# Ratio Estimation

► Ratio estimators involve two characteristics:

  ► Y: a characteristic we are interested in/ study variable

  ► X: a characteristic that is related with Y/ auxiliary variable

► **Why we consider ratio estimation?**

- ► In practic, knowledge of the ratio of a characteristic is as important as or sometimes more important than that of population totals and means.
- ► For instance: In socio-economic surveys, when we want to estimate average yield of crops per acre or per capita expenditure.
- ► We also use this estimator to increase the precision of estimated means or totals.

- ► Let Y be the study variable and X be an auxiliary variable which is correlated with Y. The population total X must be known by some earlier surveys.

- ► Let $(x_1, y_1), (x_2, y_2), \ldots\ldots, (x_n, y_n)$ be the r. s. of size n of the paired variable (X, Y) drawn by SRS from a population of size N.

- ► The ratio estimate of population total Y is gives as:

$$Y_R = \frac{\bar{y}}{\bar{x}} X$$

where $\bar{y} = \frac{1}{n} \sum_{i=1}^{n} y_i$

$\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$

$\hat{R} = \frac{\bar{y}}{\bar{x}}$: Ratio estimator in SRS.

## BIAS OF $\hat{R}$

- In SRSWOR(N,n) the bias of the ratio estimator $\hat{R}$ is approximately given by

$$B(\hat{R}) = \frac{1-f}{n\bar{X}}(RS_x^2 - \rho S_x S_y)$$

where $1 - f = \frac{N-n}{N}$

$\rho$ : Correlation coefficient between Y and X.

$S_x$ : Population standard deviation of X

$S_y$ : Population standard deviation of Y

- Bias is small if
  - sample size n is large
  - Sample fraction n/N is large
  - $\bar{X}$ is large
  - $S_x$ is small
  - High positive correlation between X and Y

## ▶ MSE OF $\hat{R}$

- ▶ Now, we are interested in the precision of ratio estimator
- ▶ Since, ratio estimators are biased in general, its MSE around R would be of greater interest.
- ▶ In SRSWOR(N,n), an approximation to MSE of $\hat{R} = \frac{\bar{y}}{\bar{x}}$ around R=$\frac{Y}{X}$ is given by

$$E(\hat{R} - R)^2 = \frac{1-f}{n} \frac{\left(S_y^2 + R^2 S_x^2 - 2\rho S_x S_y\right)}{\bar{X}^2}$$

- MSE is less biased if

- sample size n is large or Sample fraction n/N is large

- High positive correlation between X and Y

► **Efficiency of ratio eatimators im comparison to SRSWOR**
WE have,

$$MSE(\hat{\bar{Y}}_R) = \bar{Y}^2 \frac{f}{n}(C_x^2 + C_y^2 - 2\rho C_x C_y)$$

$$Var_{SRS}(\bar{y}) = \bar{Y}^2 \frac{f}{n} C_y^2$$

Hence, ratio estimator $\hat{\bar{Y}}_R$ becomes more efficient than conventional $\hat{\bar{Y}}$ if

$$MSE(\hat{\bar{Y}}_R) < Var_{SRS}(\bar{y})$$

$$\Rightarrow \bar{Y}^2 \frac{f}{n}(C_x^2 + C_y^2 - 2\rho C_x C_y) < \bar{Y}^2 \frac{f}{n} C_y^2$$

$$\Rightarrow C_x^2 - 2\rho C_x C_y < 0$$

Hence, ratio estimator works better when

$$\rho > \frac{1}{2}\frac{C_x}{C_y} \ \ when \ R > 0 \ \ or \ \ \rho < -\frac{|C_x|}{|C_y|} \ \ when \ R < 0$$

- Unbiased estimators of $\hat{R}$ have been proposed by Hartley and Ross(1854),Murthy and Nanjamma(1959) and nieto de parcual(1961)
- We will be discussing about Hartley Ross estimator.

