protected from the possibility of very bad sample.
- Known precision for subgroups of the population.
- More convenient to administer.
- Stratified sampling often gives more precise (lower variance) estimates of population means and totals.

Notation

- H : number of strata.
- N : population size.
- N_h : population size in stratum h .
- $N = N_1 + N_2 + N_3 + \cdots + N_H$
- In the stratified random sampling we independently take an SRS from each stratum: n_h units are randomly selected from the N_h population units in stratum h .
- The total sample size is $n = n_1 + n_2 + n_3 + \cdots + n_H$

Notation for Stratification: The population quantities are:

y_{hj} = value of j th unit in stratum h

$$t_h = \sum_{j=1}^{N_h} y_{hj} = \text{population total in stratum } h$$

$$t = \sum_{h=1}^H t_h = \text{population total}$$

$$\bar{y}_{hU} = \frac{\sum_{j=1}^{N_h} y_{hj}}{N_h} = \text{population mean in stratum } h$$

$$\bar{y}_U = \frac{t}{N} = \frac{\sum_{h=1}^H \sum_{j=1}^{N_h} y_{hj}}{N} = \text{overall population mean}$$

$$S_h^2 = \sum_{j=1}^{N_h} \frac{(y_{hj} - \bar{y}_{hU})^2}{N_h - 1} = \text{population variance in stratum } h$$

Estimates in each Stratum

- Mean:

$$\bar{y}_h = \frac{1}{n_h} \sum_{j \in S_h} y_{hj}$$

- Total:

$$\hat{t}_h = \frac{N_h}{n_h} \sum_{j \in S_h} y_{hj} = N_h \bar{y}_h$$

- Sampling variance:

$$s_h^2 = \sum_{j \in S_h} \frac{(y_{hj} - \bar{y}_h)^2}{(n_h - 1)}$$

Estimates in Stratified Sampling

- Population total estimator:

$$\hat{t}_{str} = \sum_{h=1}^H \hat{t}_h = \sum_{h=1}^H N_h \bar{y}_h = \sum_{h=1}^H \sum_{j \in S_h} w_{hj} y_{hj}$$

- Population mean estimator:

$$\bar{y}_{str} = \frac{\hat{t}_{str}}{N} = \sum_{h=1}^H \frac{N_h}{N} \bar{y}_h = \frac{\sum_{h=1}^H \sum_{j \in S_h} w_{hj} y_{hj}}{\sum_{h=1}^H \sum_{j \in S_h} w_{hj}}$$

- $\frac{N_h}{N}$: the proportion of the population units in stratum h .
- Population proportion estimator:

$$\hat{p}_{str} = \sum_{h=1}^H \frac{N_h}{N} \hat{p}_h$$

- The properties of these estimators follow the properties of SRS estimators (unbiased, definition of Variance, Standard Errors and Confidence intervals)

- **Unbiasedness.** \bar{y}_{str} and \hat{t}_{str} are unbiased estimators of \bar{y}_U and t . An SRS is taken in each stratum, so (2.30) implies that $E[\bar{y}_h] = \bar{y}_{hU}$ and consequently

$$E\left[\sum_{h=1}^H \frac{N_h}{N} \bar{y}_h\right] = \sum_{h=1}^H \frac{N_h}{N} E[\bar{y}_h] = \sum_{h=1}^H \frac{N_h}{N} \bar{y}_{hU} = \bar{y}_U.$$

- **Variance of the estimators.** Since we are sampling independently from the strata, and we know $V(\hat{t}_h)$ from the SRS theory, the properties of expected value in Section A.2 and (2.16) imply that

$$V(\hat{t}_{\text{str}}) = \sum_{h=1}^H V(\hat{t}_h) = \sum_{h=1}^H \left(1 - \frac{n_h}{N_h}\right) N_h^2 \frac{S_h^2}{n_h}. \quad (3.3)$$

- **Standard errors for stratified samples.** We can obtain an unbiased estimator of $V(\hat{t}_{\text{str}})$ by substituting the sample estimators s_h^2 for the population parameters S_h^2 . Note that in order to estimate the variances, we need to sample at least two units from each stratum.

