

# SAMPLING AND ESTIMATION

## PART 2: COMPLEX SAMPLING DESIGNS

Stefan Zins<sup>1</sup>, Matthias Sand<sup>2</sup>, and Jan-Philipp Kolb<sup>3</sup>

February 5, 2016

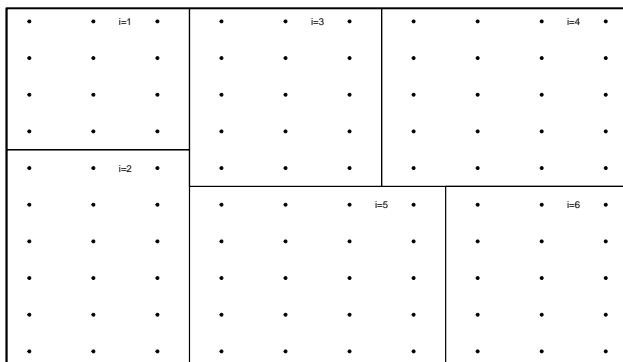
---

<sup>1</sup>Stefan.Zins@gesis.org

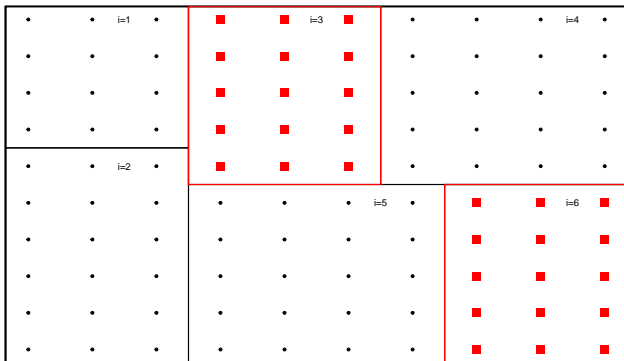
<sup>2</sup>Matthias.Sand@gesis.org

<sup>3</sup>Jan-Philipp.Kolb@gesis.org

A Population of 100 elements is clustered into  $N_l = 6$  cluster



A Population of 100 elements is clustered into  $N_l = 6$  cluster and  $n_l = 2$  clusters are selected from the population.



Sampling elementary units is often not feasible (e.g. persons or businesses). Maybe there is no uniform sampling frame available to select them from, or it would be costly to do, because the selected elements would scatter too much over a certain area and travel costs of interviewers would be too high.

Sampling elementary units is often not feasible (e.g. persons or businesses). Maybe there is no uniform sampling frame available to select them from, or it would be costly to do, because the selected elements would scatter too much over the a certain area and travel costs of interviewers would be too high.

Thus, it is very common to select clusters, so called *primary sampling units* (PSU's) that are populated by *secondary sampling units* (SSU's).

Sampling elementary units is often not feasible (e.g. persons or businesses). Maybe there is no uniform sampling frame available to select them from, or it would be costly to do, because the selected elements would scatter too much over the a certain area and travel costs of interviewers would be too high.

Thus, it is very common to select clusters, so called *primary sampling units* (PSU's) that are populated by *secondary sampling units* (SSU's).

Cluster sampling makes it still possible to obtain unbiased estimates but it can have a big influence on the variance.

$y_{ki}$	=	value of the variable of interest for the $k$ -th SSU in the $i$ -th PSU
$N_l$	=	number of PSU's in the population
$N_i$	=	number of SSU's in the $i$ -th PSU
$N$	=	total number of SSU's in the Population
$\mathcal{U}$	=	set of SSU's in the population
$\mathcal{U}_l$	=	set of PSU's in the population
$\mathcal{U}_i$	=	set of SSU's in the $i$ -th PSU
$n_l$	=	number of PSU's in the sample
$n_i$	=	number of SSU's in the sample from the $i$ -th PSU
$\Delta_l$	=	sample of PSU's
$\Delta_i$	=	sample SSU's from the $i$ -th PSU
$p_l(.)$	=	sampling design of the PSU's

All SSU's in the sampled PSU's are surveyed, thus

$$\tau_i = \sum_{k \in \mathcal{U}_i} y_{ki}$$

is known of all selected PSU's. An unbiased estimator for the population mean is

$$\bar{y}_{\text{SRCS}} = \frac{N_l}{N} \sum_{i \in \mathcal{A}_l} \frac{\tau_i}{n_l}$$

with variance

$$V(\bar{y})_{\text{SRCS}} = \frac{N_l^2}{N^2} \left( 1 - \frac{n_l}{N_l} \right) \frac{V_\tau^2}{n_l},$$

where  $V_\tau^2 = \frac{1}{N_l - 1} \sum_{i \in \mathcal{U}_l} (\tau_i - \mu_\tau)^2$  and  $\mu_\tau = \sum_{i \in \mathcal{U}_l} \frac{\tau_i}{N_l}$ .