► **FOR FINDING U. E. OF R**

  ► Correctng the estimator for bias without changing the sampling design

  ► Changing the sampling design without changing the form of R

  ► Here we will be correcting tbias i. e. the changed estimator is

$$\hat{Y}_R^* = \hat{Y}_R - B(\hat{Y}_R)$$

where $B(\hat{Y}_R)$ is an u. e. of bias of $\hat{Y}_R$.

► THE HARTLEY ROSS ESTIMATOR

    ► For SRSWR,

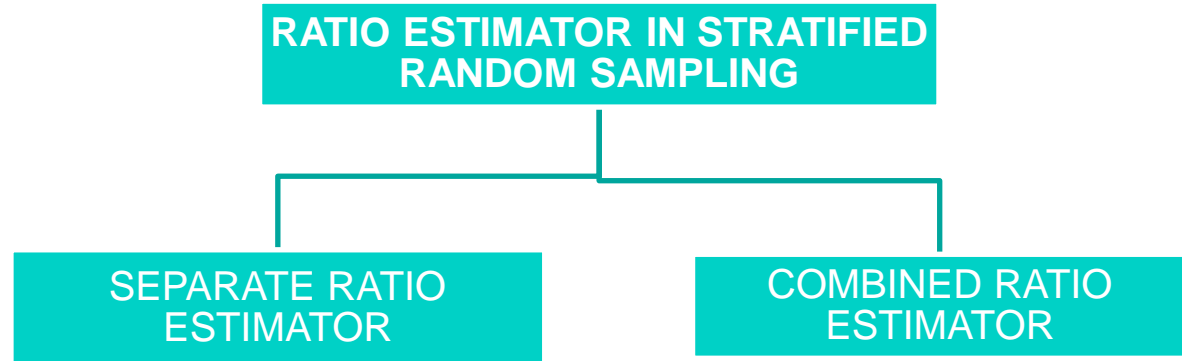$$\hat{Y}_R^* = \bar{r}X + \frac{n}{n-1}N(\bar{y} - \bar{r}\bar{x})$$

    ► For SRSWOR

$$\hat{Y}_R^* = \bar{r}X + \frac{n}{n-1}(N-1)(\bar{y} - \bar{r}\bar{x})$$

where $\bar{r} = \frac{1}{n}\sum_{i=1}^{n} r_i$, where $r_i = \frac{y_i}{x_i}$ is the ratio for each unit.

# For stratified random sampling:

**RATIO ESTIMATOR IN STRATIFIED RANDOM SAMPLING**

SEPARATE RATIO ESTIMATOR

COMBINED RATIO ESTIMATOR

► Separate ratio estimator:

Here the estimate of Y is the sum of the ratio estimate of $y_h$ over all strata.

$$\hat{y}_{RS} = \sum_{h=1}^{L} \frac{\bar{y}_h X_h}{\bar{x}_h}$$

this clearly is a biased estimator with bias

$$B(\hat{y}_{RS}) = \sum_{h=1}^{L} \frac{1-f_h}{n_h} Y_h (C_{xx}^{(h)} - \rho^{(h)} C_x^{(h)} C_y^{(h)})$$

Where, $C_{xx}^{(h)}$ denotes $(C_x^{(h)})^2 = (cv(x))^2$ for the h th stratum. now,

$$V(\hat{y}_{RS}) \simeq \sum_{h=1}^{L} \frac{N_h^2(1-f_h)}{n_h} (S_{yh}^2 + R_h^2 S_{xh}^2 - 2R_h \rho_h S_{xh} S_{yh})$$

Estimators of $V(\hat{y}_{RS})$ are (i) $\sum_{h=1}^{L} v_1(\hat{y}_{Rh})$, and (ii) $\sum_{h=1}^{L} v_2(\hat{y}_{Rh})$

for a separate ratio estimator,

$$|B(\widehat{y}_{RS})| \leq \sum_{h=1}^{L} |B(\widehat{y}_{RS})|$$

Hence,

$$\frac{|B(\widehat{y}_{RS})|}{\sigma(\widehat{y}_{RS})} \leq \frac{L \max_h |B(\hat{y}_{Rh})|}{\sqrt{L} \operatorname*{average}_h \sigma(\hat{y}_{Rh})}$$