$$\hat{V}(\hat{t}_{\text{str}}) = \sum_{h=1}^H \left(1 - \frac{n_h}{N_h}\right) N_h^2 \frac{s_h^2}{n_h} \quad (3.4)$$

$$\hat{V}(\bar{y}_{\text{str}}) = \frac{1}{N^2} \hat{V}(\hat{t}_{\text{str}}) = \sum_{h=1}^H \left(1 - \frac{n_h}{N_h}\right) \left(\frac{N_h}{N}\right)^2 \frac{s_h^2}{n_h}. \quad (3.5)$$

As always, the standard error of an estimator is the square root of the estimated variance: $\text{SE}(\bar{y}_{\text{str}}) = \sqrt{\hat{V}(\bar{y}_{\text{str}})}$.

- **Confidence intervals for stratified samples.** If either (1) the sample sizes within each stratum are large, or (2) the sampling design has a large number of strata, an approximate $100(1 - \alpha)\%$ confidence interval (CI) for the population mean \bar{y}_U is

$$\bar{y}_{\text{str}} \pm z_{\alpha/2} \text{ SE } (\bar{y}_{\text{str}}).$$

The central limit theorem used for constructing this CI is stated in Krewski and Rao (1981). Some survey software packages use the percentile of a t distribution with $n - H$ degrees of freedom (df) rather than the percentile of the normal distribution.

Sampling weights in Stratified Sampling

- In stratified sampling it is possible to have different inclusion probabilities in different strata \rightarrow weights may be unequal in different strata.
- The sampling weight is the number of units in the population represented by the sample member y_{hj} .
- the sampling weight is

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

- Stratified sampling has three major design issues:
 1. Defining the strata.
 2. Choosing the total sample size.
 3. Allocating the observations to the defined strata.

Allocating Observations to Strata

- Now we assume that the strata have already been fixed and we study methods of allocating observations to the strata.
- Different types of allocation:
 1. Equal Allocation ($n_1 = n_2 = n_3 = \dots = n_H$)
 2. Proportional Allocation
 3. Optimal Allocation
 4. Allocation for specified precision within strata
- then it is necessary to define the sample size.

Proportional Allocation

- **Goal: obtain a sample that reflects the population with respect to the stratification variable.** The final sample is a miniature version of the population \rightarrow the number of sampled units in each stratum is **proportional to the size of the stratum** $\rightarrow \frac{n_h}{n} = \frac{N_h}{N} \rightarrow n_h = \left(n \frac{N_h}{N}\right)$
- The inclusion probability $\pi_{hj} = \frac{n_h}{N_h} = \frac{n}{N}$ is **the same for all strata** as in SRS but in SRS more *bad* samples are possible.
- Each unit in the sample represent the same number of units in the population.

Proportional Allocation (ii)

- The stratified sampling estimator of the population mean (e.g.) is simply the average of all the observation.
- When the strata are large enough, the variance of the stratified estimator (e.g. mean) under proportional allocation is usually at most as large as the variance of the estimator under the SRS with the same number of observation.
- The more unequal the stratum estimates are, the more precision the researcher will gain by using proportional allocation.
- If the variances are more or less equal across the strata proportional allocation is the best option to increase precision.

Optimal Allocation

- **Idea: to sample more in larger strata and where variability of y is greater.**
- Optimal allocation works well for sampling units such as corporations, cities, hospitals which vary a lot in size.
- General optimal allocation works well when some strata are much more expensive to sample than others.
- When calculating n_h^{opt} we must take S_h as known. If they are not, it is necessary to approximate.
- It can happen that $n_h^{opt} \geq N_h$ for some h . In this case we can decide to census the stratum, and apply the optimal allocation algorithm to the remaining strata.

Optimal Allocation (ii)

- Optimal Allocation

$$n_h = \left(\frac{\frac{N_h S_h}{\sqrt{c_h}}}{\sum_{l=1}^H \frac{N_l S_l}{\sqrt{c_l}}} \right) n$$

where c_h represent the cost of taking an observation in stratum h .

- Sampling heavily within a stratum if:
 1. The stratum accounts for a large part of the population;
 2. The variance within the stratum is large: more sample size to compensate for the heterogeneity;
 3. Sampling in the stratum is not expensive.

