# ESTIMATION OF OPTIMUM SAMPLE SIZE AND NUMBER OF REPLICATIONS IN SPLIT-SPLIT PLOT DESIGN

Md. Saiful Islam[1]

## Abstract

In field experiments, it is necessary to determine the optimum sample size as well as optimum number of replications if researchers have to use sampling techniques for collecting data from such experiments. Estimates of such optimum sample size and number of replications has been determined for split-split plot design minimizing the variance for a given cost of the experiment per treatment.

Key Words: Estimation, sample size, replications, split-split plot design

## Introduction

Frequently it is not possible to measure yield and yield contributing characteristics on the whole of each experimental unit. It may be desirable to estimate the characteristics on a random sample basis. Kempthorne (1952) showed the optimum sample size and number of replications in randomized complete block design (RCBD) for equal sampling from each cell having multiple observations. Islam *et al.* (2000) also showed the optimum sample size and number of replications in RCBD with unequal observation per cell. Federer (1963) and Islam (2001) has cited some example (optimum sample size was found 5 & 9 in completely randomized design & split plot design). Usual four-way classification in Split-Split Plot Design with more than one observation per cell and equal variance are considered. Estimation of variance components are obtained from analysis of variance technique. None has yet used for split-split plot design to find out the optimum sample size and number of replications for measuring yield and yield contributing characters from experimental fields.

The main objective of this study is to determine optimum sample size and number of replications for measuring yield and yield contributing characters based on equal sampling from each sub-sub plot of split-split plot design.

[1]Senior Scientific Officer, Statistics Section, BARI, Joydebpur, Gazipur 1701, Bangladesh.

## Method

Sample size depends on the variability associated with variable and the cost of reducing that variability. For such cases, it is necessary to choose optimum sample size and number of replications. Estimation of optimum sample size and number of replications are obtained by maximizing the information for a given cost.

Suppose that we have an experiment in r replications, p main plots, q sub-plots and s sub-sub plots and that the characteristic has been observed on a random sample of n equal sampling units from each plot (sub-sub plot). The observations may be denoted by $Y_{ijklm}$, where k denotes the replications (k=1,2,..............,r); i denotes the number of main plots (i=1,2,...........,p); j denotes the sub-plots (j=1,2,.........q); l denotes the sub-sub plots (1,2,......,s) and m denotes the sampling unit (m= 1,2,............,n) and assume the following model:

$$Y_{ijklm} = \mu + \rho_k + \alpha_i + \gamma_{ik} + \beta_j + (\alpha\beta)_{ij} + \gamma_{ijk} + \delta_l + (\alpha\delta)_{il} + (\beta\delta)_{il} + (\alpha\beta\delta)_{ijl} + \varepsilon_{ijkl} + \eta_{ijklm}$$

where   $\mu$ = the general mean

   $\rho_k$ = the effect of kth replication

   $\alpha_i$ = the main plot effect due to ith level of A

   $\gamma_{ik}$  = the main plot error

   $\beta_j$ = the sub plot effect due to jth level of B

   $(\alpha\beta)_{ij}$ = the interaction effect due to ith level of A and jth level of B

   $\gamma_{ijk}$ = the sub plot error

   $\delta_l$ = the sub-sub plot effect due to lth level of C

   $(\alpha\delta)_{il}$ = the interaction effect due to ith level of A and lth level of C

   $(\beta\delta)_{jl}$ = the interaction effect due to jth level of B and lth level of C

$(\alpha\beta\delta)_{ijl}$ = the interaction effect due to ith level of A, jth level of B and      lth level of C

   $\varepsilon_{ijkl}$ = the experimental (sub-sub plot) error

   $\eta_{ijklm}$ = the sampling error of the mth observation.

We suppose for the present purpose that the  $\eta_{ijklm}$ 'S are normally and independently distributed with variance  $\sigma_n^2$, the $\varepsilon_{ijkl}$' S  are normally and independently distributed with variance  $\sigma_e^2$, $\gamma_{ik}$' S  are normally and independently distributed with variance  $\sigma_\alpha^2$ and $\gamma_{ijk}$' S  are also normally and

independently distributed with variance $\sigma_b^2$. The errors $\eta_{ijklm}$, $\varepsilon_{ijkl}$, $\gamma_{ik}$ and $\gamma_{ijk}$ are independent since the sampling is done at random.