All SSU's in the sampled PSU's are surveyed, thus

$$\tau_i = \sum_{k \in \mathcal{U}_i} y_{ki}$$

is known of all selected PSU's. An unbiased estimator for the population mean is

$$\bar{y}_{\text{SRCS}} = \frac{N_1}{N} \sum_{i \in \mathcal{A}_1} \frac{\tau_i}{n_1}$$

An unbiased variance estimator is

$$\hat{V}(\bar{y}_{\text{SRCS}})_{\text{SRS}} = \frac{N_1^2}{N^2} \left( 1 - \frac{n_1}{N_1} \right) \frac{s_\tau^2}{n_1},$$

where

$$s_\tau^2 = \frac{1}{n_1 - 1} \sum_{i \in \mathcal{A}_1} (\tau_i - \bar{\tau})^2.$$

with  $\bar{\tau} = \sum_{i \in \mathcal{A}_1} \frac{\tau_i}{n_1}$ .

TABLE: Two Variations of a Population Composed of 5 clusters of Size 5

	1	2	3	4	5	$\mu_{i.}$	$V_{i.}^2$
1	1.0	2.0	3.0	4.0	5.0	3.0	2.5
2	6.0	7.0	8.0	9.0	10.0	8.0	2.5
3	11.0	12.0	13.0	14.0	15.0	13.0	2.5
4	16.0	17.0	18.0	19.0	20.0	18.0	2.5
5	21.0	22.0	23.0	24.0	25.0	23.0	2.5
$\mu_{.j}$	11.0	12.0	13.0	14.0	15.0	13.0	
$V_{.j}^2$	62.5	62.5	62.5	62.5	62.5		54.2

We have homogeneity of means between columns and heterogeneity of means between rows.

TABLE: Two Variations of a Population Composed of 5 clusters of Size 5

	1	2	3	4	5	$\mu_{i.}$	$V_{i.}^2$
1	1.0	2.0	3.0	4.0	5.0	3.0	2.5
2	6.0	7.0	8.0	9.0	10.0	8.0	2.5
3	11.0	12.0	13.0	14.0	15.0	13.0	2.5
4	16.0	17.0	18.0	19.0	20.0	18.0	2.5
5	21.0	22.0	23.0	24.0	25.0	23.0	2.5
$\mu_{.j}$	11.0	12.0	13.0	14.0	15.0	13.0	
$V_{.j}^2$	62.5	62.5	62.5	62.5	62.5		54.2

We have homogeneity of means between columns and heterogeneity of means between rows. Sample  $n = 10$  SSU's by

- 1) SRS  $n=10$
- 2) Simple random cluster sampling  $n_l = 2$ 
  - a) columns
  - b) rows

Comparing the mean and variance of  $\bar{y}$  and  $\bar{y}_{\text{SRCS}}$  after 100,000 samples:

TABLE: Results from the Simulation Study

	SRS	SRCS a	SRCS b
mean	13.0	13.0	13.0
var	3.2	0.8	18.9

Comparing the mean and variance of  $\bar{y}$  and  $\bar{y}_{\text{SRCS}}$  after 100,000 samples:

TABLE: Results from the Simulation Study

	SRS	SRCS a	SRCS b
mean	13.0	13.0	13.0
var	3.2	0.8	18.9

True values:

$$\begin{aligned} \mu &= 13 \\ V(\bar{y})_{\text{SRS}} &= 3.25 \text{ SRS} \\ V(\bar{y}_{\text{SRCS}})_{\text{SRS}} &= 0.75 \text{ SRCS a} \\ V(\bar{y}_{\text{SRCS}})_{\text{SRCS}} &= 18.75 \text{ SRCS b} \end{aligned}$$

Bias is not an issue, however variance is.