Neyman allocation

- Neyman allocation is a special case of optimal allocation.
- It is used when the cost in the strata (but not the variances) are approx. equal.
- $n_h \propto N_h S_h$
- The formula is.

$$n_h = \left(\frac{N_h S_h}{\sum_{l=1}^H N_l S_l} \right) n$$

- If S_h are well allocated, Neyman allocation will give an estimator with smaller variance than the proportional allocation.

Lessons learned

- Stratified sampling is more efficient of simple random sampling.
We can allocate the sample efficiently if we can locate the strata where the variability of the study variable is concentrated.
- When we have positive **asymmetric populations** (many small values, few large values) is worthwhile try to intensively sample the right tail of the distribution.
- The optimal allocation is based on the consideration of a single variable as target.

Sample sizes

- If you are interested in the estimates in the strata: sample size coming from SSR.
- The different methods of allocating observations to strata give the relative sample size $\frac{n_h}{n}$.
- After the construction of strata and after that observations are allocated to strata:

$$V(\bar{y}_{\text{str}}) \leq \frac{1}{n} \sum_{h=1}^H \frac{n}{n_h} \left(\frac{N_h}{N} \right)^2 S_h^2 = \frac{v}{n},$$

where $v = \sum_{h=1}^H (n/n_h) (N_h/N)^2 S_h^2$. Thus, if the fpcs can be ignored and if the normal approximation is valid, an approximate 95% CI for the population mean will be $\bar{y}_{\text{str}} \pm z_{\alpha/2} \sqrt{v/n}$. Set $n = z_{\alpha/2}^2 v / e^2$ to achieve a desired margin of error e .

Please read Section 3.5 *Defining Strata* in your textbook Lohr (2019).

Cluster sampling

Cluster sampling

Suppose we want to find out how many bicycles are owned by residents in a community of 10,000 households. We could take a simple random sample (SRS) of 400 households, or we could divide the community into blocks of about 20 households each and sample every household (or subsample some of the households) in each of 20 blocks selected at random from the 500 blocks in the community. The latter plan is an example of cluster sampling. The blocks are the **primary sampling units** (psus), or **clusters**. (In this chapter, we use the terms cluster and psu interchangeably.) The households are the **secondary sampling units** (ssus); often the ssus are the elements in the population.

Cluster sampling

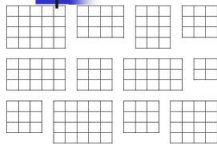
- In cluster sampling individual elements of the population are allowed in the sample only if they belong to a cluster (psu) that is included in the sample.
- The sampling unit (psu) is not the same as the observation unit (ssu)

- **Cluster sampling:** A probability sampling design in which observations are grouped into clusters (psu). A probability sample of psus is selected from the population of psus.
- **Primary sampling unit (psu):** the unit that is sampled from the population.
- **Secondary sampling unit:** a subunit that is subsampled from the selected psu.
- **One-stage cluster sampling:** a cluster sampling design in which all ssus in selected psus are observed.
- **Two-stage cluster sampling:** a cluster sampling design in which the ssus in selected psus are subsampled.

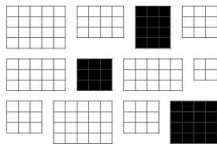
Sample of 40 elements



1-stage CS



Take an SRS of clusters; observe all elements within the clusters in the sample:

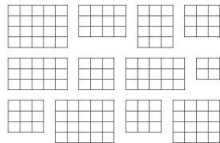


A block of cells is a cluster

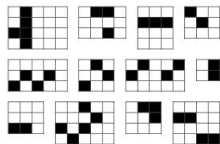
SU is a cluster

Don't sample from every cluster

ST



Take an SRS from every stratum:



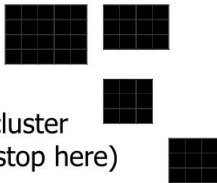
A block of cells is a stratum

SU is an element (or OU)

Sample from every stratum

1-stage vs. 2-stage cluster sampling

Sample all SSUs in sampled PSUs:

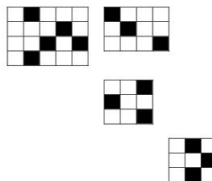


1-stage cluster
sample (stop here)

OR

Stage **1** of 2-stage
cluster sample
(select PSUs)