In short, $\widehat{y}_{RS}$ may be useful if

1.  $x_h$ is known for all h,

2.  There are reasons to believe that $R'_h s$ differ considerably from stratum to stratum.

- $\widehat{y}_{RS}$ will be a good estimator if:
  1. $n_h$ is large enough for the approximate formula for $V(\widehat{y}_{RS})$ to hold,
  2. the population in such that and the sample size is large enough so that $\sqrt{L} cv(\bar{x}_h)$ does not exceed 0.3.
  3. there are reasons to believe that the individual biases in starta should nearly cancel each other, and
  4. the circumstances are such that $\widehat{y}_{Rh}$ is a good estimator of $Y_h$ in each stratum.

► Combined ratio estimator:

If $X_h$ is not known, and if $R'_h s$ do not differ considerably from each other and if $n'_h s$ are not large enough to validate the large sample formula for $V(\widehat{y}_{RS})$, one can use the combined ratio estimator,

$$\widehat{y}_{RC} = \widehat{R}_C X$$

where, the combined ratio estimator $\widehat{R}_c$ of $R$ is

$$\widehat{R}_c = \frac{\overline{y}_{st}}{\overline{x}_{st}}$$

Hence,

$$V(\widehat{y}_{RC}) \simeq \sum_{h=1}^{L} \frac{N_h(N_h - n_h)}{n_h}(S_{yh}^2 + R^2 S_{xh}^2 - 2RS_{xyh})$$

an estimator of $V(\widehat{y}_{RC})$ is

$$v(\widehat{y}_{RC}) \simeq \sum_{h=1}^{L} \frac{N_h(N_h - n_h)}{n_h}(s_{yh}^2 + \widehat{R}_c^2 s_{xh}^2 - 2\widehat{R}_c s_{xyh})$$

$v(\widehat{y}_{RC})$ is biased estimator.

using the approach of Goodman and Hartley, we get

$$E(\widehat{R}_c) - R = -\frac{Cov(\widehat{R}_c, \overline{x}_{st})}{\overline{x}}$$

so the bias in $\widehat{R}_c$ is

$$B(\widehat{R}_c) = -\frac{\rho(\widehat{R}_c, \overline{x}_{st})\sigma(\widehat{R}_c)\sigma(\overline{x}_{st})}{\overline{X}}$$

thus,

$$\frac{|B(\widehat{R}_c)|}{\sigma(\widehat{R}_c)} \leq \frac{\sigma(\overline{x}_{st})}{\overline{X}} = cv(\overline{x}_{st})$$

the bias of $\widehat{R}_c$ is negligible as compared to its standard error provided $cv(\overline{x}_{st}) \leq 0.1$

# Comparison of Combined and Separate Ratio Estimators

So unless $R_i$ varies considerably, the use of $\hat{\bar{Y}}_{Rc}$ would provide an estimate of $\bar{Y}$ with negligible bias and the precision as good as $\hat{\bar{Y}}_{Rs}$.

- If $R_i \neq R$, $\hat{\bar{Y}}_{Rs}$ can be more precise but bias may be large.

- If $R_i \simeq R$, $\hat{\bar{Y}}_{Rc}$ can be as precise as $\hat{\bar{Y}}_{Rs}$ but its bias will be small. It also does not require knowledge of $\bar{X}_1, \bar{X}_2, ..., \bar{X}_k$.

▸ Both methods produce biased estimated of $y_U$ or t with the bias decreasing as the sample size increases. Specifically,

  ▸ For the separate ratio estimator we require that the sample sizes within each stratum be large for the bias to be small.
  ▸ For the combined ratio estimator we require that the total sample size be large for the bias to be small.
  ▸ Thus, the bias will tend to be less serious for the combined estimator than for the separate estimator.

