

Sampling Design

Some definitions :

① Survey Population >>

Let, N be a known number of units, e.g. hospitals, schools, colleges etc., each assignable identifying labels $1, 2, \dots, N$ and taking values Y_1, Y_2, \dots, Y_N respectively of a real valued random variable Y , which are initially unknown to the investigator ^{who} is interested to estimate the total $Y = \sum_{i=1}^N Y_i$ or the mean $\bar{Y} = \frac{1}{N} \sum_{i=1}^N Y_i$.

We call the sequence $U = \{1, 2, \dots, N\}$ of labels a survey population. Thus, the survey population means a finite population in which the individuals are such that they are all concentrate and tangible (tagged with labels).

② Sample >> Selection of units from the population leads to a sequence $s = \{i_1, i_2, \dots, i_n\}$ which is called a sample. Here i_1, i_2, \dots, i_n are the elements of U , not necessarily distinct from one another, but the order of its appearance is maintained. $n = n(s)$ is ~~referred~~ referred as the size of the sample s , while the effective sample size is denoted by ~~v~~ $v(s)$ and defined as $v(s) = |s|$ (ie. cardinality of s) ie. $v(s)$ is the no. of distinct elements in the sample s .

③ Survey >>

By survey, we mean the process of ascertaining the variate values. Census or complete enumeration is carried out to help the subsequent survey.

④ Survey Data >> Once a specific sample is chosen, we suppose that it is possible to ascertain the values Y_1, Y_2, \dots, Y_n of the variable Y with the respective ~~units~~ units in s . Then

$$d = \{ (i_1, Y_{i_1}), (i_2, Y_{i_2}), \dots, (i_n, Y_{i_n}) \}$$
$$= \{ (i, Y_i) ; i \in s \}$$
 is called ^{the} _^ survey data.

⑤ Estimator >>

An estimator t is a real-valued function $t(d)$, which is free of Y_i for $i \notin S$ but may or may not involve Y_i for $i \in S$.

$t(d)$ can be written alternatively as $t(d) = t(S, \underline{y})$

where $\underline{y} = (Y_1, Y_2, \dots, Y_N)$

An estimator of the special importance is the sample mean \bar{y} (for estimating the population mean) defined as,

$$t(S, \underline{y}) = \frac{1}{n(S)} \sum_{i=1}^N f_{si} Y_i \quad [f_{si} = 0 \text{ if } i \notin S]$$

where, f_{si} is the frequency of the i th unit in the sample S , such that $\sum_{i=1}^N f_{si} = n(S)$ & $N\bar{y}$ is called the expansion estimator of the population total.

⑥ Linear and homogeneous linear estimator >>

Consider an estimator of the form

$$t(S, \underline{y}) = b_S + \sum_{i=1}^N b_{si} Y_i$$

with $b_{si} = 0$ if $i \notin S$ and b_S and b_{si} are free of Y . It is called a linear estimator and keeping $b_S = 0$, we obtain a homogeneous linear estimator.

It should be noted that $t(S, \underline{y})$ is linear or homogeneously linear in Y_i , $i \in S$, but it may be a non-linear function of two random variables, e.g. taking $b_S = 0$, $b_{si} = \frac{x}{\sum_{i=1}^N f_{si} X_i}$ where X_i be

the value of some other variable X on $i \in U$. This gives

$$t(S, \underline{y}) = \sum_{i=1}^N \left(\frac{x}{\sum_{i=1}^N f_{si} X_i} Y_i \right) \quad (\text{Ratio estimator})$$

This is not linear in x and Y .

⑤ Sampling Design \Rightarrow

A sample s is chosen by assigning a selection probability to the sample s , denoted by $p(s)$, $p(s)$ satisfying the following properties.

$$(i) \quad 0 \leq p(s) \leq 1$$

$$(ii) \quad \sum_{s \in \mathcal{S}} p(s) = 1$$

where \mathcal{S} is the totality of all samples that can be drawn from the population. p is called a sampling design. Thus, by a sampling design, we mean a probability measure defined on \mathcal{S} , satisfying (i) and (ii).

⑥ Informative and non-informative design \Rightarrow

A design may depend on the related variables X, Y etc. But we assume that p is free of Y unless explicitly mentioned. To emphasize this freedom, p is often referred in the literature as a non-informative design. On the other hand if p involves any component of Y , it is called an informative design.

⑦ WOR design and WR design \Rightarrow

A design p is WOR design if no repetitions occur in any sample s with $p(s) > 0$, otherwise p is called WR design.

⑧ FS and FES design \Rightarrow

A design p is of fixed size n (fixed effective size n) if $p(s) > 0$ implies that s is of size n (effective size n). With respect to the WOR design, there is of course no difference between FS and FES.

⑨ SRSWOR and SRSWR design:

A design p is called SRSWOR if $p(s) = \begin{cases} \frac{1}{\binom{N}{n}} & \forall s \text{ with } n(s) = n \\ 0, \text{ o.w.} \end{cases}$

and it is called SRSWR if $p(s) = \begin{cases} \frac{1}{N^n} & \forall s \text{ with } n(s) = n \\ 0 & \text{o.w.} \end{cases}$

⑫ Sampling Strategy >>

The combination (p, t) denoting an estimator t based on s chosen according to the sampling design p , is called ^{the} ~~an~~ sampling strategy.

⑬ Inclusion Probability >>

For any arbitrary sampling design p , the probability $\pi_i = \sum_{s \ni i} p(s)$ is called the first order inclusion probability of the unit i .

Again, the probability $\pi_{ij} = \sum_{s \ni \{i, j\}} p(s)$ is called the 2nd order inclusion probability of the units i and j .

$\sum_{i \in s}$:- sum over the units i belonging to s (i varies, s fixed)

$\sum_{s \ni i}$:- sum over ^{the} ~~the~~ samples s containing the unit i (i fixed, s varies)

$s \ni i \rightarrow s$ includes i
 $i \in s \rightarrow i$ belongs to s .

Example

$N = 5 ; n = 3$

SRSWOR

Total no. of possible samples = $\binom{5}{3} = 10$

- | | | | |
|---------------------|---------------------|---------------------|------------------------|
| $s_1 = \{1, 2, 3\}$ | $s_4 = \{1, 3, 4\}$ | $s_7 = \{2, 3, 4\}$ | $s_{10} = \{3, 4, 5\}$ |
| $s_2 = \{1, 2, 4\}$ | $s_5 = \{1, 3, 5\}$ | $s_8 = \{2, 3, 5\}$ | |
| $s_3 = \{1, 2, 5\}$ | $s_6 = \{1, 4, 5\}$ | $s_9 = \{2, 4, 5\}$ | |

$p(s_1) = \dots = p(s_{10}) = \frac{1}{10}$

$$\begin{aligned}\pi_1 &= P(\text{unit 1 is included in the sample}) = \sum_{s \ni 1} p(s) \\ &= p(s_1) + p(s_2) + \dots + p(s_4) \\ &= \frac{6}{10} = \frac{3}{5}\end{aligned}$$

$$\text{Similarly, } \pi_2 = \pi_3 = \pi_4 = \pi_5 = \frac{3}{5}$$

$$\sum_{i=1}^5 \pi_i = 3 = \text{the fixed sample size} = n$$

The inclusion probabilities do not add up to 1.

Expectation & MSE

Whatever Y may be, $E_p(t) = \sum_s p(s) t(s, X)$ is called the expectation of t and

$M_p(t) = E_p [t(s, X) - Y]^2$ is called the MSE of t for estimating Y (the population total)

If $E_p(t(s, X)) = Y$, then t is called a p -unbiased estimator of Y .