Take an SRS of m SSUs in sampled PSU i :



Stage **2** of 2-stage
cluster sample
(select SSUs w/in PSUs)

7

Notation Cluster Sampling

- The universe U is the population on N psus.
- N here is the number of psus not of observation units.
- S is the sample of psus.
- S_i is the sample of ssus chosen from the i^{th} psu.
- y_{ij} measurement for j^{th} element in i^{th} psu

Notation Cluster Sampling (ii)

psu Level—Population Quantities

N = number of psus in the population

M_i = number of ssus in psu i

$M_0 = \sum_{i=1}^N M_i$ = total number of ssus in the population

$t_i = \sum_{j=1}^{M_i} y_{ij}$ = total in psu i

$t = \sum_{i=1}^N t_i = \sum_{i=1}^N \sum_{j=1}^{M_i} y_{ij}$ = population total

$S_t^2 = \frac{1}{N-1} \sum_{i=1}^N \left(t_i - \frac{t}{N} \right)^2$ = population variance of the psu totals

Notation Cluster Sampling (iii)

ssu Level—Population Quantities

$$\bar{y}_U = \sum_{i=1}^N \sum_{j=1}^{M_i} \frac{y_{ij}}{M_0} = \text{population mean}$$

$$\bar{y}_{iU} = \sum_{j=1}^{M_i} \frac{y_{ij}}{M_i} = \frac{t_i}{M_i} = \text{population mean in psu } i$$

$$S^2 = \sum_{i=1}^N \sum_{j=1}^{M_i} \frac{(y_{ij} - \bar{y}_U)^2}{M_0 - 1} = \text{population variance (per ssu)}$$

$$S_i^2 = \sum_{j=1}^{M_i} \frac{(y_{ij} - \bar{y}_{iU})^2}{M_i - 1} = \text{population variance within psu } i$$

Notation Cluster Sampling (iv)

Sample Quantities

n = number of psus in the sample

m_i = number of ssus in the sample from psu i

$\bar{y}_i = \sum_{j \in S_i} \frac{y_{ij}}{m_i}$ = sample mean (per ssu) for psu i

$\hat{t}_i = \sum_{j \in S_i} \frac{M_i}{m_i} y_{ij}$ = estimated total for psu i

$\hat{t}_{\text{unb}} = \sum_{i \in S} \frac{N}{n} \hat{t}_i$ = unbiased estimator of population total

$$s_t^2 = \frac{1}{n-1} \sum_{i \in S} \left(\hat{t}_i - \frac{\hat{t}_{\text{unb}}}{N} \right)^2$$

Notation Cluster Sampling (v)

$$s_i^2 = \sum_{j \in \mathcal{S}_i} \frac{(y_{ij} - \bar{y}_i)^2}{m_i - 1} = \text{sample variance within psu } i$$

w_{ij} = sampling weight for ssu j in psu i

One-stage cluster sampling

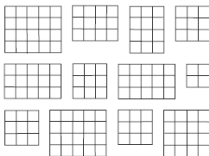
- One-stage cluster sampling is used in many surveys in which **the cost of sampling ssus is small compared to the cost of sampling psus**.
- In the simplest design, we take an SRS of n psus from the population and measure our variable of interest in **every element in the sampled psus** $\rightarrow M_i = m_i$
- One-stage cluster sampling with an SRS of psus produces a **self-weighted sample**: a sample in which all probabilities of inclusion π_i are equal, so that all sampling weights w_i are the same.
- $w_{ij} = \frac{N}{n}$

Two-stage cluster sampling

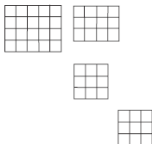
- It may be expensive to measure ssus relative to the cost of sampling psus.
- Two-stage cluster sampling:
 1. Select an SRS S OF n psus from the population of N psus.
 2. Select an SRS of ssus from each selected psus. The SRS of m_i elements from the i^{th} psu denote S_i .