▸ Unless the ratio relationship is similar across the strata, the separate estimator will be more efficient (have smaller variability) than the combined estimator.

▸ The lower efficiency of the combined estimator, however, is often offset by smaller bias and the fact that we do not need to know the separate $\overline{x}_{Uh}$ stratum means

**EXAMPLE 4.4**  Table 4.2 shows the area under paddy ($y$) and cultivated area ($x$) for the villages in a zone divided into two strata. From each stratum draw a random sample of size 3 without replacement and hence obtain a separate ratio and a combined ratio estimate of total area under paddy for all the villages along with their variance estimates.

**Table 4.2**  Data Showing the Area under Paddy ($y$) and Cultivated Area ($x$) for 18 Villages (divided into two strata)

| | Stratum I | | | Stratum II | |
|---|---|---|---|---|---|
| Sl. no. | x | y | Sl. no. | x | y |
| 1 | 630 | 250 | 1 | 1012 | 340 |
| 2 | 729 | 248 | 2 | 1181 | 416 |
| 3 | 865 | 359 | 3 | 780 | 247 |
| 4 | 305 | 129 | 4 | 815 | 306 |
| 5 | 569 | 223 | 5 | 1120 | 403 |
| 6 | 427 | 335 | 6 | 659 | 271 |
| 7 | 326 | 412 | 7 | 897 | 357 |
| 8 | 481 | 503 | 8 | 783 | 295 |
| | | | 9 | 689 | 218 |
| | | | 10 | 1217 | 398 |

Suppose the samples drawn from Stratum I and II are, respectively, (3, 5, 8) and (2, 6, 9).

$$s_{x_1}^2 = 4470, \qquad s_{x_2}^2 = 85908, \qquad s_{y_1}^2 = 19605, \qquad s_{y_2}^2 = 10392.5,$$

$$s_{xy1} = -6614.7, \qquad s_{xy2} = 28417.5,$$

$$\hat{y}_{RS} = 5570.4, \qquad \hat{y}_{RC} = 6065.3, \qquad Y = 5709, \qquad V(\hat{y}_{RS}) = 362104.144,$$

$$V(\hat{y}_{RC}) = 301661.913, \qquad v(\hat{y}_{RS}) = 559144.6, \qquad v(\hat{y}_{RC}) = 486595.238.$$

Here both $v(\hat{y}_{RS})$ and $v(\hat{y}_{RC})$ overestimate the respective variances. $\hat{y}_{RC}$ is a better estimator than $\hat{y}_{RS}$.

For cluster sampling:

➢ Let, there are N cluster and $M_i$ be the size of $i^{th}$ cluster.

➢ $M_o = \sum_{i=1}^{N} M_i$

$\bar{M} = \frac{1}{N} \sum_{i=1}^{N} M_i$

$\bar{y}_i = \frac{1}{M_i} \sum_{j=1}^{M_i} y_{ij} : mean\ of\ ith\ cluster$

$\bar{Y} = \frac{1}{M_o} \sum_{i=1}^{N} \sum_{j=1}^{M_i} y_{ij} = \frac{1}{N} \sum_{i=1}^{N} \frac{M_i}{\bar{M}} \bar{y}_i$

Estimator based on ratio method of estimation:

➢ consider the study variable $U_i$ and auxiliary variable $V_i$ as

$$U_i = \frac{M_i \overline{y_i}}{\overline{M}}$$

$$V_i = \frac{M_i}{\overline{M}}, i = 1, 2, \dots, N$$

$$\bar{u} = \frac{1}{n} \sum_{i=1}^{n} u_i$$

$$\bar{v} = \frac{1}{n} \sum_{i=1}^{n} V_i$$

➢ The ratio estimator based on U and V is

$$\widehat{\overline{Y_R}} = \frac{\bar{u}}{\bar{v}} \bar{V} = \frac{\sum_{i=1}^{n} M_i \overline{y_i}}{\sum_{i=1}^{n} M_i}$$