Some Results Regarding Inclusion Probabilities:

$$\textcircled{I} \quad E_p(v(s)) = \sum_{i=1}^N \pi_i$$

Proof \gg We know that $\pi_i = \sum_{s \ni i} p(s)$

$$\text{Define, } I_{si} = \begin{cases} 1 & \text{if } s \ni i \\ 0 & \text{o.w.} \end{cases}$$

$$\begin{aligned}\therefore E_p(I_{si}) &= \sum_s p(s) I_{si} \quad \left[\sum_s = \sum_{s \ni i} + \sum_{s \not\ni i} \right] \\ &= \sum_{s \ni i} p(s) I_{si} \quad [\because I_{si} = 0 \text{ for } s \not\ni i] \\ &= \sum_{s \ni i} p(s) \quad [\because I_{si} = 1 \text{ for } s \ni i] \\ &= \pi_i\end{aligned}$$

$$\therefore \sum_{i=1}^N \pi_i = \sum_{i=1}^N E_p(I_{si}) = E_p\left(\sum_{i=1}^N I_{si}\right) = E_p(v(s))$$

$[\because \sum_{i=1}^N I_{si} = v(s) \text{ is the effective size of the sample}]$

Special Case:

For a FES (n) design, $v(\lambda) = n \quad \forall \lambda$

Hence, $\sum_{i=1}^N \pi_i = n$ (Example SRSWOR (n))

$$\textcircled{\text{II}} \quad \sum_{j \neq i} \pi_{ij} = E_p (v(\lambda)) - \pi_i$$

Proof Define $I_{\lambda ij} = \begin{cases} 1 & \text{if } \lambda \ni i \text{ \& } \lambda \ni j \\ 0 & \text{o.w.} \end{cases}$

$$\begin{aligned} E_p (I_{\lambda ij}) &= \sum_{\lambda} p(\lambda) I_{\lambda ij} \\ &= \sum_{\lambda \ni \{i, j\}} p(\lambda) \quad \left[\begin{array}{l} \because I_{\lambda ij} = 0 \text{ for } \lambda \neq \{i, j\} \\ = 1 \text{ for } \lambda \ni \{i, j\} \end{array} \right] \\ &= \pi_{ij} \end{aligned}$$

$$\text{Now } I_{\lambda ij} = \begin{cases} 1 & \text{iff } I_{\lambda i} = 1 \text{ \& } I_{\lambda j} = 1 \\ 0 & \text{o.w.} \end{cases}$$

If $\lambda \ni i$, $I_{\lambda i} = 1$

If $\lambda \ni j$, $I_{\lambda j} = 1$

If $\lambda \ni i$ \& $\lambda \ni j$ $I_{\lambda i} I_{\lambda j} = 1 = I_{\lambda ij}$

$$\therefore \boxed{I_{\lambda ij} = I_{\lambda i} \cdot I_{\lambda j}}$$

Hence, $E_p (I_{\lambda i} I_{\lambda j}) = \pi_{ij}$

or, $\sum_{j \neq i} E_p (I_{\lambda i} \cdot I_{\lambda j}) = \sum_{j \neq i} \pi_{ij}$

or, $E_p \left(\sum_{j \neq i} I_{\lambda i} I_{\lambda j} \right) = \sum_{j \neq i} \pi_{ij}$

or, $E_p \left(I_{\lambda i} \sum_{j \neq i} I_{\lambda j} \right) = \sum_{j \neq i} \pi_{ij}$

or, $E_p \left[I_{\lambda i} (v(\lambda) - I_{\lambda i}) \right] = \sum_{j \neq i} \pi_{ij} \quad \left[\because \sum_i I_{\lambda i} = v(\lambda) \right]$

$$\text{or, } \bullet \quad E_p \left[\sum_{j \neq i} \pi_{ij} v(\lambda_j) - \sum_{j \neq i} \pi_{ij}^2 \right] = \sum_{j \neq i} \pi_{ij}$$

$$\text{or, } \quad E_p \left[v(\lambda) \sum_{j \neq i} \pi_{ij} \right] - E_p \left[\sum_{j \neq i} \pi_{ij}^2 \right] = \sum_{j \neq i} \pi_{ij}$$

$$\text{or, } \quad E_p \left[v(\lambda) \right] - E_p \left[\sum_{j \neq i} \pi_{ij} \right] = \sum_{j \neq i} \pi_{ij}$$

$$\text{or, } \quad \boxed{\sum_{j \neq i} \pi_{ij} = E_p (v(\lambda)) - \pi_i}$$

$$\begin{aligned} \sum_{j \neq i} \pi_{ij} &= \sum_{\lambda} p(\lambda) v(\lambda) \sum_{j \neq i} \pi_{ij} - \pi_i \\ &= \sum_{\lambda} p(\lambda) v(\lambda) - \pi_i \end{aligned}$$

Special Case : $v(\lambda) = n \quad \forall \lambda$

$$\begin{aligned} \text{Then, } \sum_{j \neq i} \pi_{ij} &= n \sum_{\lambda} p(\lambda) - \pi_i \\ &= n \pi_i - \pi_i \\ &= (n-1) \pi_i \end{aligned}$$

$$\textcircled{III} \quad \sum_{i=1}^N \sum_{j \neq i}^N \pi_{ij} = v_p (v(\lambda)) + E_p (v(\lambda)) [E_p (v(\lambda)) - 1]$$

Proof \gg Already proved,

$$\sum_{j \neq i} \pi_{ij} = \sum_{\lambda} p(\lambda) v(\lambda) \sum_{j \neq i} \pi_{ij} - \pi_i$$

$$\begin{aligned} \therefore \sum_{i=1}^N \sum_{j \neq i}^N \pi_{ij} &= \sum_{i=1}^N \left[\sum_{\lambda} p(\lambda) v(\lambda) \sum_{j \neq i} \pi_{ij} \right] - \sum_{i=1}^N \pi_i \\ &= \sum_{\lambda} p(\lambda) v(\lambda) \left[\sum_{i=1}^N \sum_{j \neq i} \pi_{ij} \right] - E_p (v(\lambda)) \\ &= \sum_{\lambda} p(\lambda) \{v(\lambda)\}^2 - E_p (v(\lambda)) \\ &= v_p (v(\lambda)) + [E_p (v(\lambda))]^2 - E_p (v(\lambda)) \end{aligned}$$

$$\therefore \sum_{i=1}^N \sum_{j \neq i} \pi_{ij} = V_p(u(z)) + E_p(u(z)) [E_p(u(z)) - 1]$$

