# TWO DIMENSIONAL SYSTEMATIC SAMPLING AND THE ASSOCIATED STRATIFIED AND RANDOM SAMPLING

*By* A. C. DAS

*Statistical Laboratory, Calcutta*

## INTRODUCTION

Systematic sampling has now come to be of much interest both in theoretical and applied statistics. This type of sampling with its various modifications is well known. Madow (1944), Cochran (1946) and Yates (1948) have discussed some particular types in one dimension. This paper deals with the two dimensional generalisation of Cochran's results for the uni-dimensional case.

Let the universe to be sampled consist of $ml$ rows, each of $nk$ cells in $mn$ strata, each stratum containing $lk$ cells. The cells are square or rectangular as the case may be. Let the co-ordinates $(x,y)$ be taken in such a manner that, $y = $ constant, gives a row and $x = $ constant, gives a column. Thus in each row $n$ sets of $k$ cells are along $x$ direction and in each column $m$ sets of $l$ cells are along $y$ direction. In a particular stratum, $v$ cells out of $lk$ cells are chosen at random. Then every $l$th row of cells (along $y$ direction), beginning from the one in which each of these chosen cells occurs, is selected. On each of these selected rows, cells in same column with the randomly chosen ones and every $k$th cell (along $x$ direction) are taken, thus comprising a sample of size $mnv$, which is a kind of systematic sampling in two dimensions. This type of systematic sampling may also be defined as: $v$ cells out of $lk$ cells are taken at random from a particular stratum; then in every other stratum a similar set of $v$ cells are chosen, having the same configuration within the set and also with respect to the stratum as in the first stratum. We may also take a random sample of size $v$ from each of the strata forming a total stratified sample of size $mnv$. Alternatively, $mnv$ cells may be taken at random out of all $mnlk$ cells without stratification, which will constitute a perfectly random sample.

We assume this universe to be homogeneous in the sense that $E(z_i) = \mu$, where $z$'s are the observables of the cells, and that observables have got common variance and different space-correlations. Let $\rho(u,v)$ be the space-correlation between points separated by distances $u$ and $v$ along $x$ and $y$ directions respectively, *i.e.*, if $P$ be the cell $(x,y)$, we shall pair its yield against the yield of the cell $(x+u, y+v)$ and from the set of all such possible pairs shall calculate $\rho(u, v)$, which may have certain fluctuations due to the finiteness of the universe. If this universe be infinite, such an effect will be absent. If, however, the universe be finite it may be considered to be a sample from an infinite population, as has been done by Cochran (1946), where $E(z) = \mu$, $V(z) = \sigma^2$, and $E[(z_{i+u)j+v} - \mu)(z_{ij} - \mu)] = \rho(u,v)$ for the infinite population. This $\rho(u,v)$ for different values of $u$ and $v$ will form a solid figure which may be called a correlopiped. This figure will be a continuous one for infinitesimal sizes of the cells.

In this paper the expected values of the variances of the grand sample mean, as averaged over all finite populations from the infinite population, under random, stratified and systematic sampling will be denoted by $\sigma_r{}^2$, $\sigma_{st}{}^2$ and $\sigma_{sy}{}^2$ respectively following the symbols used by Cochran (1946). $\Delta_1, \Delta_2, \delta_1, \delta_2$ are the operators, such that

$$\Delta_1\rho(u,v) = \rho(u+1,v) - \rho(u,v),$$
$$\Delta_2\rho(u,v) = \rho(u,v+1) - \rho(u,v), \quad \delta_1\rho(u,v) = \rho(u+\tfrac{1}{2},v) - \rho(u-\tfrac{1}{2},v),$$
$$\delta_1{}^2\rho(u,v) = \rho(u+1,v) + \rho(u-1,v) - 2\rho(u,v), \quad \Delta_1\Delta_2\rho(u,v) = \Delta_1\rho(u,v+1) - \Delta_1\rho(u,v)$$

and similar definitions for $\delta_2$, $\delta_2{}^2$ and $\delta_1{}^2\delta_2{}^2$. We shall also use $E_1$, $E_2$ such that $E_1\rho(u,v) = \rho(u+1,v)$ and $E_2\rho(u,v) = \rho(u,v+1)$. These $E_1$'s are different from $E$ (with no subscript) which denotes "expectation". It is also to be noted that $\rho(u,v) = \rho(-u,-v)$ and $\rho(u,-v) = \rho(-u,v)$ are different and we shall denote $\rho(u,v) + \rho(-u,+v)$ by $\psi(u,v)$ which is symmetrical with respect to $u$ and $v$. $\Delta_1\psi$ and $\Delta_2\psi$ will always denote $\Delta_1\psi(u,v)$ and $\Delta_2\psi(u,v)$ respectively.

In section 1, $\sigma_r{}^2$, $\sigma_{st}{}^2$ and $\sigma_{sy}{}^2$ have been calculated and it has been shown that their relative proportions are independent of $v$, and hence the relative efficiencies of the three types of sampling. In section 2 a set of sufficient conditions under which stratified sampling is more efficient than random sampling has been stated in a number of theorems, and in section 3, a sufficient condition under which systematic sampling is more efficient than stratified sampling has been stated in Theorem 4. From the results of these two sections, we can find out the set of sufficient conditions under which systematic sampling is more efficient than random sampling; though the two have not been directly compared. Cochran's results come out as special cases. But there are some improvements over Cochran (1946) that the results are true whatever may be $\rho$, positive or negative, and that $v$ may be greater than unity.