▶ Since the ratio estimator is biased, so $\bar{y}_c^{**}$ is also a biased estimator.

$$Bias(\bar{y}_c^{**}) = \frac{N-n}{Nn}\left(\frac{S_v^2}{\bar{V}^2} - \frac{S_{uv}}{\bar{U}\bar{V}}\right)\bar{U}$$

where $\bar{U} = \frac{1}{N}\sum_{i=1}^{N} U_i = \frac{1}{N\bar{M}}\sum_{i=1}^{N} M_i\bar{y}_i$

$$S_v^2 = \frac{1}{N-1}\sum_{i=1}^{N}\left(\frac{M_i}{\bar{M}} - 1\right)^2$$

$$S_{uv} = \frac{1}{N-1}\sum_{i=1}^{N}\left(\frac{M_i\bar{y}_i}{\bar{M}} - \frac{1}{N\bar{M}}\sum_{i=1}^{N} M_i\bar{y}_i\right)\left(\frac{M_i}{\bar{M}} - 1\right)$$

▶ An estimator of $MSE$ can be obtained as

$$\widehat{MSE}(\bar{y}_c^{**}) = \frac{N-n}{Nn}\frac{1}{n-1}\sum_{i=1}^{n}\left(\frac{M_i}{\bar{M}}\right)^2(\bar{y}_i - \bar{y}_c^{**})^2$$

The estimator $\bar{y}_c^{**}$ is biased but consistent.

- Since $\bar{\bar{y}}_c = \frac{1}{n} \sum\limits_{i=1}^{n} \bar{y}_i$ is a biased estimator of population mean and

$$Bias\left(\bar{\bar{y}}_c\right) = -\left(\frac{N-1}{M_o}\right) S_{m\bar{y}}$$

- Since SRSWOR is used, so

$$s_{m\bar{y}} = \frac{1}{n-1} \sum\limits_{i=1}^{n} \left(M_i - \bar{m}\right)\left(\bar{y}_i - \bar{\bar{y}}_c\right), \quad \bar{m} = \frac{1}{n} \sum\limits_{i=1}^{n} M_i$$

is an unbiased estimator of

$$S_{m\bar{y}} = \frac{1}{N-1} \sum\limits_{i=1}^{N} \left(M_i - \bar{M}\right)\left(\bar{y}_i - \bar{Y}\right),$$

i.e., $E\left(s_{m\bar{y}}\right) = S_{m\bar{y}}$

► So it follow that

$$E\left(\bar{\bar{y}}_c\right) - \bar{Y} = -\left(\frac{N-1}{N\bar{M}}\right) E\left(s_{m\bar{y}}\right)$$

or $E\left[\bar{\bar{y}}_c + \left(\frac{N-1}{N\bar{M}}\right) S_{m\bar{y}}\right] = \bar{Y}$

so,

$$\bar{\bar{y}}_c^{**} = \bar{\bar{y}}_c + \left(\frac{N-1}{N\bar{M}}\right) S_{m\bar{y}}$$

is an unbiased estimator of the population mean $\bar{Y}$

This estimator is based on unbiased ratio type estimator. This can be obtained by replacing the study variable (earlier $y_i$) by $\frac{M_i}{M}\bar{y}_i$ and auxiliary variable (earlier $x_i$) by $\frac{M_i}{M}$.