SRSWR (N, n)

$$p(z) = \frac{1}{N^n} \quad \forall z$$

$$\pi_i = \sum_{z \ni i} p(z)$$

$$= \frac{1}{N^n} \sum_{z \ni i} 1$$

$$= \frac{1}{N^n} [\text{No. of samples including } i]$$

$$= \frac{1}{N^n} [\text{Total no. of possible samples} - \text{No. of samples excluding } i]$$

$$= \frac{1}{N^n} [N^n - (N-1)^n]$$

$$= 1 - \left(1 - \frac{1}{N}\right)^n \quad \forall i = 1(1)N$$

$$E_p(u(z)) = \sum_{i=1}^N \pi_i = N \left[1 - \left(1 - \frac{1}{N}\right)^n\right]$$

$$\pi_{ij} = \sum_{z \ni \{i, j\}} p(z) = P(z \ni i \cap z \ni j) = P(A_i \cap A_j)$$

$$= 1 - P(A_i \cap A_j)^c, \text{ where } A_i = \{z \ni i\}$$

$$A_j = \{z \ni j\}$$

$$= 1 - P(A_i^c \cup A_j^c)$$

$$= 1 - P(A_i^c) - P(A_j^c) + P(A_i^c \cap A_j^c)$$

$$P(A_i^c) = 1 - \pi_i = \left(1 - \frac{1}{N}\right)^n$$

$$P(A_j^c) = 1 - \pi_j = \left(1 - \frac{1}{N}\right)^n$$

$$P(A_i^c \cap A_j^c) = \frac{\text{No. of samples excluding both } i \text{ \& } j}{N^n}$$

$$= \frac{(N-2)^n}{N^n} = \left(1 - \frac{2}{N}\right)^n$$

$$\therefore \pi_{ij} = 1 - \left(1 - \frac{1}{N}\right)^n - \left(1 - \frac{1}{N}\right)^n + \left(1 - \frac{2}{N}\right)^n$$

$$= 1 - 2 \left(1 - \frac{1}{N}\right)^n + \left(1 - \frac{2}{N}\right)^n$$

$$\sum_{i \neq j} \pi_{ij} = N(N-1) \left[1 - 2 \left(1 - \frac{1}{N}\right)^n + \left(1 - \frac{2}{N}\right)^n \right]$$

$$V_P(U(\lambda)) + E_P(U(\lambda)) [E_P(U(\lambda)) - 1] = N(N-1) \left[1 - 2 \left(1 - \frac{1}{N}\right)^n + \left(1 - \frac{2}{N}\right)^n \right]$$

$$\therefore V_P(U(\lambda)) = N(N-1) \left[1 - 2 \left(1 - \frac{1}{N}\right)^n + \left(1 - \frac{2}{N}\right)^n \right]$$

$$- N^2 \left[1 - \left(1 - \frac{1}{N}\right)^n \right]^2 + N \left[1 - \left(1 - \frac{1}{N}\right)^n \right]$$

$$= N(N-1) \left[1 - \left(1 - \frac{1}{N}\right)^n \right]^2 + N(N-1) \left[\left(1 - \frac{2}{N}\right)^n - \left(1 - \frac{1}{N}\right)^{2n} \right]$$

$$- N^2 \left[1 - \left(1 - \frac{1}{N}\right)^n \right]^2 + N \left[1 - \left(1 - \frac{1}{N}\right)^n \right]$$

$$= -N \left[1 - \left(1 - \frac{1}{N}\right)^n \right]^2 + N^2 \left[\left(1 - \frac{2}{N}\right)^n - \left(1 - \frac{1}{N}\right)^{2n} \right]$$

$$- N \left[\left(1 - \frac{2}{N}\right)^n - \left(1 - \frac{1}{N}\right)^{2n} \right] + N \left[1 - \left(1 - \frac{1}{N}\right)^n \right]$$

$$= N \left[-1 + 2 \left(1 - \frac{1}{N}\right)^n - \left(1 - \frac{1}{N}\right)^{2n} - \left(1 - \frac{2}{N}\right)^n + \left(1 - \frac{1}{N}\right)^{2n} \right]$$

$$+ N^2 \left[\left(1 - \frac{2}{N}\right)^n - \left(1 - \frac{1}{N}\right)^{2n} \right]$$

$$= N(N-1) \left[\left(1 - \frac{2}{N}\right)^n \right] + N \left(1 - \frac{1}{N}\right)^n - N^2 \left(1 - \frac{1}{N}\right)^{2n}$$

$$= \frac{(N-1)(N-2)^n}{(N-1)^n} + \frac{(N-1)^n}{(N-1)^n} - \frac{(N-1)^{2n}}{(N-1)^{2n}}$$

Result:

Suppose, $t(\lambda, Y) = \sum_{i \in \lambda} b_{\lambda i} Y_i$ be a homogeneous linear estimator unbiased for Y (w.r. to the sampling design p).

Then,
$$\boxed{\sum_{\lambda \ni i} p(\lambda) b_{\lambda i} = 1} \quad \forall i$$

Proof \gg Since, t is unbiased for Y ,

$$E_p(t(\lambda, Y)) = Y$$

$$\text{or, } \sum_{\lambda} p(\lambda) t(\lambda, Y) = Y$$

$$\text{or, } \sum_{\lambda} p(\lambda) \left[\sum_{i \in \lambda} b_{\lambda i} Y_i \right] = Y$$

$$\text{or, } \sum_{i=1}^N Y_i \left[\sum_{\lambda \ni i} p(\lambda) b_{\lambda i} \right] = \sum_{i=1}^N Y_i$$

$$\therefore \sum_{\lambda \ni i} p(\lambda) b_{\lambda i} = 1 \quad \forall i = 1(1)N$$

Previous Example

$N=5, n=3$, SRSWOR

$\binom{5}{3} = 10$ samples

$\lambda_1, \lambda_2, \dots, \lambda_{10}$

$$\begin{aligned} & \sum_{\lambda} p(\lambda) \left[\sum_{i \in \lambda} b_{\lambda i} Y_i \right] \\ &= p(\lambda_1) \left[\sum_{i \in \lambda_1} b_{\lambda_1 i} Y_i \right] + p(\lambda_2) \left[\sum_{i \in \lambda_2} b_{\lambda_2 i} Y_i \right] + \dots + p(\lambda_{10}) \left[\sum_{i \in \lambda_{10}} b_{\lambda_{10} i} Y_i \right] \\ &= p(\lambda_1) [b_{\lambda_1 1} Y_1 + b_{\lambda_1 2} Y_2 + b_{\lambda_1 3} Y_3] + p(\lambda_2) [b_{\lambda_2 1} Y_1 + b_{\lambda_2 2} Y_2 + b_{\lambda_2 4} Y_4] \\ &+ p(\lambda_3) [b_{\lambda_3 1} Y_1 + b_{\lambda_3 2} Y_2 + b_{\lambda_3 5} Y_5] + p(\lambda_4) [b_{\lambda_4 1} Y_1 + b_{\lambda_4 3} Y_3 + b_{\lambda_4 4} Y_4] \\ &+ p(\lambda_5) [b_{\lambda_5 1} Y_1 + b_{\lambda_5 3} Y_3 + b_{\lambda_5 5} Y_5] + p(\lambda_6) [b_{\lambda_6 1} Y_1 + b_{\lambda_6 4} Y_4 + b_{\lambda_6 5} Y_5] \\ &+ p(\lambda_7) [b_{\lambda_7 2} Y_2 + b_{\lambda_7 3} Y_3 + b_{\lambda_7 4} Y_4] + p(\lambda_8) [b_{\lambda_8 2} Y_2 + b_{\lambda_8 3} Y_3 + b_{\lambda_8 5} Y_5] \\ &+ p(\lambda_9) [b_{\lambda_9 2} Y_2 + b_{\lambda_9 4} Y_4 + b_{\lambda_9 5} Y_5] + p(\lambda_{10}) [b_{\lambda_{10} 3} Y_3 + b_{\lambda_{10} 4} Y_4 + b_{\lambda_{10} 5} Y_5] \end{aligned}$$