### Section 1

It can be proved that $\sum\limits_{i=1}^{n}(z_i - \bar{z})^2 = \frac{1}{n}\sum\limits_{i=1}^{n}\sum\limits_{i<j}^{n}(z_i - z_j)^2$     ... (1.1)

*Lemma* 1. In a plot containing $M$ rows of $N$ cells each, with $z_{ij}$ as the stochastic variate corresponding to the $(i,j)$th cell, $i = 1, 2, \ldots, N$; $j = 1, 2, \ldots, M$.

$$E\left\{\sum_{i=1}^{N}\sum_{j=1}^{M}(z_{ij} - \bar{z})^2\right\} = \frac{2\sigma^2}{MN}\left[\frac{MN(MN-1)}{2} - N\sum_{v=1}^{M-1}(M-v)\rho(0,v) - M\sum_{u=1}^{N-1}(N-u)\rho(u,0) - \sum_{u=1}^{N-1}\sum_{v=1}^{M-1}(N-u)(M-v)\psi(u,v)\right]$$

*Proof.* By (1.1), $\sum\limits_{i=1}^{N}\sum\limits_{j=1}^{M}(z_{ij} - \bar{z})^2 = \frac{1}{MN}\Sigma(z_i - z_{i'})^2$ where the summation is to be extended over all combinations of $z_i$ and $z_{i'}$ in the finite plot, which is a sample from an infinite plot of the type mentioned in the introduction. There are ${}^{MN}c_2$ combinations in all splitted as: (i) within each row ${}^{N}c_2$ combinations, being $M$. ${}^{N}c_2$ in all, (ii) within each

column $^M c_2$, in all $N \, ^M c_2$ and (iii) between $(x, y)$ and $(x+u, y+v)$, $(u \neq 0, v \neq 0)$, giving $2 \, ^M c_2 \, ^N c_2$.

Now,
$$E\left\{ \sum_{i,i'} (z_i - \bar{z})^2 \right\} = \frac{1}{MN} E\{ \Sigma(z_i - \mu + \mu - z_{i'})^2 \} = \frac{2\sigma^2}{MN} \Sigma(1 - \rho(z_i, z_{i'})) \qquad \dots (1.2)$$

where $\rho(z_i, z_{i'})$ means the correlation between $z_i$ and $z_{i'}$ and the summation is to be carried over all combinations of $i$ and $i'$ (i.e., $^{MN}c_2$ in all). Let us first find $\Sigma\rho(z_i, z_{i'})$. It will be the sum of three expressions corresponding to the three cases (i), (ii) and (iii) mentioned above. In case (i) there are $(N-u)$ cells $(u,0)$–distance apart in each row; therefore, $\Sigma\rho(z_i, z_{i'})$ of this type over all rows is $N \sum\limits_{u=1}^{N-1} (N-u)\rho(u,0)$. Similarly for case

(ii), $\Sigma\rho(z_i, z_{i'})$ is $N \sum\limits_{v=1}^{M-1} (M-v)\rho(0, v)$. In case (iii) considering the four points $(x, y), (x+u, y)$, $(x, y+v)$ and $(x+u, y+v)$ with $u \neq 0$ and $v \neq 0$, it is obvious that there will be only 2 effective types of distances of case (iii) among these four points, namely, $(u,v)$ and $(-u, v)$. The distance $(-u, -v)$ is equivalent to $(u, v)$, since the space correlation will be the same for both the distances; similarly, the distance $(u, -v)$ is equivalent to $(-u, v)$. Therefore, summing $\Sigma\rho(z_i, z_{i'})$ over this type we get $\sum\limits_{u=1}^{N-1} \sum\limits_{v=1}^{M-1} (N-u)(M-v)\psi(u,v)$;

beacause there are $(N-u)(M-v)$ sets of such quadruplets. It is a point to note that $\rho(u, 0) = \rho(-u, 0)$ and $\rho(0, v) = \rho(0, -v)$.

Again $\Sigma 1$ over all combinations is $^{MN}c_2$.

Then the lemma follows from (1.2), which can be written as:

$$E\left\{ \sum_{i=1}^{N} \sum_{j=1}^{M} (z_{ij} - \bar{z})^2 \right\} = (MN-1)\sigma^2\{1 - \phi(M, N)\} = (MN-1)\sigma^2 L(M, N) \qquad \dots (1.3)$$

where $\phi(M, N) = \dfrac{2}{M(MN-1)} \sum\limits_{v=1}^{M-1} (M-v)\rho(0, v) + \dfrac{2}{N(MN-1)} \sum\limits_{u=1}^{N-1} \rho(u, 0)$

$$+ \frac{2}{MN(MN-1)} \sum_{u=1}^{N-1} \sum_{v=1}^{M-1} (N-u)(M-v)\psi(u, v) \qquad \dots (1.4)$$

and
$$L(M, N) = 1 - \phi(M, N) \qquad \dots (1.5)$$

∴ putting $M = ml$ and $N = nk$ in (1.3), we have for the whole finite population:

$$E\{\Sigma(z - \bar{z})^2\} = (mnlk - 1)\sigma^2\{1 - \phi(ml, nk)\} \qquad \dots (1.6)$$

Similarly, for each stratum, putting $M = l$ and $N = k$ in (1.3)

$$E\{\Sigma(z - \bar{z})^2\} = (lk - 1)\sigma^2\{1 - \phi(l, k)\} \qquad \dots (1.7)$$

Therefore, for $mn$ strata, $E$ (s.s. within strata) $= mn(lk - 1)\sigma^2\{1 - \phi(l, k)\} \qquad \dots (1.8)$

13

For each of the systematic samples for $\nu = 1$,

$E\{\Sigma(z-\bar{z})^2\} = (mn-1)\sigma^2\{1-\phi_{\nu}(m,n)\}$ putting $M=m$ and $N=n$, replacing $\rho(u,v)$ by $\rho(uk, vl)$ in (1.3) and (1.4) and replacing $\phi$ by $\phi_{\nu}$.

$\therefore$ for $\nu=1$, $E$ (s.s. within systematic sample) $= lk(mn-1)\sigma^2\{1-\phi_{\nu}(m,n)\}$ ... (1.9)

Now, for a random sample size $n'$ from a finite population of size $N$, the variance of sample mean $= \dfrac{1}{n'} \dfrac{N-n'}{N-1} \left\{ \dfrac{1}{N} \Sigma(z-\bar{z})^2 \right\}$, summation being taken over whole finite population.