The least square estimates of parameters and errors are obtained as follows:

$$\hat{\mu} = \bar{y}.....$$

$$\hat{\alpha}_i = \bar{y}_i.... - \bar{y}......$$

$$\hat{\beta}_j = \bar{y}._j... - \bar{y}......$$

$$\hat{\rho}_k = \bar{y}.._k .. - \bar{y}.......$$

$$\left(\hat{\alpha}\hat{\beta}\right)_{ij} = \bar{y}_{ij}... - \bar{y}_i.... - \bar{y}._j... + \bar{y}......$$

$$\hat{\gamma}_{ik} = \bar{y}_i._k.. - \bar{y}_i.... - \bar{y}.._k.. + \bar{y}......$$

$$\hat{\gamma}_{ijk} = \bar{y}_{ijk}.. - \bar{y}_{ij}... - \bar{y}_i._k .. + \bar{y}_i....$$

$$\hat{\delta}_l = \bar{y}..._l. - \bar{y}.......$$

$$\left(\hat{\alpha}\hat{\delta}\right)_{il} = \bar{y}_i.._l. - \bar{y}_i.... - y...._l. + y......$$

$$\left(\hat{\beta}\hat{\delta}\right)_{jl} = \bar{y}._j._l. - \bar{y}._j... - y..._l. + y......$$

$$\left(\hat{\alpha}\hat{\beta}\hat{\delta}\right)_{ijl} = \bar{y}_{ij}._l. - \bar{y}_{ij}... - \bar{y}_i.._l. - \bar{y}._j._l. + \bar{y}_i.... + \bar{y}._j ... + \bar{y}..._l. - \bar{y}......$$

$$\hat{\varepsilon}_{ijkl} = \bar{y}_{ijkl}. - \bar{y}_{ijk}.. - \bar{y}_{ij}._l. + \bar{y}_{ij}...$$

$$\hat{\eta}_{ijklm} = y_{ijklm} - \bar{y}_{ijkl}.$$

Total SS can be partitioned into component SS as follows:

$$\text{Total SS} = \sum_i\sum_j\sum_k\sum_l\sum_m \left(y_{ijklm} - \bar{y}......\right)^2 = qrsn\sum_i\left(\bar{y}_i.... - \bar{y}.....\right)^2 +$$

$$pqsn\sum_i\left(\bar{y}.._k.. - \bar{y}......\right)^2 + qsn\sum_i\sum_k\left(\bar{y}_i._k.. - \bar{y}_i.... - \bar{y}.._k.. + \bar{y}......\right)^2 + prsn$$

$$\sum_j\left(\bar{y}._j.. - \bar{y}......\right)^2 +$$

$$rsn\sum_i\sum_j\left(\bar{y}_{ij}... - \bar{y}_i.... - \bar{y}._j... + \bar{y}.....\right)^2 + sn\sum_i\sum_j\sum_k\left(\bar{y}_{ijk}.. - \bar{y}_{ij}.. - \bar{y}_i._k.. + \bar{y}_i....\right)^2 +$$

$$pqrn\sum_l\left(\bar{y}..._l. - \bar{y}......\right)^2 + sn\sum_i\sum_j\sum_k\left(\bar{y}_{ijk}.. - \bar{y}_{ij}... - \bar{y}_i._k.. + \bar{y}_i....\right)^2 +$$

$$prn\sum_{j}\sum_{l}\left(\bar{y}_{\cdot j\cdot l\cdot}-\bar{y}_{\cdot j}...-\bar{y}...{}_{l\cdot}+\bar{y}.....\right)^{2}+rn\sum_{i}\sum_{j}$$

$$\sum_{l}\left(\bar{y}_{ij\cdot l\cdot}-\bar{y}_{ij}...-\bar{y}_{i}..{}_{l\cdot}+\bar{y}_{i}....+\bar{y}_{\cdot j}...+\bar{y}...{}_{l\cdot}-\bar{y}....\right)^{2}+$$

$$n\sum_{i}\sum_{j}\sum_{k}\sum_{l}\left(\bar{y}_{ijkl\cdot}-\bar{y}_{ijk}..-\bar{y}_{ij\cdot l}.+\bar{y}_{ij}...\right)^{2}+\sum_{i}\sum_{j}\sum_{k}\sum_{l}\sum_{m}\left(y_{ijklm}-\bar{y}_{ijkl\cdot}\right)^{2}$$

= Main plot SS + Sub plot SS + Error$_1$ SS + Replication SS + Interaction (MP∞SP) SS + Error$_2$ SS + Sub-Sub plot SS + Interaction (MP × SSP) SS + Interaction (SP × SSP) SS + Interaction (MP × SP×SSP) SS + Experimental error SS + Sampling error SS

= SS (A) + SS (B) + SS (Error)$_1$ + SS (Rep) + SS(A×B) + SS(Error)$_2$ + SS(C) + SS(A×C) + SS (B×C) + SS(A×B×C) + SS (Experimental error) + SS (Sampling error).