ator

$$= \sum_{i=1}^{10} Y_i$$

$$= Y_1 [b_{\lambda 11} p(\lambda_1) + b_{\lambda 21} p(\lambda_2) + b_{\lambda 31} p(\lambda_3) + b_{\lambda 41} p(\lambda_4) + b_{\lambda 51} p(\lambda_5) + b_{\lambda 61} p(\lambda_6)] +$$

$$Y_2 [b_{\lambda 12} p(\lambda_1) + b_{\lambda 22} p(\lambda_2) + b_{\lambda 32} p(\lambda_3) + b_{\lambda 42} p(\lambda_4) + b_{\lambda 52} p(\lambda_5) + b_{\lambda 62} p(\lambda_6)] +$$

$$Y_3 [b_{\lambda 13} p(\lambda_1) + b_{\lambda 23} p(\lambda_2) + b_{\lambda 33} p(\lambda_3) + b_{\lambda 43} p(\lambda_4) + b_{\lambda 53} p(\lambda_5) + b_{\lambda 63} p(\lambda_6)] +$$

$$Y_4 [b_{\lambda 14} p(\lambda_1) + b_{\lambda 24} p(\lambda_2) + b_{\lambda 34} p(\lambda_3) + b_{\lambda 44} p(\lambda_4) + b_{\lambda 54} p(\lambda_5) + b_{\lambda 64} p(\lambda_6)] +$$

$$Y_5 [b_{\lambda 15} p(\lambda_1) + b_{\lambda 25} p(\lambda_2) + b_{\lambda 35} p(\lambda_3) + b_{\lambda 45} p(\lambda_4) + b_{\lambda 55} p(\lambda_5) + b_{\lambda 65} p(\lambda_6)]$$

$$= \sum_{i=1}^N Y_i [\sum_{\lambda \in \Omega} p(\lambda) b_{\lambda i}]$$

$$= Y_1 [\sum_{\lambda \in \Omega} p(\lambda) b_{\lambda 1}] + Y_2 [\sum_{\lambda \in \Omega} p(\lambda) b_{\lambda 2}] + \dots +$$

$$Y_5 [\sum_{\lambda \in \Omega} p(\lambda) b_{\lambda 5}]$$

$$= \sum_{i=1}^5 Y_i [\sum_{\lambda \in \Omega} p(\lambda) b_{\lambda i}]$$

Thus, for a HLUB of the population total,

$$\sum_{\lambda \in \Omega} b_{\lambda i} p(\lambda) = 1 \quad \forall i = 1(1)N$$

]

Y₄

Y₄

Y₅

Y₅

Design-Based Inference

The MSE of an estimator t of the population total Y is

$$\begin{aligned}M_p(t) &= E_p \left(t(\lambda, \mathcal{X}) - Y \right)^2 \\ &= \sum_{\lambda} p(\lambda) \left(t(\lambda, \mathcal{X}) - Y \right)^2\end{aligned}$$

Let Σ_1 and Σ_2 denote respectively the sum over the samples λ for which $|t(\lambda, \mathcal{X}) - Y| \geq k > 0$ and $|t(\lambda, \mathcal{X}) - Y| < k$

$$\text{i.e. } \Sigma_1 \equiv \sum_{\lambda: |t(\lambda, \mathcal{X}) - Y| \geq k}$$

$$\& \Sigma_2 \equiv \sum_{\lambda: |t(\lambda, \mathcal{X}) - Y| < k}$$

Hence,

$$\begin{aligned}M_p(t) &= \Sigma_1 p(\lambda) \left(t(\lambda, \mathcal{X}) - Y \right)^2 + \Sigma_2 p(\lambda) \left(t(\lambda, \mathcal{X}) - Y \right)^2 \\ &\geq \Sigma_1 p(\lambda) \left(t(\lambda, \mathcal{X}) - Y \right)^2 \\ &\geq k^2 \Sigma_1 p(\lambda) \quad \left[\because \text{for } \lambda \in \Sigma_1 \text{ we have such that } |t(\lambda, \mathcal{X}) - Y| \geq k \right] \\ &= k^2 \sum_{\lambda: |t(\lambda, \mathcal{X}) - Y| \geq k} p(\lambda) \\ &= k^2 P \left(|t(\lambda, \mathcal{X}) - Y| \geq k \right)\end{aligned}$$

$$\therefore P \left(|t(\lambda, \mathcal{X}) - Y| \geq k \right) \leq \frac{M_p(t)}{k^2}$$

$$\text{or, } P \left(|t(\lambda, \mathcal{X}) - Y| \leq k \right) \geq 1 - \frac{M_p(t)}{k^2}$$

$$\text{or, } P \left(t(\lambda, \mathcal{X}) - k \leq Y \leq t(\lambda, \mathcal{X}) + k \right) \geq 1 - \frac{M_p(t)}{k^2}, \quad k > 0$$

This alternative version of Chebyshev's inequality is due to Neyman.

The coverage probability is inversely proportional to $M_p(t)$. Thus in order to increase the coverage probability, we need to reduce both the bias and the variance of t .

Special Case: $k = 3\sqrt{V_p(t)}$

$$\begin{aligned} P\left(t(z, \alpha) - 3\sqrt{V_p(t)} \leq Y \leq t(z, \alpha) + 3\sqrt{V_p(t)}\right) &\geq 1 - \frac{M_p(t)}{9V_p(t)} \\ &= 1 - \frac{[V_p(t) + \{B_p(t)\}^2]}{9V_p(t)} \\ &= \frac{8}{9} - \frac{\{B_p(t)\}^2}{9V_p(t)} \end{aligned}$$

Illustration:

Suppose, Y is a binary random variables, taking values 0 and 1. We consider estimation of \bar{Y} .

Here, Y is binary, hence $\sum_{i=1}^N Y_i^2 = \sum_{i=1}^N Y_i$ & $\sigma_Y^2 = \frac{1}{N} \sum_{i=1}^N Y_i^2 - \bar{Y}^2$

$$\begin{aligned} &= \frac{1}{N} \sum_{i=1}^N Y_i - \bar{Y}^2 \\ &= \bar{Y} - \bar{Y}^2 = \bar{Y}(1-\bar{Y}) \end{aligned}$$

If we consider SRSWR (n);

$$\begin{aligned} E_p(\bar{y}) &= \bar{Y} \\ \& \quad V_p(\bar{y}) &= \frac{\sigma_Y^2}{n} = \frac{\bar{Y}(1-\bar{Y})}{n} \\ \& \quad B_p(\bar{y}) &= 0 \end{aligned}$$

$$\begin{aligned} P\left(t(z, \alpha) - 3\sqrt{V_p(t)} \leq \bar{Y} \leq t(z, \alpha) + 3\sqrt{V_p(t)}\right) \\ \geq \frac{8}{9} - \frac{\{B_p(t)\}^2}{9V_p(t)} \end{aligned}$$

Here, $t = \bar{y}$ & $B_p(t) = 0$

$$\therefore P\left(\bar{y} - 3\sqrt{\frac{\bar{Y}(1-\bar{Y})}{n}} \leq \bar{Y} \leq \bar{y} + 3\sqrt{\frac{\bar{Y}(1-\bar{Y})}{n}}\right) \geq \frac{8}{9}$$

$V_p(t)$ is maximum at $\bar{Y} = \frac{1}{2}$ & $\text{Max } V_p(t) = \frac{1}{4n}$

Thus, replacing $V_p(t)$ by its max value, i.e. $\frac{1}{4n}$, we have,

$$P\left(\bar{y} - \frac{3}{2\sqrt{n}} \leq \bar{Y} \leq \bar{y} + \frac{3}{2\sqrt{n}}\right) \geq \frac{8}{9}$$

□ Comparison between the estimators:

The estimator t of Y is said to be better than another estimator t' of Y if $E_p(t) = E_p(t') = Y$

$$\& \quad V_p(t) \leq V_p(t') \quad \forall \underline{y} \in \Omega \subseteq \mathbb{R}^N$$

where, $\Omega = \left\{ \underline{y} \mid -\infty < a_i < Y_i < b_i < \infty, i = 1, \dots, N \right\}$

We denote by the symbol $\left\{ t \succ t' \right\}$ that t is better than t'
 (better than)

An estimator $t_0 = t_0(\lambda, \underline{y})$ is called the best for Y if no other estimator t is better than t_0 .