Hence putting $N = mnlk$ in our case, $n' = mn\nu$, we have, using (1.6) the expected variance of sample mean as:

$$\sigma_e^2 = \frac{1}{mn\nu} \frac{mn(lk-\nu)}{(mnlk-1)} \frac{1}{mnlk} (mnlk-1)\sigma^2\{1-\phi(ml, nk)\} = \frac{\sigma^2}{mn\nu} \left( 1 - \frac{\nu}{lk} \right) \left\{ 1 - \phi(ml, nk) \right\}$$

... (1.10)

Again, there are $mn$ strata with $lk$ cells in each and let $x_{11}, x_{12}, ..., x_{1\nu}$ be the sample of size $\nu$ from the ith stratum. Here $x_{1j}$ is a stochastic variate and is not to be confused with $x$, having no subscripts, that denotes a co-ordinate in the introduction. Then

$$\bar{x} = \frac{1}{mn\nu} \sum_{i=1}^{mn} \sum_{j=1}^{\nu} x_{1j} = \frac{1}{mn} \sum_{i=1}^{mn} \bar{x}_1, \quad \text{where } \bar{x}_1 = \frac{1}{\nu} \sum_{j=1}^{\nu} x_{1j} ; \quad V(\bar{x}) = \frac{1}{m^2n^2} \sum_{i=1}^{mn} V(\bar{x}_1).$$

But $E\{V(\bar{x}_1)\} = \dfrac{1}{\nu} \dfrac{(lk-\nu)}{(lk-1)} \dfrac{(lk-1)}{lk} \sigma^2\{1-\phi(l, k)\}$, from (1.7).

Hence
$$\sigma_{s1}^2 = E\{V(\bar{x})\} = \frac{\sigma^2}{mn\nu} \left( 1 - \frac{\nu}{lk} \right) \left\{ 1 - \phi(l, k) \right\}$$
... (1.11)

Now, to find the expected variance $\sigma_{sy}^2$ of systematic sample, let $x_{11}, x_{12}, ..., x_{1\nu}$ be the sample from the first stratum. Then the samples from other strata are automatically fixed up. The grand sample mean $\bar{x}_{sy} = \dfrac{1}{mn\nu} \sum_{i=1}^{mn} \sum_{j=1}^{\nu} x_{1j}$ and $\sigma_{sy}^2 = E(\bar{x}_{sy}-\bar{x})^2$, where $\bar{x}$ is the grand mean of the finite population.

Here let $\sum_{i=1}^{mn} x_{1j}/mn$ be denoted by $_j\bar{x}_{sy}$, so that $(\bar{x}_{sy}-\bar{x}) = \dfrac{1}{\nu} \sum_{j=1}^{\nu} (_j\bar{x}_{sy}-\bar{x})$

Hence
$$(\bar{x}_{sy}-\bar{x})^2 = \frac{1}{\nu^2} \left\{ \sum_{j=1}^{\nu} (_j\bar{x}_{sy}-\bar{x})^2 + \sum_{\substack{j,j'=1 \\ (j \neq j')}}^{\nu} (_j\bar{x}_{sy}-\bar{x}) (_{j'}\bar{x}_{sy}-\bar{x}) \right\}$$
... (1.12)

We shall take the average of (1.12) over $^{lk}c_r$ systematic samples.

If $x_{ij}$ be a random point (observable) from the first stratum, $_jx_{sy}$ is the mean of the corresponding systematic sample. Then for summation of (1.12) over $^{lk}c_s$ systematic samples, any $_jx_{sy}$ $(j = 1, 2, ..., mn)$, will occur $^{lk-1}c_{r-1}$ times, because only $(r-1)$ will be required from the remaining $(lk-1)$ to complete a sample of size $r$; and $_jx_{sy}$ and $_{j'}x_{sy}$, $(j \neq j')$, will occur together $^{lk-2}c_{r-2}$ times; therefore in the expression (1.12), when summed over $^{lk}c_r$ samples, $(_jx_{sy}-\bar{x})$, $j = 1, 2, ..., lk$, will occur $^{lk-1}c_{r-1}$ times and

$$(_jx_{sy}-\bar{x})(_{j'}x_{sy}-\bar{x})(j \neq j'), \ ^{lk-2}c_{r-2} \text{ times. But } \sum_{i=1}^{lk} (_ix_{sy}-\bar{x}) = 0,$$

whence $\sum_{i=1}^{lk} (_jx_{sy}-\bar{x})(_ix_{sy}-\bar{x}) = 0$; and $\sum_{j=1}^{lk} \sum_{\substack{i=1 \\ i \neq j}}^{lk} (_jx_{sy}-\bar{x})(_ix_{sy}-\bar{x}) = - \sum_{j=1}^{lk} (_jx_{sy}-\bar{x})^2$

Hence in the summation, $(_jx_{sy}-\bar{x})^2$ occurs $(^{lk-1}c_{r-1} - ^{lk-2}c_{r-2}) = ^{lk-2}c_{r-1}$ times; and the required average of (1.12) is

$$\frac{^{lk-2}c_{r-1}}{^{lk}c_r} \frac{1}{r^2} \sum_{i=1}^{lk} (_jx_{sy}-\bar{x})^2 = \frac{1}{mnr} \frac{(lk-r)}{lk(lk-1)} \sum_{i=1}^{lk} mn(_jx_{sy}-\bar{x})^2$$

$$= \frac{1}{mnr} \frac{(lk-r)}{lk(lk-1)} \left\{ \sum_{i=1}^{mn} \sum_{j=1}^{lk} (x_{ij}-\bar{x})^2 - (\text{a.s. within systematic sample with } r = 1) \right\} \quad ...(1.13)$$

Hence taking expectation

$$\sigma_{sy}^2 = \frac{\sigma^2}{mnr} \frac{(lk-r)}{lk(lk-1)} \left\{ (mnlk-1)-(mnlk-1)\phi(ml, nk)-lk(mn-1)+lk(mn-1)\phi_{sy}(m,n) \right\}$$
$$\text{from (1.8) and (1.9)}$$

$$= \frac{\sigma^2}{mnr} \left( 1-\frac{r}{lk} \right) \left\{ 1-\frac{(mnlk-1)}{(lk-1)} \phi(ml,nk) + \frac{lk(mn-1)}{(lk-1)} \phi_{sy}(m,n) \right\} \quad ... (1.14)$$

and from (1.10) and (1.11):

$$\sigma_r^2 = \frac{\sigma^2}{mnr} \left( 1-\frac{r}{lk} \right) \left\{ 1-\phi(ml,nk) \right\} \quad ... (1.15)$$

$$\sigma_{st}^2 = \frac{\sigma^2}{mnr} \left( 1-\frac{r}{lk} \right) \left\{ 1-\phi(l,k) \right\} \quad ... (1.16)$$

where $\phi(M, N)$ is given in (1.4), and

$$\phi_{sy}(m, n) = \frac{2}{m(mn-1)} \sum_{v=1}^{m-1} (m-v)\rho(0, vl) + \frac{2}{n(mn-1)} \sum_{u=1}^{n-1} \rho(uk, 0)$$
$$+ \frac{2}{mn(mn-1)} \sum_{u=1}^{n-1} \sum_{v=1}^{m-1} (n-u)(m-r)\phi(uk, vl)$$

It is thus seen that the relative efficiencies $\frac{1}{\sigma_r^2} : \frac{1}{\sigma_{s1}^2} : \frac{1}{\sigma_{r7}^2}$ are independent of $\nu$ and are dependent on the nature of stratification and the correlopiped, of which, if the latter be known, the former can be adjusted to suit the purpose.

