

二重標本抽出에서의 危險率에 의한 두가지 要因의 比較調査

金 益 贊*

Two-Factor Comparative Surveys with Risks in Double Sampling

Kim Ik-chan

요 약

二重 標本抽出에서의 最適 設計로서 2×2 table 로 주어지는 두가지 要因을 比較分析함에 있어서, 비용을 一般化한 危險率에 주어질때 이들 要因들의 同一한 精度를 最大化하기 爲한 標本の 크기 및 最小 分散을 算出하였다.

1. Double sampling and required precision

In many sample surveys the principle objective is to compare several sectors of a finite population. Specially, there may be several factors of interest and each of these factors may have been divided into several categories.

If the elements (N_{ij}) represented by the cells in a 2×2 table are not identifiable in advance, one cannot sample independently in each of them. However, one may select a large preliminary sample (n') and identify the subpopulation to which each sampled element (n'_{ij}) belongs. Then, for each sub-population, a sub-sample (n_{ij}) is selected for analytical surveys. Such a double sampling procedure is useful if the risk of identifying an element is small relative to the risk of securing the necessary information in the main

* 사범대학 수학교육과

survey. Now consider the optimal design for two-factor comparative surveys. The two factor α and τ are represented by a 2×2 table with (i, j) th cell denoting i th category of α and j th category of τ .

The two categories for each factor are compared by considering

$$\begin{aligned} D\alpha &= W_{.1}(\bar{Y}_{11} - \bar{Y}_{21}) + W_{.2}(\bar{Y}_{12} - \bar{Y}_{22}) \\ D\tau &= W_{1.}(\bar{Y}_{11} - \bar{Y}_{12}) + W_{2.}(\bar{Y}_{21} - \bar{Y}_{22}) \end{aligned} \tag{1}$$

where N_{ij} = total number of units in the (i, j) th cell, $W_{ij} = N_{ij}/N$, $W_{i.} = \sum_j W_{ij}$, $W_{.j} = \sum_i W_{ij}$, \bar{Y}_{ij} = true mean for (i, j) th cell and $N = \sum_i \sum_j N_{ij}$, the size of population.

Using the double sampling method, the unbiased estimators are given by

$$\begin{aligned} \widehat{D\alpha} &= \frac{n_{.1}'}{n'}(\bar{y}_{11} - \bar{y}_{21}) + \frac{n_{.2}'}{n'}(\bar{y}_{12} - \bar{y}_{22}) \\ \widehat{D\tau} &= \frac{n_{1.}'}{n'}(\bar{y}_{11} - \bar{y}_{12}) + \frac{n_{2.}'}{n'}(\bar{y}_{21} - \bar{y}_{22}) \end{aligned} \tag{2}$$

note that $n'_{i.} = \sum_j n'_{ij}$, $n'_{.j} = \sum_i n'_{ij}$ are obtained from the preliminary sample n' and sample mean \bar{y}_{ij} from n_{ij} .

Let $n_{ij} = n_{ij}' \nu_{ij}$, $0 < \nu_{ij} \leq 1$ and $w_{ij} = \frac{n_{ij}'}{n'}$ then

n_{ij} , w_{ij} , \bar{y}_{ij} are random variables and

Lemma 1. $E(n_{ij}) = E(n_{ij}' \nu_{ij}) = E(w_{ij} n' \nu_{ij}) = n' \nu_{ij} E(w_{ij}) = n' \nu_{ij} W_{ij}$

Lemma 2. $E\left(\frac{1}{n_{ij}}\right) \simeq \frac{1}{E(n_{ij})}$

Proof : See []