Theorem I:

No best unbiased estimator of Y exists based on any non-census design.

Proof \gg A census design is defined as

$$p(\lambda) = \begin{cases} 1, & \text{if } \lambda \text{ includes every element of } U \\ 0, & \text{otherwise} \end{cases}$$

Any design other than a census design is called non-census design.

$$\Omega = \left\{ \underline{y} : -\infty < a_i < Y_i < b_i < \infty; i = 1, \dots, N \right\}$$

If possible, let \exists some unbiased estimator $t = t(\lambda, X)$ of Y based on the sampling design p , with the least possible variance.

Let us choose a point $A = (A_1, A_2, \dots, A_N) \in \Omega$ and define

$$A = \sum_{i=1}^N A_i$$

Define a new estimator $t_A(\lambda, X) = t(\lambda, X) - t(\lambda, A) + A$

Here, $E_p(t_A(\lambda, X)) = E_p(t(\lambda, X)) - E_p(t(\lambda, A)) + A$

$$= Y - Y + A$$

$$= Y$$

Thus, $t_A(\lambda, X)$ is another unbiased estimator of Y .

Again,

$$V_p(t_A(\lambda, X)) = E_p[t_A(\lambda, X) - Y]^2$$

$$= E_p[t(\lambda, X) - t(\lambda, A) + A - Y]^2$$

$$= E_p[\{t(\lambda, X) - Y\} - \{t(\lambda, A) - A\}]^2$$

$$= E_p[t(\lambda, X) - Y]^2 +$$

$$E_p[t(\lambda, A) - A]^2$$

$$- 2 E_p[\{t(\lambda, X) - Y\}\{t(\lambda, A) - A\}]$$

Note that,

$$V_p(t_A(\lambda, X)) \Big|_{Y=A} = 0$$

But since the point A is arbitrarily chosen, it follows that

$$V_p(t_A(\lambda, X)) = 0$$

$$\Rightarrow V_p(t(\lambda, X)) = 0$$

$$\Rightarrow t(\lambda, X) = E_p(t(\lambda, X)) \text{ with probability } 1$$

$$\Rightarrow t(\lambda, X) = Y \text{ w.p. } 1.$$

This is not possible for any ~~non-census~~ non-census design.

Hence, no best n.e. exists for any non-census design.

Theorem 2 :

A necessary and sufficient condition for the existence of an u.e. of the survey population total is

$$\pi_i > 0 \quad \forall i = 1(1)N.$$

Proof >> If Part

Suppose, $\pi_i > 0 \quad \forall i = 1(1)N$

Then, $\frac{y_i}{\pi_i}$ is defined $\forall i = 1(1)N$

Define
$$I_{\lambda i} = \begin{cases} 1 & , \text{ if } \lambda \ni i \\ 0 & , \text{ o.w.} \end{cases}$$

$$\therefore E(I_{\lambda i}) = \sum_{\lambda \ni i} p(\lambda) = \pi_i$$

$$\therefore E(I_{\lambda i}) = \sum_{\lambda} p(\lambda) I_{\lambda i} = \sum_{\lambda \ni i} p(\lambda) I_{\lambda i} + \sum_{\lambda \not\ni i} p(\lambda) I_{\lambda i}$$

$$= \sum_{\lambda \ni i} p(\lambda) = \pi_i$$

Define,
$$t = \sum_{i=1}^N \frac{y_i}{\pi_i} \cdot I_{\lambda i}$$

$$E_p(t) = \sum_{i=1}^N \frac{y_i}{\pi_i} E(I_{\lambda i}) = \sum_{i=1}^N \frac{y_i}{\pi_i} \cdot \pi_i = \sum_{i=1}^N y_i = Y$$

\therefore There exists an u.e. t of Y .

Only If Part

Suppose, $\pi_i = 0$ for some $i \in \{1, 2, \dots, N\}$ & ~~there exists some~~
~~u.e.~~ let, $t = t(\lambda, Y)$ be an u.e. of the population total Y .

$$\therefore \pi_i = \sum_{\lambda \ni i} p(\lambda) = 0 \quad \text{--- (1)}$$

$$E_p(t) = Y \Rightarrow \sum_{\lambda} p(\lambda) t(\lambda, Y) = \sum_{i=1}^N y_i \quad \text{--- (2)}$$

From (1), it follows that $p(\lambda) > 0 \quad \forall \lambda \in I$ since $p(\lambda)$ are non-negative quantities.

From (2),
$$\sum_{\lambda \in I} p(\lambda) t(\lambda, X) + \sum_{\lambda \notin I} p(\lambda) t(\lambda, X) = \sum_{i=1}^N \gamma_i$$

$$\Rightarrow \sum_{\lambda \notin I} p(\lambda) t(\lambda, X) = \sum_{i=1}^N \gamma_i \quad [\because p(\lambda) = 0 \quad \forall \lambda \in I]$$

This is a contradiction, since L.H.S. doesn't involve γ_i but R.H.S. involves γ_i .

Hence, $\pi_i > 0 \quad \forall i = 1(1)N$.

Theorem 3.1 (Due to Godambe)

Among the class of H.L.U.E.s, no ~~one~~^{one} exists with uniformly minimum variance for any usual design.

Proof \gg Let, $t = t(\lambda, X) = \sum_{i \in \lambda} b_{\lambda i} \gamma_i$ be a H.L.U.E. of population total

Since, t is unbiased for γ , it follows that
$$\sum_{\lambda \in I} p(\lambda) b_{\lambda i} = 1 \quad (*)$$

 $\forall i = 1(1)N$

$$\begin{aligned} V_p(t) &= E_p(t^2(\lambda, X)) - [E_p(t(\lambda, X))]^2 \\ &= \sum_{\lambda} p(\lambda) t^2(\lambda, X) - \gamma^2 \\ &= \sum_{\lambda} p(\lambda) \left(\sum_{i \in \lambda} b_{\lambda i} \gamma_i \right)^2 - \gamma^2 \end{aligned}$$

We minimize $V_p(t)$ subject to the constraints (*)

The Lagrangian function we minimize, is

$$V^* = V_p(t) - \sum_{i=1}^N \lambda_i \left(\sum_{\lambda \in I} p(\lambda) b_{\lambda i} - 1 \right), \text{ where,}$$

$\lambda_1, \lambda_2, \dots, \lambda_N$ are the Lagrangian multipliers.