<div style="text-align:center">SECTION 2</div>

Here we shall find out a set of sufficient conditions under which stratified sampling is more efficient than random sampling, $i.e.$, under which $\sigma_r^2 \geqslant \sigma_{s1}^2$. We shall at first prove a few important lemmas.

*Lemma II.* For all values of $l$ and $k$, $L(l+1, k) \geqslant L(l, k)$, if the following conditions be satisfied ; (i) $\Delta_1\psi(u, v) \leqslant 0$, (ii) $\Delta_2\psi(u, v) \leqslant 0$, and (iii) $2\psi(u+1, v) \geqslant \psi(u, v+1) + \psi(u+1, v+1)$.

*Proof.* Replacing $M$ by $l$ and $N$ by $k$ in (1.4) and (1.5), we get $L(l,k)$. Then, putting $l' = l+1, A = lk-1, A' = l'k-1,$ and

$$\beta_{0v} = 2\left(\frac{l-v}{lA} - \frac{l'-v}{l'A'}\right),\ \alpha_{u0} = \frac{k-u}{AA'},\ \alpha_{uv} = \frac{(k-v)}{k}2\left(\frac{l-v}{lA} - \frac{l'-v}{l'A'}\right) \quad \ldots\ (2.1)$$

where $u \neq 0$ and $v \neq 0$   we have

$$L(l+1; k)-L(l,k) = \sum_{v=1}^{l}\beta_{0v}\rho(0,v)+\sum_{u=1}^{l}\sum_{v=0}^{l}\alpha_{uv}\psi(u,v) \quad \ldots\ (2.2)$$

Now because $\psi(u,v) = -\Delta_2\psi(u,v) + \psi(u,v+1),$   therefore

$$\sum_{v=0}^{l}\alpha_{uv}\psi(u,v) = -\sum_{v=0}^{l-1}S_{uv}\Delta_2\psi(u,v)+S_{ul}\psi(u,l) \quad \ldots\ (2.3)$$

where

$$S_{uv} = \sum_{j=0}^{v}\alpha_{uj}$$

And

$$\sum_{v=1}^{l}\beta_{0v}\rho(0,v) = -\sum_{v=1}^{l-1}S'_{0v}\Delta_2\rho(0,v)+S'_{0l}\rho(0,l) \quad \ldots\ (2.4)$$

where

$$S'_{0v} = \sum_{j=1}^{v}\beta_{0j}$$

Then $S'_{0v} = v\left\{y\left(\frac{1}{lA} - \frac{1}{l'A'}\right) + \frac{(k+1)}{l'AA'}\right\}$, where $y = l-1-v.$

Hence for $y \geqslant 0,$ i.e., $v \leqslant (l-1), S'_{0v} \geqslant 0$ $\quad \ldots\ (2.5)$

and $\qquad\qquad\qquad S'_{0l} = -\frac{(k-1)}{AA'}$ $\quad \ldots\ (2.6)$

Now $S_{uv} = \alpha_{u0}+\sum_{j=1}^{v}\alpha_{uj} = \frac{k-u}{AA'} + \frac{k-u}{k}S'_{0v}$, which from (2.5), is positive for $v \leqslant (l-1),$

and using (2.6), $\qquad\qquad\qquad S_{ul} = \frac{(k-u)}{kAA'}$ $\quad \ldots\ (2.7)$

Then from (2.3), (2.4), (2.5), (2.6) and (2.7), the equation (2.2) leads to

$$L(l+1, k) - L(l, k) = -\sum_{u=1}^{k} \sum_{v=0}^{l-1} \frac{(k-u)}{AA'} \Delta_2 \psi(u, v) - \sum_{u=1}^{k} \sum_{v=0}^{l-1} \frac{(k-u)}{k} S'_{0v} \Delta_1 \psi(u, v)$$

$$- \sum_{v=1}^{l-1} S'_{0v} \Delta_2 \psi(0, v) + \sum_{u=1}^{k} \frac{(k-u)}{kAA'} \psi(u, l) - \frac{(k-1)}{2AA'} \psi(0, l) \qquad \dots \text{ (2.8)}$$

Now $\sum_{u=1}^{k} \frac{(k-u)}{kAA'} \psi(u, l) - \frac{(k-1)}{2AA'} \psi(0, l) = \sum_{u=1}^{k} \frac{(k-u)}{kAA'} \left\{ \psi(u, l) - \psi(0, l) \right\} = \sum_{u=1}^{k} \frac{(k-u)}{kAA'} \sum_{i=0}^{u-1} \Delta_1 \psi(i, l)$

$$\dots \text{ (2.9)}$$

$$= \sum_{u=0}^{k} \frac{(1-u)(k-u)}{2kAA'} \Delta_1 \psi(u-1, l) + \sum_{u=1}^{k} \frac{(k-u)}{2AA'} \Delta_1 \psi(u-1, l) \qquad \dots \text{ (2.10)}$$

For $\Delta_1 \psi \leqslant 0$, the first part of (2.10) is greater than or equal to zero, but the other part is not positive ; therefore for $\Delta_2 \psi \leqslant 0$, and $\Delta_1 \psi \leqslant 0$, we have from (2.8):

$$L(l+1, k) - L(l, k) = Q + \sum_{u=1}^{k} \frac{(k-u)}{2AA'} \Delta_1 \psi(u-1, l) - \sum_{u=1}^{k} \sum_{v=0}^{l-1} \frac{(k-u)}{AA'} \Delta_2 \psi(u, v) \qquad \dots \text{ (2.11)}$$

where $Q \geqslant 0$.

Again $\sum_{v=0}^{l-1} \Delta_2 \psi(u, v) = \psi(u, l) - \psi(u, 0)$.