The exact variance of this estimate is complicated and does not reduces to a simple form. The approximate variance upto first order of approximation is

$$var(\bar{\bar{y}}_c^{**}) = \frac{1}{n(N-1)} \sum_{i=1}^{N} \left[ \left( \frac{M_i}{\bar{M}}(\bar{y}_i - \bar{Y}) \right) - \left( \frac{1}{N\bar{M}} \sum_{i=1}^{N} \bar{y}_i \right) (M_i - \bar{M}) \right]^2$$

► A consistent estimate of this variance is

$$\widehat{var}(\bar{\bar{y}}_c^{**}) = \frac{1}{n(n-1)} \sum_{i=1}^{N} \left[ \left( \frac{M_i}{\bar{M}}(\bar{y}_i - \bar{y}_c) \right) - \left( \frac{1}{N\bar{M}} \sum_{i=1}^{N} \bar{y}_i \right) \left( M_i - \frac{\sum_{i=1}^{n} M_i}{n} \right) \right]^2$$

The variance of $\bar{\bar{y}}_c^{**}$ will be smaller than that of $\bar{y}_c^{**}$ provided the regression coefficient of $\frac{M_i \bar{y}_i}{\bar{M}}$ on $\frac{M_i}{\bar{M}}$ is nearer to $\frac{1}{N} \sum_{i=1}^{N} \bar{y}_i$ than

to $\frac{1}{N} \sum_{i=1}^{N} M_i \bar{y}_i$.

► **Ratio Estimator Under Probability Proportional To Size With Replacement Sampling Scheme**

A sampling scheme with replacement in which each sampling unit has unequal probability of selection, the probability being proportional to the size of the auxiliary variable associated with the particular unit , is called probability proportional to size and with replacement (PPSWR) sampling scheme.

Let Y be a study variable and X be an auxiliary variable. For example, consider we want to estimate the population in the villages of a particular district. Then we would choose as our auxiliary variable a variable on which we have information , e.g.:

(a) Area of each village of the district (correlation with a study variable $= 0.70$, say);

(b) Number of households in each village of the district (correlation with a study variable $= 0.85$, say) .

On the basis of the above information , we would choose the auxiliary variable which has maximum correlation with the study variable. Thus the variable at ( b ) may be a more useful auxiliary variable when selecting a sample using PPSWR sampling.

In this case Ratio Estimator ($\hat{Y}_R = \frac{\hat{Y}}{\hat{X}}X$) would be efficient if the set of probabilities used for selection is appropriate for both the study and auxiliary variables. And also $y/p$ and $x/p$ should be positively related. (Since we know that $V(\hat{Y}_{PPS}) < V(\hat{Y}_{SRS})$ if $y^2/x$ and $x$ are positively correlated i.e. if $y$ and $x$ are positively correlated.)

Now for PPSWR sampling scheme, the expressions for $V(\hat{X})$ , $V(\hat{Y})$ , $Cov(\hat{X}, \hat{Y})$ are given by

$$V(\hat{X}) = \frac{1}{n}\left(\sum_{i=1}^{N} \frac{X_i^2}{P_i} - X^2\right) , \quad V(\hat{Y}) = \frac{1}{n}\left(\sum_{i=1}^{N} \frac{Y_i^2}{P_i} - Y^2\right) \text{ and}$$

$$Cov(\hat{X}, \hat{Y}) = \frac{1}{n}\left(\sum_{i=1}^{N} \frac{X_i Y_i}{P_i} - XY\right)$$

$X =$ Total Value of the Auxiliary Variable, $Y =$ Total Value of the Study Variable and $P_i =$ Probability that i th population unit is included in the sample

Hence the expression for bias of the ratio estimator $\hat{Y}_R$ is given by

$$\hat{B}_R = \frac{1}{X^2}(RV(\hat{X}) - Cov(\hat{X}, \hat{Y})) \text{ , where } R = \frac{Y}{X}$$

$$= \frac{1}{X^2}\left(\frac{Y}{X} \cdot \frac{1}{n}\left(\sum_{i=1}^{N} \frac{X_i^2}{P_i} - X^2\right) - \frac{1}{n}\left(\sum_{i=1}^{N} \frac{X_i Y_i}{P_i} - XY\right)\right)$$