The normal equations
$$\frac{\partial V^*}{\partial b_{\lambda i}} = 0, \quad i = 1(1)N$$

$$2 p(\lambda) \left(\sum_{i \in \lambda} b_{\lambda i} \gamma_i \right) \gamma_i - \lambda_i p(\lambda) = 0, \quad i = 1(1)N$$

$$\Rightarrow \lambda_i = 2 \left(\sum_{i \in \lambda} b_{\lambda i} \gamma_i \right) \gamma_i \Rightarrow \sum_{i \in \lambda} b_{\lambda i} \gamma_i = \frac{\lambda_i}{2}$$

If we assume that $\gamma_i \neq 0$ and $\gamma_j = 0 \forall j \neq i$,

$$b_{xi} \gamma_i = \frac{\lambda_i}{2\gamma_i}$$

$$\Rightarrow b_{xi} = \frac{\lambda_i}{2\gamma_i^2} \text{ (free of } \lambda)$$

Putting the value of b_{xi} in (*),

$$\sum_{\lambda \in \Omega} b_{xi} p(\lambda) = 1$$

$$\Rightarrow \sum_{\lambda \in \Omega} \frac{\lambda_i}{2\gamma_i^2} p(\lambda) = 1 \quad \forall i = 1(1)N$$

$$\therefore \frac{\pi_i \lambda_i}{2\gamma_i^2} = 1 \quad \forall i = 1(1)N$$

$$\Rightarrow \frac{\lambda_i}{2\gamma_i^2} = \frac{1}{\pi_i} \quad \forall i = 1(1)N$$

$$\therefore \sum_{i \in \Omega} \frac{\gamma_i}{\pi_i} = \sum_{i \in \Omega} \gamma_i \frac{\lambda_i}{2\gamma_i^2} = \sum_{i \in \Omega} \frac{\lambda_i}{2\gamma_i} = \frac{\lambda_i}{2\gamma_i} [\because \gamma_j = 0 \forall j \neq i]$$

So, if $\lambda_1 \in \Omega$, $\lambda_2 \in \Omega \Rightarrow p(\lambda_1) > 0$ & $p(\lambda_2) > 0$

then

$$\sum_{i \in \Omega_1} \frac{\gamma_i}{\pi_i} = \sum_{i \in \Omega_2} \frac{\gamma_i}{\pi_i} \text{ for the existence of UMVUE in the class}$$

of HLUUE. This can't be designed unless the design p satisfies the conditions that for λ_1, λ_2 such that $p(\lambda_1) > 0, p(\lambda_2) > 0$, either (i) $\lambda_1 \cap \lambda_2 = \emptyset$ or (ii) $\lambda_1 \sim \lambda_2$ (equivalent). λ_1 and λ_2 are equivalent in the sense that both contain an identical set of distinct elements of U . A design p satisfying the property is called 'uni-cluster design' (UCD). For example, the design corresponding to systematic sampling, any design other than this non-uni-cluster design (NUCD) (not satisfying these stringent conditions). For an UCD it is possible to realize

$$\sum_{i \in \Omega_1} \frac{\gamma_i}{\pi_i} = \sum_{i \in \Omega_2} \frac{\gamma_i}{\pi_i}, \text{ uniformly in } \gamma \text{ but not for NUCD.}$$

Hence for a NUED no. UMVUE exists.

Rao-Blackwellization:

An estimator $t = t(s, Y)$ may depend on either the order in which the units appear in s and may depend on the multiplicity of the appearance of units in s .

Example: 1 Let, P_i [$0 < P_i < 1$, $\sum_{i=1}^N P_i = 1$] be known numbers associated with the unit i of U and on the first draw a unit i be chosen with probability P_i and on the second draw from U (~~leaving~~ ^{leaving} aside unit i if chosen in the 1st draw), a unit j ($\neq i$) be chosen with probability $\frac{P_j}{1-P_i}$. Consider

$$t = t(i, j) = \frac{1}{2} \left[\frac{Y_i}{P_i} + \frac{Y_i + Y_j}{P_j / (1-P_i)} \right]. \text{ It can be shown that}$$

$$E_p \left[\frac{Y_i}{P_i} \right] = Y \text{ and } E_p \left[Y_i + \frac{Y_j}{P_j / (1-P_i)} \right] = Y$$

and hence, $E_p(t(i, j)) = Y$ but $t(i, j) \neq t(j, i)$
(Desh Raj's Estimator)

Example: 2 Consider the estimator $t = \frac{1}{n} \sum_{i=1}^N \frac{Y_i}{P_i} f_{si}$ (Hansen & Hurwitz's Estimator), where f_{si} is the multiplicity of the unit i in the sample s . This estimator depends upon the multiplicity of units.

With an arbitrary sample $s = \{i_1, i_2, \dots, i_n\}$ we associate the sample $s^* = \{j_1, j_2, \dots, j_k\}$ which is equivalent to s (denoted by $s \sim s^*$) and satisfies $j_1 < j_2 < \dots < j_k$. Hence, all the units of s are listed in s^* , without repetitions and without taking into account the order in which they appear. So s^* may be called an unordered sample and should be identified without an equivalence class of ordered samples for which repetitions are allowed.

Example: $\lambda_1 = \{3, 3, 1, 5, 5\}$; $d_1 = \{(3, 75), (3, 75), (1, 50), (5, 100), (5, 100)\}$

$\lambda_2 = \{1, 3, 1, 1, 5\}$; $d_2 = \{(1, 50), (3, 75), (1, 50), (1, 50), (5, 100)\}$

$\lambda_3 = \{5, 3, 3, 1, 5\}$; $d_3 = \{(5, 100), (3, 75), (3, 75), (1, 50), (5, 100)\}$

Thus, $\lambda^* = \{1, 3, 5\}$

$d^* = \{(1, 50), (3, 75), (5, 100)\}$

Let, $\Omega =$ Parametric space

Consider any design p , yielding the survey data $d = \{(i, Y_i) : i \in \lambda\}$, compatible with the subset $\Omega_d = \{\lambda \in \Omega \mid Y_i \text{ is observed for } i \in \lambda\}$ of the parametric space.

Define,
$$I_d(\lambda) = \begin{cases} 1 & \text{if } \lambda \in \Omega_d \\ 0 & \text{if } \lambda \notin \Omega_d \end{cases}$$

Then the likelihood of λ given d is proportional to the probability of observing d when λ is the underlying parametric point ~~ie.~~ ie.

$$L_d(\lambda) = p(\lambda) I_d(\lambda) = P_\lambda(d) \text{ (say)}$$

The reduced data is $d^* = \{i; Y_i \mid i \in \lambda^*\}$

Then $\forall d$,

$$I_d(\lambda) = I_{d^*}(\lambda) \quad \& \quad L_d(\lambda) = p(\lambda^*) I_{d^*}(\lambda) = P_\lambda(d^*)$$

For simplicity, we write $P(d)$ instead of $P_\lambda(d)$

$$\frac{P(d)}{P(d^*)} = \frac{P(d \cap d^*)}{P(d^*)} = P(d | d^*)$$

$$\frac{p(\lambda) I_d(\lambda)}{p(\lambda^*) I_{d^*}(\lambda)} = P(d | d^*)$$

$$\Rightarrow P(d | d^*) = \frac{p(\lambda)}{p(\lambda^*)} \text{ free of } d^* (= \text{constant})$$

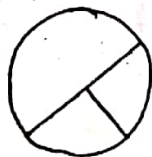
So, d^* is a sufficient statistic, given d .