Hence, from (2.11), $L(l+1, k) - L(l, k) = Q + \sum_{u=1}^{k} \frac{(k-u)}{2AA'} \{2\psi(u, 0) - \psi(u, l) - \psi(u-1, l)\}$

$$\dots \text{ (2.12)}$$

Now, $2\psi(u, 0) - \psi(u, l) - \psi(u-1, l) \geqslant 2\psi(u, l-1) - \psi(u, l) - \psi(u-1, l)$, for $\Delta_2 \psi \leqslant 0$.

$\therefore$ if $2\psi(u+1, v) \geqslant \psi(u, v+1) + \psi(u+1, v+1)$, for all values of $v$,

$2\psi(u, l-1) - \psi(u-1, l) - \psi(u, l) \geqslant 0$, i.e., expression (2.12) and hence (2.11) is greater than or equal to zero.

Hence, the lemma is proved.

In a similar way, the following lemma can be proved.

*Lemma III.* For all values of $l$ and $k$, $L(l, k+1) \geqslant L(l, k)$, if the following conditions be satisfied: (i) $\Delta_1 \psi(u, v) \leqslant 0$, (ii) $\Delta_2 \psi(u, v) \leqslant 0$, and (iii) $2\psi(u+1, v) \geqslant \psi(u+1, v) + \psi(u+1, v+1)$.

*Theorem 1.* For all infinite populations in which (i) $\Delta_1 \psi \leqslant 0$, (ii) $\Delta_2 \psi \leqslant 0$, (iii) $2\psi(u, v+1) \geqslant \psi(u+1, v) + \psi(u+1, v+1)$ and (iv) $2\psi(u+1, v) \geqslant \psi(u, v+1) + \psi(u+1, v+1)$, $\sigma_{sn}^2 \leqslant \sigma_r^2$ for any size of the sample ; and $\sigma_{sl}^2 < \sigma_r^2$ unless equality holds in each of the above four cases.

*Proof.* These four conditions satisfy all the conditions of Lemma II and Lemma III;

$$\therefore \qquad L(l, k) \leqslant L(l+1, k) \leqslant L(l+1, k+1) \leqslant \dots \leqslant L(ml, nk)$$

*i.e.,* $\qquad 1 - \phi(l, k) \leqslant 1 - \phi(ml, nk)$, from (1.5).

Hence, $\quad \sigma_{s_l}{}^2 \leqslant \sigma_r{}^2$, from (1.10) and (1.11).

Thus the theorem is proved.

Again, if $\Delta_1 \psi = \Delta_2 \psi$, the conditions (iii) and (iv) of the above theorem are satisfied. Hence, a set of more stringent conditions can be stated as:

*Corollary.* For all infinite populations in which $\Delta_1 \psi = \Delta_2 \psi \leqslant 0$, $\sigma_{s_l}{}^2 \leqslant \sigma_r{}^2$ for any size of the sample, and $\sigma_{s_l}{}^2 < \sigma_r{}^2$ unless equality holds in each of the above cases.

*Theorem 2.* For all infinite populations in which stratification is made by parallel strips along u-direction (say), $\sigma_r{}^2 \geqslant \sigma_{s_l}{}^2$ for any size of the sample, if the following conditions be satisfied : (i) $\Delta_1 \psi \leqslant 0$, (ii) $\Delta_2 \psi \leqslant 0$, and (iii) $2\psi(u+1, v) \geqslant \psi(u, v+1) + \psi(u+1, v+1)$; and $\sigma_r{}^2 > \sigma_{s_l}{}^2$ unless equality holds in each of the above three cases.

*Proof.* These three conditions satisfy all the conditions of Lemma II.

Then $\qquad L(l, k) \leqslant L(l+1, k) \leqslant \dots \leqslant L(ml, k)$, *i.e.,* $1 - \phi(l, k) \leqslant 1 - \phi(ml, k)$

Here, of course, $k$ represents the total number of cells in a stratum; therefore

$\sigma_{s_l}{}^2 \leqslant \sigma_r{}^2$ and the theorem follows.

Now, condition (iii) is satisfied, if $\Delta_1 \psi \geqslant \Delta_2 \psi$. Hence a set of more stringent conditions can be stated as:

*Corrolary 1.* For all infinite populations in which stratification is made by parallel strips along u-direction (say), $\sigma_r{}^2 \leqslant \sigma_{s_l}{}^2$ if $\Delta_2 \psi \leqslant \Delta_1 \psi \leqslant 0$; and $\sigma_r{}^2 > \sigma_{s_l}{}^2$ unless equality holds in each of these cases.

*Corrolary 2.* For $k = 1$ in the above case, *i.e.,* for the one-dimensional field along v-direction, $\sigma_r{}^2 \geqslant \sigma_{s_l}{}^2$, if $\Delta_2 \psi(u.v) \leqslant 0$, or, $\Delta_2 \rho(0,v) \leqslant 0$; because here $u = 0$. This case has been discussed by Cochran (1946) for $v = 1$. All these theorems give the conditions irrespective of the values of $l,m,n,k$ and hence any pattern of stratification is more efficient under them. But the author (1949) has given another set of conditions that depends on some relations of the parameters. We shall now discuss it. Here also, we have got some degrees of freedom in stratification. It can be stated in the following theorem.

*Theorem 3.* For all patterns of stratification in which $\sqrt{l} \leqslant k \leqslant l^2$ and $(n-1)k = (m-1)l$, $\sigma_{s_l}{}^2 \leqslant \sigma_r{}^2$, if the following conditions be satisfied: (i) $\Delta_1 \psi \leqslant 0$, (ii) $\Delta_2 \psi \leqslant 0$ and (iii) $\Delta_1 \Delta_2 \psi \geqslant 0$ and $\sigma_{s_l}{}^2 < \sigma_r{}^2$ unless equality holds in each of these three cases.

The parametric relation is, of course, satisfied, if $k = l$ and $m = n$. This theorem will be proved with the help of lemma IV.

*Lemma IV.* The expression $\sum_{i=0}^{k}\sum_{j=0}^{l}\alpha_{ij}\psi(i,j)\geqslant 0$, if the following conditions be satisfied:

(i) $\sum_{i=0}^{u}\sum_{j=0}^{v}\alpha_{ij}\geqslant 0$ for $u \leqslant k$, $v \leqslant l$; but $\sum_{i=0}^{k}\sum_{j=0}^{l}\alpha_{ij}=0$.