If equal precision is desired for $\widehat{D\alpha}$ and $\widehat{D\tau}$, we use the objective function ;

$$\begin{aligned} \bar{V} &= \frac{1}{2} \{ V(\widehat{D\alpha}) + V(\widehat{D\tau}) \} \\ &= \frac{1}{2} \left[E \left\{ \sum \sum \left(\frac{(n_{.j}')^2 + (n_{i.}')^2}{(n')^2} \right) \cdot \frac{S_{ij}^2}{n_{ij}} \right\} + V \left\{ \frac{n_{.1}'}{n'} (\bar{Y}_{11} - \bar{Y}_{21}) \right. \right. \\ &\quad \left. \left. + \frac{n_{.2}'}{n'} (\bar{Y}_{12} - \bar{Y}_{22}) \right\} + V \left\{ \frac{n_{1.}'}{n'} (\bar{Y}_{11} - \bar{Y}_{12}) + \frac{n_{2.}'}{n'} (\bar{Y}_{21} - \bar{Y}_{22}) \right\} \right] \end{aligned}$$

where S_{ij}^2 is the true variance in (i, j) th cell. Using lemma 1, 2 and the approximation

$$\begin{aligned} E \left\{ \frac{(n_{.j}')^2 + (n_{i.}')^2}{n'^2 n_{ij}} \right\} &\simeq \frac{\{E(n_{.j}')\}^2 + \{E(n_{i.}')\}^2}{n'^2 E(n_{ij})} = \frac{W_{.j}^2 + W_{i.}^2}{n' W_{ij} \nu_{ij}} \\ &= \frac{\tilde{g}_{ij}}{n' \nu_{ij}} \end{aligned}$$

then the objective function reduces to

$$\bar{V} = E \left(\sum \frac{g_{ij}^2}{n_{ij}} \right) \simeq \sum \frac{g_{ij}^2}{n' W_{ij} \nu_{ij}} \quad (3)$$

where $2g_{ij}^2 = \tilde{g}_{ij} W_{ij} S_{ij}^2$

2. Determine the risks in double sampling

Let Ω be a parameter space on random variable X and U be a numerical function defined on Ω whose value we wish to estimate on the basis of the outcome of an experiment $x \in X$.

Let A be the space of actions on real line R^1 and a non-randomized decision function δ^* defined on X be a numerical function specifying for each x the number $a \in A$ which will be chosen to estimate U when that x is observed. Then the loss function $L(U, \delta^*)$ defined on $\Omega \times A$ is the loss incurred when U is estimated by δ^* .

If we define the loss;

$$L(U, \delta^*) = |\delta^* - U| + \sum_i \sum_j C_{ij} n_{ij} + C' n'$$

and replace δ^* with strata mean $\bar{y}_{st} = \sum \sum w_{ij} \bar{y}_{ij}$, U with population mean \bar{Y} , then the risk function R is defined by

$$R(U, \delta^*) = E |\bar{y}_{st} - \bar{Y}| + \sum \sum C_{ij} \lambda_{ij} + C' n' \quad (4)$$

where C' is the cost of classification per unit and C_{ij} the cost of measuring a unit in (i, j) th cell.

Lemma 3. An estimator \bar{y}_{st} is unbiased estimator of $\bar{Y} = \sum \sum W_{ij} \bar{Y}_{ij}$

$$\begin{aligned} \text{Proof: } E(\bar{y}_{st}) &= E[E(\Sigma \Sigma w_{ij} \bar{y}_{ij} | w_{ij})] = E(\Sigma \Sigma w_{ij}) E_2(\bar{y}_{ij}) \\ &= E(\Sigma \Sigma w_{ij} \bar{Y}_{ij}) = \bar{Y} \end{aligned}$$

where the subscript 2 refer to an average over all random sample of n_{ij} units that can be drawn from a given n'_{ij} units.