Comparing the mean and variance of  $\bar{y}$  and  $\bar{y}_{\text{SRCS}}$  after 100,000 samples:

TABLE: Results from the Simulation Study

	SRS	SRCS a	SRCS b
mean	13.0	13.0	13.0
var	3.2	0.8	18.9

True values:

$$\begin{aligned} \mu &= 13 \\ V(\bar{y})_{\text{SRS}} &= 3.25 \text{ SRS} \\ V(\bar{y}_{\text{SRCS}})_{\text{SRS}} &= 0.75 \text{ SRCS a} \\ V(\bar{y}_{\text{SRCS}})_{\text{SRCS}} &= 18.75 \text{ SRCS b} \end{aligned}$$

Bias is not an issue, however variance is.

If the cluster were strata, which stratification would you use, columns or rows?

Comparing the mean and variance of  $\bar{y}$  and  $\bar{y}_{\text{SRCS}}$  after 100,000 samples:

TABLE: Results from the Simulation Study

	SRS	SRCS a	SRCS b
mean	13.0	13.0	13.0
var	3.2	0.8	18.9

True values:

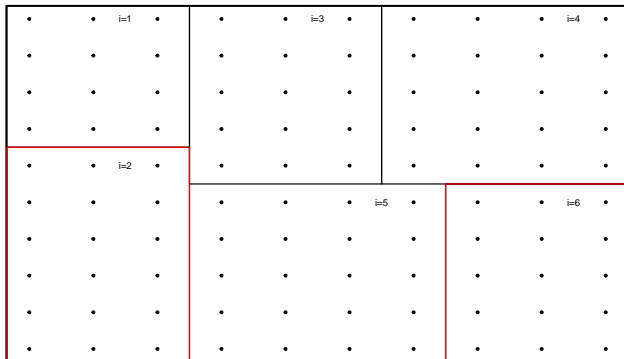
$$\begin{aligned} \mu &= 13 \\ V(\bar{y})_{\text{SRS}} &= 3.25 \text{ SRS} \\ V(\bar{y}_{\text{SRCS}})_{\text{SRS}} &= 0.75 \text{ SRCS a} \\ V(\bar{y}_{\text{SRCS}})_{\text{SRCS}} &= 18.75 \text{ SRCS b} \end{aligned}$$

Bias is not an issue, however variance is.

If the cluster were strata, which stratification would you use, columns or rows?

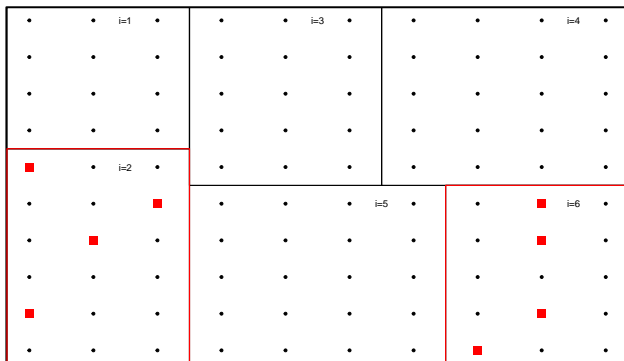
What good is for stratified sampling, i.e. low SSW, is bad of cluster sampling and vice versa.

A Population of 100 elements is clustered into  $N_1 = 6$  clusters and  $n_1 = 2$  clusters (PSU) are selected at the first sampling stage





A Population of 100 elements is clustered into  $N_1 = 6$  clusters and  $n_1 = 2$  clusters (PSU) are selected at the first sampling stage and  $n_j = 4$  elements are selected from each sampled cluster.



**First stage** A sample  $\delta_I$  of PSU's is drawn from  $\mathcal{U}_I$  according to some sampling design  $p_I(\cdot)$

**Second stage** For every  $i \in \delta_I$  a sample  $\delta_i$  of SSU's is selected from  $\mathcal{U}_i$  according to some design  $p_i(\cdot | \delta_I)$

The resulting sample of SSU's is denote  $\delta = \bigcup_{i \in \delta_I} \delta_i$ . In general, samples  $\delta_i$  are selected independently of each other, thus, the inclusion probability of a element  $k \in \mathcal{U}_i$  is

$$\pi_k = \pi_{Ii} \pi_{k|i} ,$$

where  $\pi_{Ii}$  is the probability of selecting the  $i$ -th PSU and  $\pi_{k|i}$  the probability of selecting the  $k$ -th SSU in the  $i$ -th PSU.