(ii) $\Delta_2\psi(k,r)\leqslant 0$, (iii) $\Delta_1\psi(u,l)\leqslant 0$, and (iv) $\Delta_1\Delta_2\psi(u,r)\geqslant 0$.

*Proof.* $\sum_{i=0}^{k}\sum_{j=0}^{l}\alpha_{ij}\psi(i,j) = -\sum_{j=0}^{l}\sum_{i=0}^{k-1}S_{ij}\Delta_1\psi(u,j) + \sum_{j=0}^{l}S_{kj}\psi(k,j)$, where $S_{uj} = \sum_{i=0}^{u}\alpha_{ij}$

$$= \sum_{u=0}^{k-1}\sum_{r=0}^{l-1}T_{ur}\Delta_1\Delta_2\psi(u,r) - \sum_{u=0}^{k-1}T_{ul}\Delta_1\psi(u,l) - \sum_{r=0}^{l-1}T_{kr}\Delta_2\psi(k,r)+T_{kl}\psi(k,l)$$

where $T_{ur} = \sum_{j=0}^{r}S_{uj} = \sum_{i=0}^{u}\sum_{j=0}^{r}\alpha_{ij}$ whence the lemma can be proved.

Now, we shall prove theorem 3.

$$L(l+1,k+1)-L(l,k) = \sum_{r=1}^{l}\alpha_{0r}\psi(0,r)+\sum_{u=1}^{k}\alpha_{u0}\psi(u,0)+\sum_{u=1}^{k}\sum_{r=1}^{l}\alpha_{ur}\psi(u,r) = \sum_{i=0}^{k}\sum_{j=0}^{l}\alpha_{ij}\psi(i,j);$$

where putting $k' = k+1$, $l' = l+1$, $A = lk-1$, and $A' = l'k'-1$

$$\alpha_{00} = 0, \quad \alpha_{0j} = \frac{l-j}{lA} - \frac{l'-j}{l'A'}, \quad \alpha_{i0} = \frac{k-i}{kA} - \frac{k'-i}{k'A'}$$

and $\quad \alpha_{ij} = \frac{2(k-i)(l-j)}{lkA} - \frac{2(k'-i)(l'-j)}{l'k'A'}$, for $i \neq 0$, $j \neq 0$.

Then $\quad T_{ur} = \sum_{i=1}^{u}\alpha_{i0} + \sum_{j=1}^{r}\alpha_{0j} + \sum_{i=1}^{u}\sum_{j=1}^{r}\alpha_{ij}$, whence

$$2T_{ur} = -\frac{u(k+x-1)}{kA} + \frac{r(l+y-1)}{lA} + \frac{ur(k+x-1)(l+y-1)}{lkA} - \frac{u(k'+x)}{k'A'} - \frac{r(l'+y)}{l'A'}$$
$$- \frac{ur(k'+x)(l'+y)}{l'k'A'}$$

where $x = (k-u)$ and $y = (l-v)$, from which it can be proved that $T_{ur} > 0$, if $\sqrt{l} < k < l^2$, and $T_{kl}$ is always zero ; i.e., if $\sqrt{l} < k < l^2$, $\sum_{i=0}^{k}\sum_{j=0}^{l}\alpha_{ij}=0$

and $\sum_{i=0}^{u}\sum_{j=0}^{v}\alpha_{ij} \geqslant 0$.

Hence, applying lemma IV

$L(l+1,k+1)-L(l,k) \geqslant 0$, if (i) $\Delta_1\psi(u,l) \leqslant 0$, (ii) $\Delta_2\psi(u,l) \leqslant 0$ and (iii) $\Delta_1\Delta_2\psi(u,r) \geqslant 0$

whence  $L(l, k) \leqslant L(l+1, k+1) \leqslant ... \leqslant L(ml, nk)$ if (i) $\Delta_1\dot\psi \leqslant 0$, (ii) $\Delta_2\dot\psi \leqslant 0$,

(iii) $\Delta_1\Delta_2\dot\psi(u,v) \geqslant 0$, (iv) $(m-1)l = (n-1)k$, and (v) $\sqrt{l} \leqslant k \leqslant l^2$.

Hence the theorem.

## SECTION 3

We shall now find out a set of sufficient conditions under which the systematic sampling is more efficient than the stratified sampling.

*Theorem 4.* For all infinite populations, in which (i) $\delta_1{}^2 \rho(u,v) \geqslant 0$, and (ii) $\delta_2{}^2\rho (u,v) \geqslant 0$, $\sigma_{sy}{}^2 \leqslant \sigma_{st}{}^2$, for any size of the sample and $\sigma_{sy}{}^2 < \sigma_{st}{}^2$ unless equality holds in each of the above two cases.

*Proof.* As the relative efficiency $\sigma_{sy}{}^2/\sigma_{st}{}^2$ is independent of $v$, we can take it to be unity, i.e., take one sample from each of $mn$ strata.

From (1.13) for $v=1$, average $(\bar{x}_{sy} - \bar{x})^2$, the estimate of $\sigma_{sy}{}^2$ is

$$\frac{1}{mnlk}\left\{ \sum_{i=1}^{mn} \sum_{j=1}^{lk} (x_{ij}-\bar{x})^2 - \frac{1}{mn} \times \text{(average s.s. within systematic samples)} \right\} \quad ... \quad (3.1)$$

Let $x_{1j}, x_{2j}, ..., x_{mn, j}$ be the stratified samples from the $mn$ strata, $\bar{x}_{si}$ the mean ; $x_{ij}$ being the sample from the $i$th stratum.

Then $_j\bar{x}_{si} = \sum_{i=1}^{mn} x_{ij}/mn$. Let $\bar{x}_i$ be the mean of the $i$th stratum and $\bar{x}$ the grand mean.

Then $\bar{x}_i = \frac{1}{lk} \sum_{j=1}^{lk} x_{ij}$, $\bar{x} = \frac{1}{mn} \sum_{i=1}^{mn} \bar{x}_i$, $\sigma_{st}{}^2 = E(\bar{x}_{si}-\bar{x})^2$. Now, $\sum_{i=1}^{mn} (x_{ij}-\bar{x})^2 = \sum_{i=1}^{mn} (x_{ij}-\bar{x}_{si})^2 + mn(\bar{x}_{si}-\bar{x})^2$.