Now the specified cost of taking the sample is generalized by the risk R and

Theorem 1.

$$R(U, \delta^*) \leq n' MD_p + \Sigma \Sigma W_{ij} MD_{ij} | n_{ij} - n'_{ij} | + \Sigma \Sigma C_{ij} n_{ij} \nu + C'n' \quad (5)$$

where MD_p is the true man deviation and MD_{ij} the (i, j) th cell mean deviation.

$$\begin{aligned} \text{Proof : } R(U, \delta^*) &= E | \bar{y}_{st} - \bar{Y} | + \Sigma \Sigma C_{ij} n_{ij} + C'n' \\ | \bar{y}_{st} - \bar{Y} | &= | \Sigma \Sigma w_{ij} \bar{y}_{ij} + \bar{Y} | \\ &= | \Sigma \Sigma w_{ij} \bar{y}_{ij}' + \Sigma \Sigma w_{ij} (\bar{y}_{ij} - \bar{y}_{ij}') - \bar{Y} | \\ &\leq | \Sigma \Sigma w_{ij} \bar{y}_{ij}' - \bar{Y} | + | \Sigma \Sigma w_{ij} (\bar{y}_{ij} - \bar{y}_{ij}') | \end{aligned}$$

And Since

$$\begin{aligned} MD_p &= \frac{1}{N} \Sigma_i \Sigma_j | y_{ij} - \bar{Y} | \quad \text{and} \quad E | y_{ij} - \bar{Y} | \leq \frac{n'}{N} \Sigma \Sigma | y_{ij} - \bar{Y} | \\ &= n' \cdot MD_p \end{aligned}$$

$$\begin{aligned} \text{So } E_2 | \Sigma \Sigma w_{ij} (\bar{y}_{ij} - \bar{y}_{ij}') | &= \Sigma \Sigma W_{ij} E_2 | (\bar{y}_{ij} - \bar{Y}_{ij}) - (\bar{y}_{ij}' - \bar{Y}_{ij}) | \\ &= \Sigma \Sigma W_{ij} | E_2 (\bar{y}_{ij} - \bar{Y}_{ij}) - E_2 (\bar{y}_{ij}' - \bar{Y}_{ij}) | \\ &\leq \Sigma \Sigma W_{ij} | n_{ij} MD_{ij} - n_{ij}' MD_{ij} | \end{aligned}$$

$$\text{Therefore } E(\bar{y}_{st} - \bar{Y}) \leq n' MD_p + \Sigma \Sigma W_{ij} MD_{ij} | n_{ij} - n_{ij}' |$$

This completes the proof.

Theorem 2.

Let the expected risk be $R^* = \{ R(U, \delta^*) \}$, then

$$R^* \leq n'B + n' \Sigma \Sigma v_{ij} W_{ij} D_{ij} \quad (6)$$

where $B = C' + MD_p + \sum \sum W_{ij}^2 MD_{ij}$ and $D_{ij} = C_{ij} - W_{ij} MD_{ij}$

$$\begin{aligned} \text{Proof : since } E(\sum \sum W_{ij} MD_{ij} | n_{ij} - n'_{ij} |) &= \sum \sum W_{ij} MD_{ij} E(n_{ij} | 1 - \frac{1}{\nu_{ij}} |) \\ &= \sum \sum W_{ij} MD_{ij} n' \nu_{ij} W_{ij} | \frac{1}{\nu_{ij}} - 1 | \\ &= \sum \sum n' MD_{ij} W_{ij}^2 (1 - \nu_{ij}) \end{aligned}$$

Hence

$$\begin{aligned} R^* = E(R) &\leq E(n' MD_p) + E(\sum \sum W_{ij} MD_{ij} | n_{ij} - n'_{ij} |) + n' \sum \sum C_{ij} \nu_{ij} W_{ij} + C' n' \\ &= n' MD_p + n' \sum \sum MD_{ij} W_{ij}^2 (1 - \nu_{ij}) + n' \sum \sum C_{ij} \nu_{ij} W_{ij} + C' n' \\ &= n'(C' + MD_p + \sum \sum W_{ij}^2 MD_{ij}) + n' \sum \sum \nu_{ij} W_{ij} (C_{ij} - W_{ij} MD_{ij}) \\ &= n'B + n' \sum \sum \nu_{ij} W_{ij} D_{ij} \end{aligned}$$

3. Optimum design for two factor comparative surveys with specified risk

We consider the optimum design to find those values of the preliminary sample size n' and main sample size n which maximize, for a given risk, the equal precision of comparisons of two-factor with categories.