# ESTIMATION SIMPLE RANDOM

## TWO STAGE SAMPLING

Designs  $p_1(.)$  and  $p_i(.|\Delta_1)$  are both SRS. Since not all SSU's in the sampled PSU's are surveyed  $\tau_i$  has to be estimated by

$\hat{\tau}_i = \frac{N_i}{n_i} \sum_{k \in \Delta_i} y_{ki}$ . An unbiased estimator for the population mean is

$$\bar{y}_{2\text{SRS}} = \frac{N_1}{N} \sum_{i \in \Delta_1} \frac{\hat{\tau}_i}{n_1}$$

# ESTIMATION SIMPLE RANDOM

## TWO STAGE SAMPLING

Designs  $p_1(.)$  and  $p_i(.|\mathcal{A}_1)$  are both SRS. Since not all SSU's in the sampled PSU's are surveyed  $\tau_i$  has to be estimated by  $\hat{\tau}_i = \frac{N_i}{n_i} \sum_{k \in \mathcal{A}_i} y_{ki}$ . An unbiased estimator for the population mean is

$$\bar{y}_{2\text{SRS}} = \frac{N_1}{N} \sum_{i \in \mathcal{A}_1} \frac{\hat{\tau}_i}{n_i}$$

with variance

$$V(\bar{y}_{2\text{SRS}})_{\text{SRS}} = \frac{1}{N^2} \left( N_1^2 \left( 1 - \frac{n_1}{N_1} \right) \frac{V_\tau^2}{n_1} + \frac{N_1}{n_1} \sum_{i \in \mathcal{U}_1} N_i^2 \left( 1 - \frac{n_i}{N_i} \right) \frac{V_i^2}{n_i} \right),$$

where  $V_i^2 = \frac{1}{N_i-1} \sum_{k \in \mathcal{U}_i} (y_{ki} - \mu_i)^2$  with  $\mu_i = \sum_{k \in \mathcal{U}_i} \frac{y_{ki}}{N_i}$ .

# ESTIMATION SIMPLE RANDOM

## TWO STAGE SAMPLING

Designs  $p_1(\cdot)$  and  $p_i(\cdot|\delta_i)$  are both SRS. Since not all SSU's in the sampled PSU's are surveyed  $\tau_i$  has to be estimated by

$\hat{\tau}_i = \frac{N_i}{n_i} \sum_{k \in \delta_i} y_{ki}$ . An unbiased estimator for the population mean is

$$\bar{y}_{2\text{SRS}} = \frac{N_1}{N} \sum_{i \in \delta_1} \frac{\hat{\tau}_i}{n_i}$$

An unbiased variance estimator is given by

$$\hat{V}(\bar{y}_{2\text{SRS}})_{\text{SRS}} = \frac{1}{N^2} \left( N_1^2 \left( 1 - \frac{n_1}{N_1} \right) \frac{s_{\hat{\tau}}^2}{n_1} + \frac{N_1}{n_1} \sum_{i \in \delta_1} N_i^2 \left( 1 - \frac{n_i}{N_i} \right) \frac{s_i^2}{n_i} \right),$$

where  $s_{\hat{\tau}}^2 = \frac{1}{n_1 - 1} \sum_{i \in \delta_1} (\hat{\tau}_i - \bar{\hat{\tau}})^2$  with  $\bar{\hat{\tau}} = \sum_{i \in \delta_1} \frac{\hat{\tau}_i}{n_1}$  and

$s_i^2 = \frac{1}{n_i - 1} \sum_{k \in \delta_i} (y_{ki} - \bar{y}_i)^2$  with  $\bar{y}_i = \sum_{k \in \delta_i} \frac{y_{ki}}{n_i}$ .

There are good reasons to deviate from the simple selection procedure that gives every unit the same inclusion probability. If good prior information is available, its incorporation into the sampling design can dramatically improve the efficiency of an estimator.

An optimal allocation would be favorable to a proportional allocation.

Selecting the elements proportional to a variable that is correlated to the variable of interest can greatly improve the quality of estimates.

There are many techniques (i.e. sampling algorithms) to select elements with unequal probabilities [Tillé, 2006].

A design unbiased estimator for the total  $\tau = \sum_{k \in \mathcal{U}} y_k$  is given by

$$\hat{\tau}_{\pi} = \sum_{k \in \mathcal{d}} \frac{y_k}{\pi_k},$$

which is also known as *Horvitz-Thompson* (HT) or  $\pi$ -estimator.