The degrees of freedom for different sum of squares are:

| Sum of Squares | degrees of freedom |
|---|---|
| Total | (pqrsn-1) |
| Replication | (r-1) |
| Main plot (A) | (p-1) |
| Error-1 | (r-1) (p-1) |
| Sub plot (B) | (q-1) |
| Interaction (A×B) | (p-1) (q-1) |
| Error-2 | p(q-1)(r-1) |
| Sub-sub plot (C) | (s-1) |
| Interaction (A×C) | (p-1)(s-1) |
| Interaction (B×C) | (q-1)(s-1) |
| Interaction (A×B×C) | (p-1)(q-1)(s-1) |
| Experimental error (Error-3) | pq(s-1)r-1) |
| Sampling error | pqrs(n-1) |

The analysis of variance of table for the analysis on a sample basis is given below:

**Table 1.**

| Source of variation | d. f | Sum of squares | Mean square | Expectation of mean square* |
|---|---|---|---|---|
| Replication | (r-1) | $pqsn\sum_k\left(\bar{y}_{..k..}-\bar{y}_{......}\right)^2$ | R | |
| Main plot (A) | (p-1) | $qrsn\sum_i\left(\bar{y}_{i.....}-\bar{y}_{......}\right)^2$ | MP | |
| Error-1 | (r-1)(p-1) | $qsn\sum_i\sum_k\left(\bar{y}_{i.k..}-\bar{y}_{i....}-\bar{y}_{..k..}+\bar{y}_{......}\right)^2$ | $E_1$ | $qsn\sigma^2_a+sn\sigma^2_b+n\sigma^2_e+\sigma^2_n$ |
| Sub plot (B) | (q-1) | $prsn\sum_j\left(\bar{y}_{.j....}-\bar{y}_{......}\right)^2$ | SP | |
| Interaction (A×B) | (p-1)(q-1) | $rsn\sum_i\sum_j\left(\bar{y}_{ij...}-\bar{y}_{i....}-\bar{y}_{.j...}+\bar{y}_{......}\right)^2$ | MP∞SP | |
| Error-2 | p(q-1)(r-1) | $sn\sum_i\sum_j\sum_k\left(\bar{y}_{ijk.}-\bar{y}_{ij...}-\bar{y}_{i.k..}+\bar{y}_{i....}\right)^2$ | $E_2$ | $sn\sigma^2_b+n\sigma^2_e+\sigma^2_n$ |
| Sub-sub plot (C) | (s-1) | $pqrn\sum_l\left(\bar{y}_{...l.}-\bar{y}_{......}\right)^2$ | SSP | |
| Interaction (A×C) | (p-1)(s-) | $qrn\sum_i\sum_l\left(\bar{y}_{i..l.}-\bar{y}_{i....}-\bar{y}_{...l.}+\bar{y}_{......}\right)^2$ | MP∞SSP | |
| Interaction (B×C) | (q-1)(s-1) | $prn\sum_j\sum_l\left(\bar{y}_{.j.l.}-\bar{y}_{.j...}-\bar{y}_{...l.}+\bar{y}_{......}\right)^2$ | SP∞SSP | |
| Interaction (A×B×C) | (p-1)(q-1)(s-1) | $rn\sum_i\sum_j\sum_l\left(\bar{y}_{ij.l.}-\bar{y}_{ij...}-\bar{y}_{i..l.}-\bar{y}_{.j.l.}+\bar{y}_{i....}+\bar{y}_{.j...}+\bar{y}_{...l.}-\bar{y}_{......}\right)^2$ | MP x SP x SSP | |
| Experimental error (Error-3) | pq(s-1)(r-1) | $n\sum_i\sum_j\sum_k\sum_l\left(\bar{y}_{ijkl}-\bar{y}_{ijk.}-\bar{y}_{ij.l}+\bar{y}_{ij..}\right)^2$ | $E_3$ | $n\sigma^2_e+\sigma^2_n$ |
| Sampling error | pqrs(n-1) | $\sum_i\sum_j\sum_k\sum_l\sum_m\left(y_{ijklm}-y_{ijkl}\right)^2$ | S | $\sigma^2_n$ |
| Total | (pqrsn-1) | $\sum_i\sum_j\sum_k\sum_l\sum_m\left(y_{ijklm}-\bar{y}_{......}\right)^2$ | - | - |

\* The expression are given only for error mean square, sources Kempthorne (1952, pp213) and Islam *et al.*, (2000, pp 91)

R, MP, $E_1$ SP, MPxSP, $E_2$, SSP, MPxSSP, SPxSSP, MPxSPxSSP, $E_3$, & S are the mean squares of replication, main plot, error$_1$, sub plot, MPxSP interaction, error$_2$, sub-sub plot, MPxSSP interaction, SPxSSP interaction, MPxSPxSSP interaction, experimental error and sampling error, respectively.