Without loss the generality, we can assume that inequality in (6) change to equality, therefore

$$R^* = n'B + n' \sum \sum \nu_{ij} W_{ij} D_{ij} \tag{7}$$

Let find the values of n' and ν_{ij} which maximize (3)

$$\bar{V} = \sum \sum \frac{a_{ij}}{n' W_{ij} \nu_{ij}}$$

subject to (7) and $0 < \nu_{ij} \leq 1$, where $a_{ij} = 2g_{ij}^2$ are known constants. We determine first the optimal ν_{ij} for a given n' and then the optimal n' .

By Cauchy inequality ;

$$\sum \alpha_h^2 \sum \beta_h^2 - (\sum \alpha_h \beta_h)^2 = \sum_{ij>i} (\alpha_i \beta_j - \alpha_j \beta_i)^2$$

then $(\sum \alpha_h)^2 (\sum \beta_h)^2 \geq (\sum \alpha_h \beta_h)^2$

And if $\frac{\beta_1}{\alpha_1} = \frac{\beta_2}{\alpha_2} = \dots = \frac{\beta_2}{\alpha_2} = \text{constant}$ the

equality holds

Now let $R' = R^* - n'B = \sum \sum n' \nu_{ij} W_{ij} D_{ij}$, (8)

then the product $R' \bar{V} = \left(\sum \sum \frac{a_{ij}^2}{n' W_{ij} \nu_{ij}'} \right) \left(\sum \sum n' \nu_{ij} W_{ij} D_{ij} \right)$

Using the Cauchy inequalities,

Put $\alpha_h = \frac{a_{ij}}{\sqrt{n' W_{ij} \nu_{ij}}}$, $\beta_h = \sqrt{n' \nu_{ij} W_{ij} D_{ij}}$

then $\frac{\beta_h}{\alpha_h} = \frac{n' W_{ij} \nu_{ij} \sqrt{D_{ij}}}{a_{ij}}$ and

$$\frac{n' W_{11} \nu_{11} \sqrt{D_{11}}}{a_{11}} = \frac{n' W_{12} \nu_{12} \sqrt{D_{12}}}{a_{12}} = \dots = \frac{\sum \sum n' W_{ij} \nu_{ij} D_{ij}}{\sum \sum a_{ij} \sqrt{D_{ij}}}$$

$$= \frac{R^* - B}{\sum \sum a_{ij} \sqrt{D_{ij}}}$$

Hence $\frac{n' W_{ij} \nu_{ij} \sqrt{D_{ij}}}{a_{ij}} = \frac{R^* - B}{\sum \sum a_{ij} \sqrt{D_{ij}}}$

So the optimal ν_{ij} for fixed n' is given by

$$n' W_{ij} \nu_{ij} = \frac{a_{ij} (R^* - n'B)}{\sqrt{D_{ij}} \sum \sum a_{ij} \sqrt{D_{ij}}} \tag{9}$$

Provided $n' W_{ij} \nu_{ij} \leq n' W_{ij}$ for all i, j ; that is

$$\frac{a_{ij} (R^* - n'B)}{\sqrt{D_{ij}} \sum \sum a_{ij} \sqrt{D_{ij}}} \leq n' W_{ij}$$

Hence $n' \geq \frac{R^*}{B + W_{ij} \cdot \frac{\sqrt{D_{ij}}}{a_{ij}} \sum \sum a_{ij} \sqrt{D_{ij}}}$; $i, j = 1, 2$

$$= [B + W_{(1,1)} \cdot \frac{\sqrt{D_{(1,1)}}}{a_{(1,1)}} \cdot \sum \sum a_{ij} \sqrt{D_{ij}}]^{-1} \cdot R^* \equiv m_{11}' \tag{10}$$

where (1, 1) denotes the group with the smallest value of $W_{ij} \sqrt{D_{ij}} / a_{ij}$

The minimum value of \bar{V} for $n' \geq m'_{11}$ after substituting the optimal ν_{ij} in (3), is given by

$$\bar{V}_{11}(n') = \frac{(\sum\sum a_{ij} \sqrt{D_{ij}})^2}{R^* - n'B} \quad (11)$$

so that the minimum occurs at the value $m_{11} = m'_{11}$.