A design unbiased estimator for the total  $\tau = \sum_{k \in \mathcal{U}} y_k$  is given by

$$\hat{\tau}_{\pi} = \sum_{k \in \mathcal{d}} \frac{y_k}{\pi_k},$$

which is also known as *Horvitz-Thompson* (HT) or  $\pi$ -estimator. The variance of  $\hat{\tau}_{\pi}$  is

$$V(\hat{\tau}_{\pi}) = \sum_{k \in \mathcal{U}} \sum_{l \in \mathcal{U}} (\pi_{kl} - \pi_k \pi_l) \frac{y_k}{\pi_k} \frac{y_l}{\pi_l},$$

which can be estimated by

$$\hat{V}(\hat{\tau}_{\pi})_1 = \sum_{k \in \mathcal{d}} \sum_{l \in \mathcal{d}} \frac{(\pi_{kl} - \pi_k \pi_l)}{\pi_{kl}} \frac{y_k}{\pi_k} \frac{y_l}{\pi_l}.$$



A design unbiased estimator for the total  $\tau = \sum_{k \in \mathcal{U}} y_k$  is given by

$$\hat{\tau}_{\pi} = \sum_{k \in \mathcal{d}} \frac{y_k}{\pi_k},$$

which is also known as *Horvitz-Thompson* (HT) or  $\pi$ -estimator. For a fixed size design, we may write the variance of  $\hat{\tau}_{\pi}$  as

$$V(\hat{\tau}_{\pi}) = -\frac{1}{2} \sum_{k \in \mathcal{U}} \sum_{l \in \mathcal{U}} (\pi_{kl} - \pi_k \pi_l) \left( \frac{y_k}{\pi_k} - \frac{y_l}{\pi_l} \right)^2,$$

which can be estimated by

$$\hat{V}(\hat{\tau}_{\pi})_2 = -\frac{1}{2} \sum_{k \in \mathcal{d}} \sum_{l \in \mathcal{d}} \frac{(\pi_{kl} - \pi_k \pi_l)}{\pi_{kl}} \left( \frac{y_k}{\pi_k} - \frac{y_l}{\pi_l} \right)^2.$$

A design unbiased estimator for the total  $\tau = \sum_{k \in \mathcal{U}} y_k$  is given by

$$\hat{\tau}_{\pi} = \sum_{k \in \mathcal{d}} \frac{y_k}{\pi_k},$$

which is also known as *Horvitz-Thompson* (HT) or  $\pi$ -estimator. For a fixed size design, we may write the variance of  $\hat{\tau}_{\pi}$  as

$$V(\hat{\tau}_{\pi}) = -\frac{1}{2} \sum_{k \in \mathcal{U}} \sum_{l \in \mathcal{U}} (\pi_{kl} - \pi_k \pi_l) \left( \frac{y_k}{\pi_k} - \frac{y_l}{\pi_l} \right)^2,$$

which can be estimated by

$$\hat{V}(\hat{\tau}_{\pi})_2 = -\frac{1}{2} \sum_{k \in \mathcal{d}} \sum_{l \in \mathcal{d}} \frac{(\pi_{kl} - \pi_k \pi_l)}{\pi_{kl}} \left( \frac{y_k}{\pi_k} - \frac{y_l}{\pi_l} \right)^2.$$

Provided that  $\pi_{kl} > 0$  for all  $k \neq l \in \mathcal{U}$  both variance estimators are unbiased. Nevertheless both variance estimators can become negative!

If there is some prior information available in the form of a variable  $\mathcal{X} = \{x_1, x_2, \dots, x_k, \dots, x_N\}$ , which is correlated to our variable of interest  $\mathcal{Y}$ , we can select elements proportional to it

$$\nu_k = \frac{x_k}{\sum_{k \in \mathcal{U}} x_k} n .$$

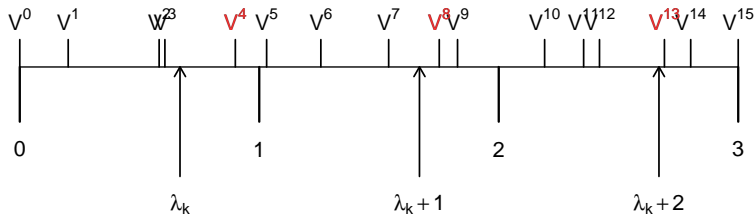
Unequal inclusion probabilities can reduce the variance of an estimator, *if* they are related to the variable of interest. We may have  $\nu_k = \pi_k$  but in general  $\nu_k$  can be greater than 1.