One-Stage

Population of N psu's:



Take an SRS of n psu's:

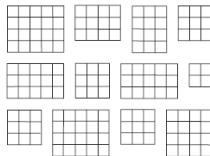


Sample all ssu's in sampled psu's:

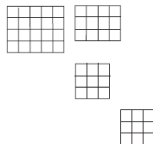


Two-Stage

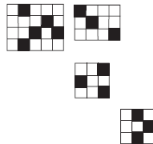
Population of N psu's:



Take an SRS of n psu's:



Take an SRS of m_i ssu's in sampled psu i :



Inclusion probability and weight



$$\begin{aligned}\pi_{ij} &= P(j^{th} \text{ ssu in } i^{th} \text{ psu is selected}) = \\ &= P(i^{th} \text{ psu selected}) \times P(j^{th} \text{ ssu selected} | i^{th} \text{ psu selected}) = \\ &= \frac{n}{N} \frac{m_i}{M_i}\end{aligned}$$



$$w_{ij} = \frac{1}{\pi_{ij}} = \frac{NM_i}{nm_i}$$

- ssu j in psu i represents $\frac{1}{\pi_{ij}} = \frac{NM_i}{nm_i}$ ssus in the population: itself and $\frac{1}{\pi_{ij}} = \frac{NM_i}{nm_i} - 1$ units which are not sampled.

Cluster sampling: issues

1. What overall precision is needed? → common in all the survey designs
 2. What size should psus be? → often natural.
 3. How many ssus should be sampled in each psu selected for the sample?
 4. How many psus should be sampled?
- **Intraclass correlation coefficient (ICC)**: the pearson correlation coefficient of all pairs of unit within the same cluster.

Choosing Subsampling sample size

- Goal: Most information possible for the least cost and inconvenience.
- Suppose that we want to conduct a two-stage cluster survey when all the psus have the same number, equal to M , of ssus.
- Consider the simple total cost function

$$C = c_1n + c_2nm$$

where

- c_1 cost per psu (not including c_2)
- c_2 cost of measuring each ssu

Choosing Subsampling sample size (ii)

- The optimal m and n are the values that minimize the variance for the fixed total cost C .
- We have

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

and

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

- R_a^2 is a measure of homogeneity defined as

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

and MSW is the pooled value of the within-cluster variances.

Although we discussed only designs where all M_i 's are equal, we can use these methods with unequal M_i 's as well: just substitute \bar{M} for M in the above work, and decide the average subsample size \bar{m} to take. Then either take \bar{m} observations in every cluster, or allocate observations so that

$$\frac{m_i}{M_i} = \text{constant}.$$

As long as the M_i 's do not vary too much, this should produce a reasonable design.

Number of psus (choosing the sample size)

After the psu size is determined and the subsampling fraction set, we then look at the number of psus to sample, n . Like any survey design, design of a cluster sample is an iterative process: (1) Determine a desired precision, (2) choose the psu and subsample sizes, (3) conjecture the variance that will be achieved with that design, (4) set n to achieve the precision, and (5) iterate (adding stratification and auxiliary variables to use in ratio estimation) until the cost of the survey is within your budget.

If clusters are of equal size and we ignore the psu-level fpc, (5.30) implies that

$$V(\hat{y}_{\text{unb}}) \leq \frac{1}{n} \left[\frac{\text{MSB}}{M} + \left(1 - \frac{m}{M}\right) \frac{\text{MSW}}{m} \right] = \frac{1}{n} v.$$

Number of psus (choosing the sample size)

An approximate $100(1 - \alpha)\%$ CI will be

$$\hat{y}_{\text{unb}} \pm z_{\alpha/2} \sqrt{\frac{1}{n} v}.$$

Thus, to achieve a desired CI half-width e , set $n = z_{\alpha/2}^2 v / e^2$.

Non-probability sampling

Non-probability sampling

- Non-probability techniques, relying on the judgment of the researcher or on accident, cannot generally be used to make generalizations about the whole population.
- Non-probability sampling represents a group of sampling techniques that help researchers to select units from a population that they are interested in studying.
- A core characteristic of non-probability sampling techniques is that samples are selected based on the **subjective judgement** of the researcher, rather than random selection (i.e., probabilistic methods)

Principles of non-probability sampling

- Non-probability sampling represents a valuable group of sampling techniques that can be used in research that follows
 - qualitative,
 - mixed methods,
 - and even quantitative research designs.
- Researchers following a quantitative research design often feel that they are forced to use non-probability sampling techniques because of some inability to use probability sampling (e.g., the lack of access to a list of the population being studied). However, this is not the case for researchers following a **qualitative research design**.