$$= \frac{1}{nX^2}\left(\sum_{i=1}^{N} \frac{X_i}{P_i}(RX_i - Y_i)\right)$$

Also the expression for MSE of $\hat{Y}_R$ is given by

$$\hat{M}_R = \frac{1}{X^2}(V(\hat{Y}) - 2RCov(\hat{X}, \hat{Y}) + R^2 V(\hat{X})) \text{ , where } R = \frac{Y}{X}$$

$$= \frac{1}{X^2}\left(\frac{1}{n}\left(\sum_{i=1}^{N} \frac{Y_i^2}{P_i} - Y^2\right) - 2\frac{Y}{X} \cdot \frac{1}{n}\left(\sum_{i=1}^{N} \frac{X_i Y_i}{P_i} - XY\right) + \frac{Y^2}{X^2} \cdot \frac{1}{n}\left(\sum_{i=1}^{N} \frac{X_i^2}{P_i} - X^2\right)\right)$$

$$= \frac{1}{nX^2}\left(\sum_{i=1}^{N} \frac{1}{P_i}(RX_i - Y_i)^2\right)$$

In previous slide, we have found an unbiased estimator of $\hat{R}$ and as well as an unbiased estimator of $\hat{Y}$. That was the Hartley-Ross Estimator(1954). Now we will try to find another unbiased estimator of $\hat{R}$ as well as another unbiased estimator of $\hat{Y}$. To find these, we will take help of Probability Proportional to Aggregate Size (PPAS) sampling proposed by Midzuno (1952) using the unbiased ratio estimates of Lahiri (1951) .

## ▶ PPAS Sampling Method

Let us denote the probability allotted to the sample $'s'$ to be $P(s)$, where $s \in S$ being the sample space.

Then the sampling design is given by $(s, P(s), s \in S)$.

Now the ratio estimator is given by

$\hat{Y}_R = \hat{R}_s X$, where $\hat{R}_s = \frac{\hat{Y}_s}{\hat{X}_z}$.

$\hat{Y}_s$ is an unbiased estimator of Y.

$\hat{X}_s$ is an unbiased estimator of X.

Now $E_P(\hat{R}_s) \neq R$ i.e. it is a biased estimator.

To get an unbiased estimator we modify the sampling design as $(s, P'(s), s \in S)$ where $P'(s) \propto P(s) \hat{X}_s$.

This sampling design is called <u>Probability Proportional to Aggregate Size (PPAS) Sampling Design.</u>

Now $\sum_{s \in S} P'(s) = 1$ .

$\implies \sum_{s \in S} kP(s)\hat{X}_s = 1$ , $k$ being the constant of proportionality.

$\implies k \sum_{s \in S} P(s)\hat{X}_s = 1$

$\implies kE_P(\hat{X}) = 1$

$\implies kX = 1$

$\implies kX = 1$

$\implies k = \frac{1}{X}$

- Under this new sampling design

$$E_{P'}(\hat{R}_s)$$

$$= \sum_{s \in S} \hat{R}_s P'(s)$$

$$= \sum_{s \in S} \frac{\hat{Y}_s}{\hat{X}_s} \cdot \frac{1}{X} \cdot P(s) \cdot \hat{X}_s$$

$$= \frac{1}{X} \cdot \sum_{s \in S} \hat{Y}_s \cdot P(s)$$

$$= \frac{1}{X} \cdot Y$$

$$= R \ \forall R$$

$\therefore \hat{R}_S$ becomes unbiased for $R$ under $P'$.

$E_{P'}(\hat{Y}_R)$

$= E_{P'}(\hat{R}_s X)$

$= RX$

$= Y \; \forall Y$

$\therefore \; \hat{Y}_R$ becomes unbiased for $Y$ under $P'$.

In PPAS, first unit is selected using PPS and the remaining $(n-1)$ units are selected using simple random sampling without replacement.

Gauran and Poblador (2012) used sampling with PPAS to estimate total production area of top cereals and root crops across Philippine regions.

# Thank you