For instance, to estimate the sales in an industry, sampling companies or business with equal probabilities might be bad idea. It would be better to sample companies proportional to a variable that is related to their sales, say their number of employee. That way there is a much higher chance to included the biggest companies into the sample that accumulate the major share of the sales.

Note: In the extreme case if  $\pi_k = \alpha y_k$ , with  $\alpha \in \mathbb{R}$ , for all  $k \in \mathcal{U}$  and we have a fixed size design  $V(\hat{\tau}_\pi)$  would even be zero.

# SYSTEMATIC SAMPLING WITH UNEQUAL INCLUSION PROBABILITIES

Again the elements of the population are brought into a specific ordered and  $V^i = \sum_{k=1}^i \pi_k$ . Then  $\lambda_k$  is drawn from a uniform distribution between 0 and 1.



Systematic selection remains popular because of its simplicity. Also it can easily be applied to the case where an element can be selected more than one time, i.e.  $\pi_k \neq \nu_k > 1$ . Then we would use  $V^i = \sum_{k=1}^i \nu_k$ .

A two-stage Sampling Design:

**First Stage** Municipalities are the PSU's. The sampling design for the PSU's is a stratified design with an allocation proportional to the population within each stratum (not number of PSU's). Within the strata PSU's are sampled proportional to their population size.

**Second Stage** Persons are the SSU's. The SSU's are selected from the population register of the municipalities by a simple systematic sample.

Very large municipalities, (e.g. Berlin), are selected with certainty, this happens if  $\nu_i = \frac{N_i}{N} n_i > 1$ . The integer part of  $\nu_i$  indicates how many sampling points are *at least* associated with a municipalities. A sampling point is here a multiplier, indicating how many times  $n_i$  SSU's are selected from the  $i$ -th PSU, where  $n_i$  is usually fix for all PSU's.

For instance,  $\nu_i = 3.4$ , means that the  $i$ -th PSU will always be in the sample with at least 3 sampling points, but with a probability of 0.4 it can be in the sample with 4 sampling points.

# A TYPICAL SAMPLE OF PERSONS

IN GERMANY

Has this design equal inclusion probabilities?

Has this design equal inclusion probabilities?

Yes, if for each sampling point the same number of SSU's  $n_*$  is sampled. Because

$$\frac{N_j}{N} n_l \times \frac{n_*}{N_j} = \frac{n_l n_*}{N} .$$

Note, that  $n_l$  is not the size of the PSU sample, but the number of sampling points, which can be higher.



Selecting PSU's or clusters proportional to some size measure is very common. This however does not mean that the inclusion probabilities of elementary units are unequal.

The concept of two-stage designs can also be extended to three, four, or more stages. The principle of such multi-stage design remains the same, select clusters then select again within clusters.

There are many ways to optimize the sampling design with respect to one particular goal, i.e. the estimation of a specific statistic. However, it becomes difficult to optimize a design and at the same time retain a balance for a maximum of possible applications, which is a problem when planning a multipurpose survey that has a multitude of variables and covers different topics. Thus, simple designs, such as SRS or StrSRS, are justifiable, as these designs are robust towards any possible analysis of the sample data.

Multi-stage sampling is usually not a matter of choice, but done out of necessity.

Most importantly, the same design weights ( $\pi_k^{-1}$ ) do *not* imply the same sampling variance. Different designs can be used to select samples with same  $\pi_k$ 's, however their  $\pi_{kl}$ 's might be very different and so is their associated sampling variance.

The design effect compares strategies, i.e. a combination of a sampling design and an estimator.

If  $p(\cdot)$  is some other design than SRS, however with  $\sum_{i=1}^N \pi_k$  equal to the sample size  $n$  of the SRS design, then the *design effect* for strategy  $(p(\cdot), \hat{\tau}_\pi)$  can be defined as

$$\text{deff}(p, \hat{\tau}_\pi) = \frac{V(\hat{\tau}_\pi)_p}{V(\hat{\tau}_\pi)_{\text{SRS}}} = \frac{\sum_{k=1}^N \sum_{l=1}^N (\pi_{kl} - \pi_k \pi_l) \frac{y_k}{\pi_k} \frac{y_l}{\pi_l}}{N^2 \left(1 - \frac{1}{N}\right) \frac{V^2}{n}}.$$

The design effect  $\text{deff}(p, \hat{\tau}_\pi)$  expresses how well a design  $p(\cdot)$  fares in comparison to reference design SRS.