Note that $\nu_{(1,1)} = 1$ when $n' = m_{11}$

To examine values of n' smaller than m_{11} , set $\nu_{11} = 1$ and use the Cauchy inequality to obtain the remaining ν_{ij} .

This gives

$$n' W_{ij} \nu_{ij} = \frac{a_{ij}}{\sqrt{D_{ij}}} \left\{ \left(\frac{R^* - n'B - n' W_{(1,1)} \cdot D_{(1,1)}}{\sum\sum a_{ij} \sqrt{D_{ij}}} \right) / \sum\sum_{(i,j) \neq (1,1)} a_{ij} \sqrt{D_{ij}} \right\} \quad (i, j) \neq (1, 1) \quad (12)$$

Provided

$$n' \geq \left\{ B + D_{(1,1)} W_{(1,1)} + \left(W_{(1,2)} \cdot \frac{\sqrt{D_{(1,2)}}}{a_{(1,2)}} \right) \sum\sum_{(i,j) \neq (1,1)} a_{ij} \sqrt{D_{ij}} \right\}^{-1} \cdot R^* \equiv m'_{12}$$

where Σ denotes the summation over $i, j \neq (1)$, and (1, 2) denotes the group with the second smallest values of $W_{ij} \sqrt{D_{ij}} / a_{ij}$

Therefore, the minimum value of \bar{V} , for n' in the range $m'_{12} \leq n' \leq m'_{11}$ is given by

$$\bar{V}_{12}(n') = \frac{a_{(1,1)}^2}{n' W_{(1,1)}} + \frac{(\sum\sum_{(i,j) \neq (1,1)} a_{ij} \sqrt{D_{ij}})^2}{R^* - n'(B + D_{(1,1)} \cdot W_{(1,1)})} \quad (13)$$

From this \bar{V}_{12} , to find the optimal n' over the range

$m'_{12} \leq n' \leq m'_{11}$, we put $\frac{d\bar{V}_{12}(n')}{dn'} = 0$, then

$$n' = [B + D_{(1,1)} \cdot W_{(1,1)} + \frac{\{W_{(1,2)} (B + D_{(1,1)} \cdot W_{(1,1)})\}^{\frac{1}{2}}}{a_{(1,1)}} \cdot \sum\sum_{(i,j) \neq (1,1)} a_{ij} \cdot \sqrt{D_{ij}}]^{-1} \cdot R^*$$

(14)

If $d\bar{V}_{12}(n')/dn'$ does not vanish for $n' \geq m_{12}'$, we need to see ν_{11} , ν_{12} , and so forth, = 1 in turn until the turning point of \bar{V} is found, and note that $\bar{V}(n')$ has a unique minimum.

Literature Cited

- (1) Aggarwal, O. M. 1960. Bayes and minimax procedures in sampling from finite and infinite population- I, American Statistical Association Journal 206~218.
- (2) Booth, G and Sedransk, J. 1969. Planing some two-factor comparative surveys, American Statistical Association Journal 64, 560~573.
- (3) Cochran, W. G. 1977. Sampling Techniques 3rd ed, John Wiley and Sons, New-York.
- (4) Rao, S. N. K. 1973. On double sampling for stratification and analytical surveys, Biometrika 60, 125~133.
- (5) Sedransk, J. 1965. A double sampling scheme for analytical surveys, American Statistical Association 60, 985~1004.