Let, $\mathcal{D} = \{d\} = \text{Data space}$

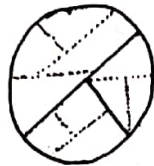
$d = \text{A specific data point}$

$t = t(d) = t(x, X) = \text{A statistic}$

Every statistic induces a partition on the data space. Partition set $t(d)$ is the collection of sets which are mutually exclusive and exhaustive. A pair of statistics t_1 and t_2 may be compared if the corresponding partitions are such that every partition set induced by t_1 is a union of more than one partition sets induced by t_2 (which means induced by t_2 is contained entirely within one of the partition sets induced by t_1). Then we say that t_1 induces a thicker partition than t_2 .



t_1



t_2

If both t_1 and t_2 are sufficient statistics, then we say that t_1 induces a more effective summarization of the data than t_2 , while neither sacrifices any relevant information. A sufficient statistic that induces thickest partitioning is called a minimal sufficient statistic and summarization of data is to the maximal extent, sacrificing no relevant information contained in the data.

Theorem

Given a sample $s = \{i_1, i_2, \dots, i_n\}$ based on a non-informative design p and the data point $d = \{(i_1, Y_{i_1}), \dots, (i_n, Y_{i_n})\}$, the statistic $d^* = \{(j_1, Y_{j_1}), (j_2, Y_{j_2}), \dots, (j_k, Y_{j_k})\}$ derived from d through the sample set $s^* = \{j_1, j_2, \dots, j_k\}$ obtained from s by subprasing the order and for multiplicity of the units of U , is the minimal sufficient statistic.

Construction of a new estimator

Let, $t = t(d) = t(\lambda, \mathcal{X})$ be a statistic based on d which is an unbiased estimator of γ .

λ^* be the set to which one or more samples correspond ($\lambda \sim \lambda^*$)

$\sum_{\lambda \sim \lambda^*} p(\lambda) =$ Probability of observing $\lambda^* = p^*(\lambda^*)$

Then the estimator $t^*(\lambda^*, \mathcal{X}) = \frac{\sum_{\lambda \sim \lambda^*} p(\lambda) t(\lambda, \mathcal{X})}{\sum_{\lambda \sim \lambda^*} p(\lambda)}$ is also

unbiased for γ & $V_p(t^*) \leq V_p(t)$

Proof \Rightarrow

$$\begin{aligned} E_p(t) &= \sum_{\lambda} p(\lambda) t(\lambda, \mathcal{X}) \\ &= \sum_{\lambda^*} \sum_{\lambda \sim \lambda^*} p(\lambda) t(\lambda, \mathcal{X}) \\ &= \sum_{\lambda^*} \left[t^*(\lambda^*, \mathcal{X}) \sum_{\lambda \sim \lambda^*} p(\lambda) \right] \\ &= \sum_{\lambda^*} \left[t^*(\lambda^*, \mathcal{X}) \cdot p^*(\lambda^*) \right] \\ &= E_{p^*} \left(t^*(\lambda^*, \mathcal{X}) \right) \quad \text{--- (1)} \\ &= \sum_{\lambda^*} \sum_{\lambda \sim \lambda^*} t^*(\lambda^*, \mathcal{X}) p(\lambda) \end{aligned}$$

For the samples λ with $\lambda \sim \lambda^*$,
 $t^*(\lambda^*, \mathcal{X}) = t^*(\lambda, \mathcal{X})$

$$\begin{aligned} \therefore E_p(t) &= \sum_{\lambda^*} \sum_{\lambda \sim \lambda^*} t^*(\lambda, \mathcal{X}) p(\lambda) \\ &= \sum_{\lambda} p(\lambda) t^*(\lambda, \mathcal{X}) \\ &= E_p(t^*) \end{aligned}$$

Thus, t^* is also unbiased for γ w.r. to p .

Again, $E_p(t t^*) = \sum_{\lambda} t(\lambda, \chi) t^*(\lambda^*, \chi) p(\lambda)$

$$= \sum_{\lambda^*} \sum_{\lambda \sim \lambda^*} t(\lambda, \chi) t^*(\lambda^*, \chi) p(\lambda)$$

$\lambda \sim \lambda^*$

$$= \sum_{\lambda^*} \left[\sum_{\lambda \sim \lambda^*} t(\lambda, \chi) p(\lambda) \right] t^*(\lambda^*, \chi)$$

~~$$= \sum_{\lambda^*} t^{*2}(\lambda^*, \chi) \sum_{\lambda \sim \lambda^*} p(\lambda)$$~~

$$= \sum_{\lambda^*} t^{*2}(\lambda^*, \chi) p^*(\lambda^*) \left[\because t^*(\lambda^*, \chi) = \frac{\sum_{\lambda \sim \lambda^*} t(\lambda, \chi) p(\lambda)}{p^*(\lambda^*)} \right]$$

$$= \sum_{\lambda^*} t^{*2}(\lambda^*, \chi) \left(\sum_{\lambda \sim \lambda^*} p(\lambda) \right)$$

$$= \sum_{\lambda^*} \sum_{\lambda \sim \lambda^*} t^{*2}(\lambda^*, \chi) p(\lambda)$$

$$= \sum_{\lambda^*} \sum_{\lambda \sim \lambda^*} t^{*2}(\lambda, \chi) p(\lambda) \left[\because \text{for } \lambda \sim \lambda^*, t^*(\lambda^*, \chi) = t^*(\lambda, \chi) \right]$$

$$= \sum_{\lambda} t^{*2}(\lambda, \chi) p(\lambda)$$

$$= E_p [t^{*2}] \quad \text{--- (2)}$$

Now, we know that $E_p (t - t^*)^2 \geq 0$

$$\text{or, } E_p(t^2) + E_p(t^{*2}) - 2E_p(t t^*) \geq 0$$

$$\text{or, } E_p(t^2) + E_p(t^{*2}) - 2E_p(t^{*2}) \geq 0 \quad [\text{from (2)}]$$

$$\text{or, } E_p(t^2) - E_p(t^{*2}) \geq 0$$

$$\text{or, } E_p(t^2) \geq E_p(t^{*2})$$

$$\text{or, } E_p(t^2) - \{E_p(t)\}^2 \geq E_p(t^{*2}) - \{E_p(t^*)\}^2$$

$$[\because E_p(t) = E_p(t^*)]$$

$$\text{or, } V_p(t) \geq V_p(t^*)$$

Horvitz-Thompson's Estimator of Population Total

Definition:

$$t_{HT} = \sum_{i=1}^N \frac{Y_i}{\pi_i} I_{xi} = \sum_{i \in S} \frac{Y_i}{\pi_i}$$

$$E_p [t_{HT}] = \sum_{i=1}^N \frac{Y_i}{\pi_i} E_p (I_{xi}) = \sum_{i=1}^N \frac{Y_i}{\pi_i} \cdot \pi_i \quad [\because E_p (I_{xi}) = \pi_i \text{ already proved}]$$

$$= \sum_{i=1}^N Y_i = Y$$

$$V_p (t_{HT}) = V_p \left[\sum_{i=1}^N \frac{Y_i}{\pi_i} I_{xi} \right] = \sum_{i=1}^N \frac{Y_i^2}{\pi_i^2} V_p (I_{xi}) + \sum_{i \neq j} \frac{Y_i Y_j}{\pi_i \pi_j} + \sum_{i \neq j} \frac{Y_i Y_j}{\pi_i \pi_j} \text{Cov}_p (I_{xi}, I_{xj})$$

$$\text{Now, } V_p (I_{xi}) = E_p (I_{xi}^2) - [E_p (I_{xi})]^2$$

$$= E_p (I_{xi}) - [E_p (I_{xi})]^2 \quad [\because I_{xi} \text{ is an indicator variable}]$$