Averaging over $(lk)^{mn}$ samples we have

$$\frac{1}{lk} \sum_{i=1}^{mn} \sum_{j=1}^{lk} (x_{ij}-\bar{x})^2 = \text{(average s.s. within sample)} + mn(\bar{x}_{si}-\bar{x})^2.$$

$$\therefore \quad \sigma_{st}{}^2 = E(\bar{x}_{si}-\bar{x})^2 = \frac{1}{mnlk} E \left\{ \sum_{i=1}^{mn} \sum_{j=1}^{lk} (x_{ij}-\bar{x})^2 - \frac{1}{mn} \text{ (average s.s. within sample)} \right\} (3.2)$$

If $x_1, x_2, ..., x_{mn}$ be the sample, and $x$ its mean, then s.s. within sample is

$$\sum_{i=1}^{mn} (x_i-x)^2 = \frac{1}{mn} \sum_{\substack{i,j=1 \\ j > i}}^{mn} (x_i-x_j)^2 ; \quad ... \quad (3.3)$$

which is the average s.s within the sample, systematic or stratified.

Let $\Sigma_{sy}$ and $\Sigma_{st}$ be the expectations of the average s.s. within systematic and stratified samples respectively. Then $\sigma_{st}{}^2 \geqslant \sigma_{sy}{}^2$, if $\Sigma_{sy} \geqslant \Sigma_{st}$, from (3.1) and (3.2).

From (3.3), $\Sigma_{sy}$ and $\Sigma_{st}$ are given by $\frac{1}{mn}$ multiplied by the sum of terms like $E(x_i - z_j)^2$, for all combinations of pairs of strata, $(i, j$ denoting strata).

We can, therefore, compare $\Sigma_{sy}$ and $\Sigma_{st}$ term by term.

Let $E_{sy}$ and $E_{st}$ denote expectations for systematic and stratified samples respectively. Also let $x_i$, $x_j$ belong to $(i_1, i_2)$ and $(j_1, j_2)$ strata respectively, such that $(i_1 - j_1) = u$ and $(i_2 - j_2) = v$.

Then
$$E_{sy}(x_i - x_j)^2 = E_{sy}(x_i - \mu + \mu - x_j)^2 = 2\sigma^2\{1 - \rho(uk, vl)\}, \quad \ldots \ (3.4)$$

$x_j$ being determined when $x_i$ is fixed and *vice versa*.

For $E_{st}(x_i - x_j)^2$, $l^2 k^2$ combinations of pairs are possible in the present case. We are to take mean over all these combinations.

It can be seen that $(k - |i|)$ pairs are $(ku + i)$ distance apart and $(l - |j|)$ are $(lv + j)$ distance apart.

Averaging over $l^2 k^2$ combinations we have

$$E_{st}(x_i - x_j)^2 = 2\sigma^2\left\{1 - \frac{1}{k^2 l^2}\sum_{i=-(k-1)}^{(k-1)}\sum_{j=-(l-1)}^{(l-1)}(k - |i|)(l - |j|)\rho(vk + i, vl + j)\right\} \quad \ldots \ (3.5)$$

From (3.4) and (3.5), $(\Sigma_{sy} - \Sigma_{st})$ will be proportional to the sum of terms like

$$\frac{1}{k^2 l^2}\sum_{i=-(k-1)}^{(k-1)}\sum_{j=-(l-1)}^{(l-1)}(k - |i|)(l - |j|)\rho(vk + i, vl + j) - \rho(uk, vl) = \frac{1}{k^2 l^2}(T_1 + T_2 + T_3 + T_4),$$

where, $T_1 = \sum_{i=1}^{(k-1)}\sum_{j=1}^{(l-1)}(k - i)(l - j)\{\rho(uk + i, vl + j) + \rho(uk + i, vl - j)$
$$+ \rho(uk - i, vl + j) + \rho(uk - i, vl - j)\},$$

$T_2 = k\sum_{j=1}^{(l-1)}(l - j)\{\rho(uk, vl + j) + \rho(uk, vl - j)\}$, $T_3 = l\sum_{i=1}^{k-1}(k - i)\{\rho(uk + i, vl)$
$+ \rho(uk - i, vl)\}$, and $T_4 = -lk(lk - 1)\rho(uk, vl)$.

Let us subtract $4\rho(uk, vl)$ from the terms of $T_1$, and $2\rho(uk, vl)$ from terms of $T_2$ and $T_3$ each and add the balance to $T_4$.

Then we get : $T_1' = \sum_{i=1}^{k-1}\sum_{j=1}^{l-1}(k - i)(l - j)\{\rho(uk + i, vl + j) + \rho(uk + i, vl - j) + \rho(uk - i, vl + j)$
$$+ \rho(uk - i, vl - j) - 4\rho(uk, vl)\}$$

$$T_2' = k \sum_{j=1}^{l-1} (l-j)\{\rho(uk, vl+j) + (\rho(uk, vl-j) - 2\rho(uk,vl))\}, \text{ and}$$

$$T_3' = l \sum_{i=1}^{k-1} (k-i)\{\rho(uk+i, vl) + \rho(uk-i, vl) - 2\rho(uk,vl)\}$$

We are in all subtracting $lk(lk-1)\,\rho(uk, vl)$, and adding the same to $T_0$, we get zero as resultant.