$\text{deff}(p, \hat{\tau}_\pi) > 1$  precision is lost by not using SRS

$\text{deff}(p, \hat{\tau}_\pi) < 1$  precision is gained by not using SRS

Recall from stratified sampling:

TABLE: Population ANOVA

Source	sf	Sum of Squares
Between cluster	$N_I - 1$	$SSB = \sum_{i=1}^{N_I} N_i (\mu_i - \mu)^2$
Within cluster	$N - N_I$	$SSW = \sum_{i=1}^{N_I} (N_i - 1) V_i^2$
Total	$N - 1$	$SSTO = (N - 1) V^2$

TABLE: Two Variations of a Population Composed of 5 clusters of Size 5

	1	2	3	4	5	$\mu_{i.}$	$V_{i.}^2$
1	1.0	2.0	3.0	4.0	5.0	3.0	2.5
2	6.0	7.0	8.0	9.0	10.0	8.0	2.5
3	11.0	12.0	13.0	14.0	15.0	13.0	2.5
4	16.0	17.0	18.0	19.0	20.0	18.0	2.5
5	21.0	22.0	23.0	24.0	25.0	23.0	2.5
$\mu_{.j}$	11.0	12.0	13.0	14.0	15.0	13.0	
$V_{.j}^2$	62.5	62.5	62.5	62.5	62.5		54.2

TABLE: Design Effects

	SRCS a	SRCS b
<i>deff</i>	0.23077	5.76923



$$deff(p(.), \hat{\theta}) = \frac{V(\hat{\theta})_p}{V(\hat{\theta})_{SRS}}$$

Find suitable estimators for both  
the denominator and  
the numerator

$$deff(p(.), \hat{\theta}) = \frac{V(\hat{\theta})_p}{V(\hat{\theta})_{SRS}}$$

Find suitable estimators for both  
the denominator and  
the numerator

Treat the data as if it had arisen from SRS for estimation of the  
enumerator (but using available weights).

$$deff(p(.), \hat{\theta}) = \frac{V(\hat{\theta})_p}{V(\hat{\theta})_{SRS}}$$

Find suitable estimators for both  
the denominator and  
the numerator

Treat the data as if it had arisen from SRS for estimation of the  
enumerator (but using available weights).

It is also common to use SRSWR as the reference design, which  
can simplify the design effect estimation.

Since it is so difficult to estimate the design effect, models are used to describe it. Effectively an alternative version of the design effect is defined which is easier to estimate.

Again we have  $N_i$  clusters in the population of size  $N_i$ ,  
 $i = 1, \dots, N_i$

Variable  $\mathcal{Y}$  is assumed to follow the *common parameter model*,

$$\begin{aligned} E(y_{ki})_M &= \mu \\ V(y_{ki})_M &= \sigma^2 \\ \text{COV}(y_{ki}, y_{k'i'})_M &= \begin{cases} \sigma^2 \rho & \text{for } k \neq k', i = i' \\ 0 & \text{otherwise.} \end{cases} \end{aligned}$$

$\rho$  the so called *Intra-Class Correlation Coefficient* is a central parameter to the model. It determine how similar to each other the elements from the same cluster are.

The considered estimator is  $\bar{y}_w = \frac{\sum_{k \in \delta} w_k y_k}{\sum_{k \in \delta} w_k}$ , where  $w_k$  is the survey associated with  $k$ -th element, e.g.

$$w_k = \begin{cases} \pi_k^{-1} & \text{for } k \in \delta \\ 0 & \text{else} \end{cases}.$$

Under the model-based approach, *deff* can be displayed as the product of two factors:

***deff<sub>c</sub>*** design effect due to clustering

***deff<sub>p</sub>*** design effect due to unequal survey weights

They indicate loss in precision due to cluster sampling and unequal weights, respectively.

Note that under the model-based approach,  $w_k$  is treated as independent of  $y_k$  for all  $k \in \mathcal{U}$ . Potential gains in precision by using proportional to size design are not considered, quite the contrary, unequal weight will increase the variance of  $\bar{y}_w$  in this setting.

Under cluster sampling and two-stage sampling we can use:

$$\widehat{deff}_c = 1 + (\bar{b} - 1) \hat{\rho}$$

with

$\bar{b}$  as the average cluster size  $\frac{N}{N_1} \left( \frac{n}{n_1} \right)$  or an estimator for it and  
 $\hat{\rho}$  as an appropriate estimator of  $\rho$ .