Theoretical reasons

- Unlike probability sampling, **the goal is not** to achieve objectivity in the selection of samples, or necessarily attempt **to make generalisations** (i.e., statistical inferences) from the sample being studied to the wider population of interest.
- Instead, researchers following a qualitative research design tend to be interested in the **intricacies** of the sample being studied. Whilst making generalisations from the sample to the population under study may be desirable, it is more often a secondary consideration.
- Even whether this is desired, there are additional problems of bias and transferability (or validity).

Practical reasons

- Non-probability sampling is often used because the procedures used to select units for inclusion in a sample are much easier, quicker and cheaper when compared with probability sampling.
- To sample hidden or hard-to-reach population where a list of the population simply does not exist.
- Non-probability sampling can also be particularly useful in exploratory research where the aim is to find out if a problem or issue even exists in a quick and inexpensive way.

Main types of non-probability sampling

- Quota sampling
- Convenience sampling
- Purposive sampling
- Self-selection sampling
- Snowball sampling

Basic ideas: Quota sampling

- Quota sampling is a sampling methodology wherein data is collected from a homogeneous group.
- It involves a two-step process where two variables can be used to filter information from the population.
- It can easily be administered and helps in quick comparison.
- Many persons confuse quota sampling with stratified sampling. In quota sampling, quota classes are formed that serve the role of strata, but **the survey taker uses a nonprobability sampling method such as convenience sampling to reach the desired sample size** in each quota class.

Basic ideas: convenience sampling

- A convenience sample is simply one where the units that are selected for inclusion in the sample are the **easiest to access**.

Basic ideas: purposive sampling

- Purposive sampling, also known as judgmental, selective, or subjective sampling, is a form of non-probability sampling in which researchers **rely on their own judgment** when choosing members of the population to participate in their surveys.
- This survey sampling method requires researchers to have **prior knowledge** about the purpose of their studies so that they can properly choose and approach eligible participants for surveys

Purposive vs Convenience sampling

- The terms purposive sampling and convenience sampling are often used interchangeably, but they do not mean the same thing.
- Convenience sampling is when researchers leverage individuals that can be identified and approached with as little effort as possible. These are often individuals that are geographically close to the researchers or those who have previously completed an online survey.
- Purposive sampling is when researchers thoroughly think through how they will establish a sample population, even if it is not statistically representative of the greater population at hand. As the name suggests, researchers went to this community on purpose because they think that these individuals fit the profile of the people that they need to reach.

Purposive vs Convenience sampling

- While the findings from purposive sampling do not always have to be statistically representative of the greater population of interest, they are qualitatively generalizable.
- The more prior information that researchers have about their particular communities of interest, the better the sample that they're going to select.

Basic ideas: Self-selection sampling

- Self-selection sampling is appropriate when we want to allow units or cases, whether individuals or organisations, to choose to take part in research on their own accord.
- The key component is that research subjects (or organisations) volunteer to take part in the research rather than being approached by the researcher directly.
- A sample is self-selected when the inclusion or exclusion of sampling units is determined by whether the units themselves agree or decline to participate in the sample, either explicitly or implicitly.

Basic ideas: Snowball sampling

- Snowball sampling is particularly appropriate when the population you are interested in is **hidden and/or hard-to-reach** (drug addicts, homeless people, individuals with AIDS/HIV, prostitutes ...)
- This is a sampling technique, in which **existing subjects provide referrals to recruit samples** required for a research study.
- This sampling method involves a **primary data source nominating other potential data sources** that will be able to participate in the research studies.
- Snowball sampling consists of two steps:
 1. Identify potential subjects in the population. Often, only one or two subjects can be found initially.
 2. Ask those subjects to recruit other people (and then ask those people to recruit. Participants should be made aware that they do not have to provide any other names.

Gaia Bertarelli

Department of Economics and Management in the Era of
Data Science



Department
of Excellence
2018 - 2022

EMbeDS

Economics and Management
in the era of Data Science



Sant'Anna
Scuola Universitaria Superiore Pisa



gaia.bertarelli@santannapisa.it