Hence the expression is $\qquad\qquad T_1' + T_2' + T_3' \qquad\qquad$ ... (3.6)

If as usual $\delta^2 = \Delta^2 E^{-1} = (E^{\frac{1}{2}} - E^{-\frac{1}{2}})^2$, it can be proved that

$$\sum_{j=-(l-1)}^{(l-1)} (i-|j|)\delta^2 E^j = E^i + E^{-i} - 2 \; ;$$

Then $\quad \sum_{j=-(l-1)}^{(l-1)} (i-|j|)\delta^2\phi(u+j) = \phi(u+i) + \phi(u-i) - 2\phi(u) \qquad$ ... (3.7)

Hence in $T_1'$, $\quad \rho(uk+i, vl+j) + \ldots + \rho(uk-i, vl-j) - 4\rho(uk,vl)$

$$= (E_1^i + E_1^{-i} - 2)\{\rho(uk,vl+j) + \rho(uk,vl-j)\} + 2(E_2^j + E_2^{-j} - 2)\rho(uk, vl)$$

$$= \sum_{s=-(l-1)}^{(l-1)} (i-|s|)\{\delta_1^2\rho(uk+s, vl+j) + \delta_1^2\rho(uk+s,vl-j)\}$$

$$+ 2\sum_{t=-(j-1)}^{(j-1)} (j-|t|)\delta_2^2\rho(uk,vl+t), \text{ using (3.7)}$$

Similarly, in $T_2'$

$$\rho(uk, vl+j) + \rho(uk, vl-j) - 2\rho(uk, vl) = \sum_{t=-(j-1)}^{(j-1)} (j-|t|)\,\delta_2^2\,\rho(uk, vl+t),$$

and in $T_3'$

$$\rho(uk+i, vl) + \rho(uk-i, vl) - 2\rho(uk,vl) = \sum_{s=-(l-1)}^{(l-1)} (i-|s|)\,\delta_1^2\rho(uk+s, \, vl).$$

Thus $T_1'$, $T_2'$ and $T_3'$ are each positive, if (i) $\delta_1^2\rho(u,v) \geqslant 0$ and (ii) $\delta_2^2\rho(u,v) \geqslant 0$ for all values of $u$ and $v$.

Hence under these conditions $\Sigma_{\eta} > \Sigma_{\iota_1}$ whence $\sigma_{\gamma_1}{}^2 < \sigma_{\iota_1}{}^2$.

It is evident that the difference vanishes only when equality holds in each of the above three cases. Thus the theorem is proved.

From the proof of the theorems given in this section and in previous section, we can find conditions under which $\sigma_{\gamma}{}^2 > \sigma_{\gamma\gamma}{}^2$. In view of the fact that in certain cases, we are compelled to adopt systematic of stratified types of sampling [technique, we ought to find out the degree of inefficiency also, which, though not easily expressible, can be obtained from (1.14), (1.15) and (1.16) of section I.

Now, we shall consider one example.

Let $\quad \rho(u,v) = \sum\limits_{i=1}^{p} A_i e^{-\lambda_i |u| -\mu_i |v|}$, where $\sum\limits_{i=1}^{p} A_i = 1$, and $\lambda_i$ and $\mu_i$ are positive.

Evidently, $\rho(u,v) = \rho(-u,v)$ and we can consider only positive values of $u$ and $v$. This is supposed to be a natural condition (Ghosh, 1949).

Here, $\psi(u, v) = 2\rho(u, v)$. Then $\Delta_1 \psi(u, v) = 2 \sum\limits_{i=1}^{p} A_i e^{-\lambda_i u -\mu_i v} (e^{-\lambda_i} -1)$.

But when $\quad \lambda_i < 0, e^{-\lambda_i} < 1$; therefore, $\Delta_1 \psi(u, v) < 0$, for positive correlations.

Similar is the case with $\Delta_2 \psi(u,v)$.

But without knowing the values of $\lambda_i$ and $\mu_i$ we cannot say whether conditions of theorem I are satisfied or not. Only if $\lambda_i = \mu_i$ for all values of $i$, the corollary to theorem I is applicable.[*]

Again $\Delta_1 \Delta_2 \psi(u, v) = 2 A_i e^{-\lambda_i u -\mu_i v} (e^{-\lambda_i} -1)(e^{-\mu_i} -1)$, which is positive.

Hence, theorem 3 can be applied, i.e., stratified sampling with $\sqrt{l} < k < l^2$ and $(n-1)k = (m+1)l$ is more efficient than random sampling, though the universe is homogeneous as regards the expected means.

Also, $\delta_1{}^2 \psi(u, v) = 2\{\rho(u+1, v)+\rho(u-1, v) -2\rho(u, v)\}$

$$= 2 \sum\limits_{i=1}^{p} A_i e^{-\lambda_i u -\mu_i v} (e^{\lambda_1/2} -e^{-\lambda_1/2})^2, \text{ which is positive for positive } \rho.$$

---

[*]It has been also proved by the author that for $\Delta_1 \psi(u, v) < 0$, $\Delta_1 \psi(u, v) < 0$, $\delta_1{}^2 \psi(u, v) > 0$ and $\delta_1{}^2 \psi(u, v) > 0$ the relation $\sigma_\gamma{}^2 > \sigma_{\iota_1}{}^2 > \sigma_{\iota_2}{}^2$ holds good and $\sigma_\gamma{}^2 = \sigma_{\iota_1}{}^2$ if equality holds in each of the above four cases when $\sigma_{\iota_1}{}^2 = \sigma_{\iota_2}{}^2$ also will be true. These results will be published in a forthcoming issue of *Science and Culture*.

Thus in this example $\sigma_\gamma{}^2 > \sigma_{\iota_1}{}^2 > \sigma_{\iota_2}{}^2$.

Similarly, $\delta_u^2 \psi(u, v) > 0$,

Therefore applying theorem 4, systematic sampling is more efficient than stratified sampling.

While the paper was in press, it came to the notice of the author that some of the results discussed in this paper had been also obtained by Quenouille (1949). The present author published these results in Science and Culture in 1949.

Thanks are due to Prof. B. N. Ghosh for some valuable suggestions.

REFERENCES

1.  COCHRAN, W. G. (1946):  Relative accuracy of systematic and stratified random samples for a certain class of populations.  *Annals of Mathematical Statistics*, 17, 164-177.

2.  DAS, A. C. (1949):  Two dimensional systematic, stratified and random sampling.  *Proc. Ind. Sci. Congress*, 36th Session, Part III, 6.

3.  GHOSH, B. N. (1949):  On a particular type of natural field.  *Proc. Ind. Sci. Congress*, 36th Session, Part III, 7.

4.  MADOW, W. G. and L. H. (1944):  On the theory of systematic sampling, I.  *Annals of Mathematical Statistics*, 15, 1-24.

5.  QUENOUILLE, M. H. (1949):  Problems in plane sampling.  *Annals of Mathematical Statistics*, 20, 355-375.

6.  YATES, F (1946):  Systematic sampling.  *Phil. Trans, Roy.-Soc.*, 241(A), 345-377.

*Paper received : 2 February, 1949.*