$$= \pi_i - \pi_i^2 = \pi_i (1 - \pi_i)$$

$$\text{Cov}_p (I_{xi}, I_{xj}) = E_p (I_{xi} \cdot I_{xj}) - E_p (I_{xi}) \cdot E_p (I_{xj})$$

$$= E_p (I_{xij}) - E_p (I_{xi}) E_p (I_{xj}) \quad [\because I_{xij} = I_{xi} \cdot I_{xj}]$$

$$= \pi_{ij} - \pi_i \pi_j$$

$$\therefore V_p (t_{HT}) = \sum_{i=1}^N \frac{Y_i^2}{\pi_i^2} \cdot \pi_i (1 - \pi_i) + \sum_{i \neq j} \frac{Y_i Y_j}{\pi_i \pi_j} (\pi_{ij} - \pi_i \pi_j)$$

$$= \sum_{i=1}^N \frac{Y_i^2}{\pi_i} (1 - \pi_i) + \sum_{i \neq j} \frac{Y_i Y_j}{\pi_i \pi_j} (\pi_{ij} - \pi_i \pi_j) E_p (I_{xi})$$

$$= E_p \left[\sum_{i=1}^N \frac{Y_i^2}{\pi_i^2} (1 - \pi_i) I_{xi} + \sum_{i \neq j} \frac{Y_i Y_j}{\pi_i \pi_j} \left(\frac{\pi_{ij} - \pi_i \pi_j}{\pi_j} \right) I_{xij} \right]$$

Thus, the unbiased estimator of $V_p(t_{HT})$ is

$$\sum_{i \in \lambda} \frac{y_i^2}{\pi_i^2} (1 - \pi_i) I_{\lambda} + \sum_{i \neq j} \frac{y_i y_j}{\pi_i \pi_j} \cdot \frac{(\pi_{ij} - \pi_i \pi_j)}{\pi_{ij}} I_{\lambda ij}$$

$$= \sum_{i \in \lambda} \frac{y_i^2}{\pi_i^2} (1 - \pi_i) + \sum_{i \neq j \in \lambda} \frac{y_i y_j}{\pi_i \pi_j} \cdot \frac{(\pi_{ij} - \pi_i \pi_j)}{\pi_{ij}}$$

An alternative form of $V_p(t_{HT})$ by Sen, Yates & Grundy

If $\nu(\lambda) = n$ & $\pi_{ij} > 0 \forall i, j$

$$V_p(t_{HT}) = \sum_{i < j} (\pi_i \pi_j - \pi_{ij}) \left(\frac{y_i}{\pi_i} - \frac{y_j}{\pi_j} \right)^2$$

Proof $\gg \sum_{i < j} (\pi_i \pi_j - \pi_{ij}) \left(\frac{y_i}{\pi_i} - \frac{y_j}{\pi_j} \right)^2$

$$= \sum_{i < j} (\pi_i \pi_j - \pi_{ij}) \left(\frac{y_i^2}{\pi_i^2} + 2 \frac{y_i y_j}{\pi_i \pi_j} + \frac{y_j^2}{\pi_j^2} \right)$$

$$= \sum_{i < j} \pi_j \frac{y_i^2}{\pi_i^2} + \sum_{i < j} \pi_i \frac{y_j^2}{\pi_j^2} - 2 \sum_{i < j} y_i y_j$$

$$- \sum_{i < j} \pi_{ij} \frac{y_i^2}{\pi_i^2} + 2 \sum_{i < j} \frac{y_i y_j}{\pi_i \pi_j} \pi_{ij} - \sum_{i < j} \pi_{ij} \frac{y_j^2}{\pi_j^2}$$

$$= \sum_{i=1}^N \frac{y_i^2}{\pi_i^2} \left(\sum_{j>i} \pi_j \right) + \sum_{j=1}^N \frac{y_j^2}{\pi_j^2} \left(\sum_{i<j} \pi_i \right) - \sum_{i=1}^N \frac{y_i^2}{\pi_i^2} \left(\sum_{j>i} \pi_{ij} \right)$$

$$- \sum_{j=1}^N \frac{y_j^2}{\pi_j^2} \left(\sum_{i<j} \pi_{ij} \right) + 2 \sum_{i < j} \frac{y_i y_j}{\pi_i \pi_j} (\pi_{ij} - \pi_i \pi_j)$$

Interchanging the role of i and j in term 2 and 4 to have

$$\text{LHS} = \sum_{i=1}^N \frac{y_i^2}{\pi_i^2} \left(\sum_{j>i} \pi_j \right) + \sum_{i=1}^N \frac{y_i^2}{\pi_i^2} \left(\sum_{j<i} \pi_j \right) - \sum_{i=1}^N \frac{y_i^2}{\pi_i^2} \left(\sum_{j>i} \pi_{ij} \right)$$

$$- \sum_{i=1}^N \frac{y_i^2}{\pi_i^2} \left(\sum_{j<i} \pi_{ji} \right) + 2 \sum_{i < j} \frac{y_i y_j}{\pi_i \pi_j} (\pi_{ij} - \pi_i \pi_j)$$

$$= \sum_{i=1}^N \frac{y_i^2}{\pi_i^2} \left(\sum_{j>i} \pi_j + \sum_{j<i} \pi_j \right) - \sum_{i=1}^N \frac{y_i^2}{\pi_i^2} \left(\sum_{j>i} \pi_{ij} + \sum_{j<i} \pi_{ji} \right)$$

$$+ 2 \sum_{i < j} \frac{y_i y_j}{\pi_i \pi_j} (\pi_{ij} - \pi_i \pi_j)$$

$$= \sum_{i=1}^N \frac{Y_i^2}{\pi_i} \left(\sum_{j \neq i} \pi_j \right) - \sum_{i=1}^N \frac{Y_i^2}{\pi_i^2} \left(\sum_{j \neq i} \pi_{ij} \right) + 2 \sum_{i < j} \frac{Y_i Y_j}{\pi_i \pi_j} (\pi_{ij} - \pi_i \pi_j)$$

Since, the design is fixed effective size n , hence $\sum_{j=1}^N \pi_j = n$
(FES (n))

$$\sum_{j \neq i} \pi_j = n - \pi_i \quad \& \quad \sum_{j \neq i} \pi_{ij} = (n-1) \pi_i$$

$$\text{LHS} = \sum_{i=1}^N \frac{Y_i^2}{\pi_i} (n - \pi_i) - \sum_{i=1}^N \frac{Y_i^2}{\pi_i^2} (n-1) \pi_i + 2 \sum_{i < j} \frac{Y_i Y_j}{\pi_i \pi_j} (\pi_{ij} - \pi_i \pi_j)$$

$$= n \sum_{i=1}^N \frac{Y_i^2}{\pi_i} - \sum_{i=1}^N \frac{Y_i^2}{\pi_i} \cdot \pi_i - n \sum_{i=1}^N \frac{Y_i^2}{\pi_i} + \sum_{i=1}^N \frac{Y_i^2}{\pi_i} \cdot \pi_i$$

$$+ 2 \sum_{i < j} \frac{Y_i Y_j}{\pi_i \pi_j} (\pi_{ij} - \pi_i \pi_j)$$

$$= \sum_{i=1}^N \frac{Y_i^2}{\pi_i} \pi_i (1 - \pi_i) + \sum_{i \neq j} \frac{Y_i Y_j}{\pi_i \pi_j} (\pi_{ij} - \pi_i \pi_j)$$

$$= V_p(\text{tht})$$