When  $w_k$ 's vary we have

$$\widehat{deff}_c = 1 + (b^* - 1) \hat{\rho}$$

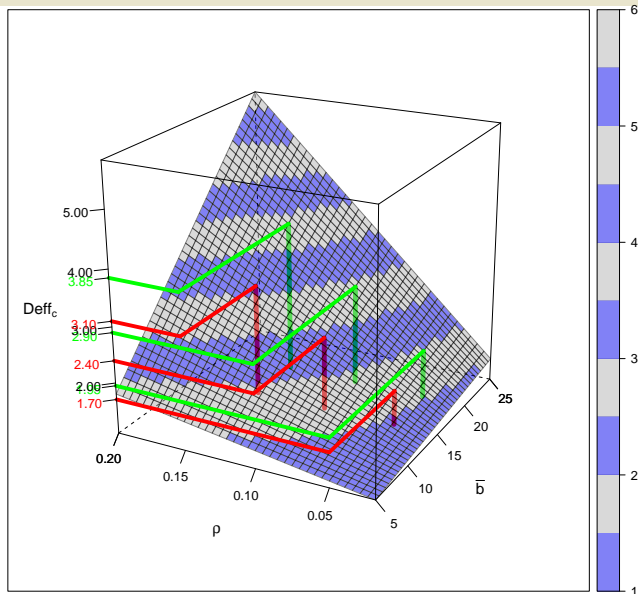
where  $b^* = \frac{\sum_{i \in \Delta_1} \left( \sum_{k \in \Delta_i} w_k \right)^2}{\sum_{k \in \Delta} w_k^2}$  is a kind of weighted average cluster size [?].

The design effect due to unequal weights is

$$deff_p = n \frac{\sum_{k \in \Delta} w_k^2}{(\sum_{k \in \Delta} w_k)^2}$$

Obviously, if  $w_k$ 's are constant,  $deff_p = 1$





There are many different ways to estimate  $\rho$  [?]. The classical ANOVA estimator is:

$$\hat{\rho} = \frac{\widehat{MSB} - \widehat{MSW}}{\widehat{MSB} + (K - 1)\widehat{MSW}},$$

where  $\widehat{MSB} = (n_l - 1)^{-1} \sum_{i \in \mathcal{I}_l} n_i (\bar{y}_i - \bar{y})^2$ ,

$\widehat{MSW} = (n - n_l)^{-1} \sum_{i \in \mathcal{I}_l} (n_i - 1) s_i^2$ , and

$$K = (n_l - 1)^{-1} \left( n - \sum_{i \in \mathcal{I}_l} \frac{n_i^2}{n} \right).$$

There are many different ways to estimate  $\rho$  [?]. The classical ANOVA estimator is:

$$\hat{\rho} = \frac{\widehat{MSB} - \widehat{MSW}}{\widehat{MSB} + (K - 1)\widehat{MSW}},$$

where  $\widehat{MSB} = (n_l - 1)^{-1} \sum_{i \in \mathcal{L}_l} n_i (\bar{y}_i - \bar{y})^2$ ,

$\widehat{MSW} = (n - n_l)^{-1} \sum_{i \in \mathcal{L}_l} (n_i - 1) s_i^2$ , and

$$K = (n_l - 1)^{-1} \left( n - \sum_{i \in \mathcal{L}_l} \frac{n_i^2}{n} \right).$$

The model-based approach is widely used, because it often presents the only option for data users to estimate a design effect.

For multi-stage design, the used model considers only the cluster effect of the PSU's and neglects any subsequent sampling stages.

# COMPARISON OF DESIGN BASED AND MODEL BASED DESIGN EFFECTS

TABLE: Intra-Class Correlation Coefficients and (Model) Design Effects

	SRCS a	SRCS b
$\rho$	-0.01626	0.99136
<i>deff</i>	0.93496	4.96546

TABLE: Design Effects

	SRCS a	SRCS b
<i>deff</i>	0.23077	5.76923



S. Gabler, S. Häder, & P. Lahiri.

A Model Based Justification of Kish's Formula for Design Effects for Weighting and Clustering.

*Survey Methodology*, 1999.



M. Ganninger.

Design Effects: Model-based versus Design-based Approach

PhD Thesis, *GESIS-Schriftenreihe Band 3*, 2009



C.-E. Särndal, B. Swensson, & J. Wretman.

Model Assisted Survey Sampling

*Springer*, 1992.



Y. Tillé.

Sampling Algorithms

*Springer Series in Statistics: Springer*, 2006.