**Estimation**

The component of variances : $\sigma^2_n, \sigma^2_e, \sigma^2_b$ *and* $\sigma^2_a$ are estimated as

$$\hat{\sigma}^2_n = S, E_3 = n\hat{\sigma}^2_e + \hat{\sigma}^2_e, E_2 = sn\hat{\sigma}^2_b + n\hat{\sigma}^2_e + \hat{\sigma}^2_n \ and \ E_1$$
$$= qsn\hat{\sigma}^2_a + sn\hat{\sigma}^2_b + n\hat{\sigma}^2_e + \hat{\sigma}^2_n$$

or, $\hat{\sigma}^2_e = \dfrac{E_3 - S}{n}, \hat{\sigma}^2_b = \dfrac{E_2 - E_3}{ns}, \hat{\sigma}^2_a = \dfrac{E_1 - E_2}{snq}$

Now $\begin{aligned} \hat{\delta}_l &= \bar{y}...{}_{l}. - \bar{y}..... \\ \hat{\delta}'_l &= \bar{y}...{}'_{l}. - \bar{y}..... \end{aligned} \qquad {}_l \neq l'$

$$V(\hat{\delta}_l - \hat{\delta}'_l) = V(\bar{y}...._l. - \bar{y}...l'.) = V(\bar{y}...{}_l.) + V(y...{}'_l.)$$

$$V(\bar{y}...{}_l.) = V(\bar{\gamma}..) + V(\bar{\gamma}...) + V(\bar{e}..._l) + V(\bar{\eta}..._l.)$$

(Where $V(\bar{y}..._l.)$ = variance of treatment mean comes sub-sub plots (treatments), formula cited from Kempthorne (1952) & Islam *et al.* (2000))

$$= \frac{\sigma_a{}^2}{pr} + \frac{\sigma^2_b}{pqr} + \frac{\sigma_e{}^2}{pqr} + \frac{\sigma_n{}^2}{npqr} = \frac{n(q\sigma_a{}^2 + \sigma^2_b + \sigma_e{}^2) + \sigma_n{}^2}{npqr}$$

$$V(\hat{\delta}_l - \hat{\delta}l') = 2\frac{n(q\sigma_a{}^2 + \sigma^2_b + \sigma_e{}^2) + \sigma_n{}^2}{npqr}$$

The estimated variacne of difference between two treatment (sub-sub plot) means on a per sample basis will be proportional to $\dfrac{n(q\hat{\sigma}_a{}^2 + \hat{\sigma}^2_b + \hat{\sigma}_e{}^2) + \hat{\sigma}_n{}^2}{npqr} = \dfrac{E_1 - E_3 + sE_3}{npqrs}$ Similarly, the estimated variance of treatment comparison with $r'$ replications, $p'$ main plot, $q'$ sub plot and

$n'$ samples per plot will be proportional to $\dfrac{\hat{\sigma}^2_n + n'(q'\hat{\sigma}^2_a + \hat{\sigma}^2_b + \hat{\sigma}^2_e)}{r'p'q'n'}$

$$= \frac{\hat{\sigma}^2_n}{r'p'q'n'} + \frac{\hat{\sigma}^2_a}{r'p'} + \frac{\hat{\sigma}^2_b}{p'q'r'} + \frac{\hat{\sigma}^2_e}{r'p'q'}$$

$$= \frac{S}{r'p'q'n'} + \frac{E_1 - E_2}{nsqr'p'} + \frac{E_2 - E_3}{snp'q'r'} + \frac{E_3 - S}{nr'p'q'} = \frac{E_1 - E_2}{snqp'r'} + \frac{E_2 - E_3}{snp'q'r'} + \frac{E_2}{np'q'r'} + \frac{S}{r'p'q'}\left[\frac{1}{n'} - \frac{1}{n}\right]$$

The relative information (RI) for $r'$ replication, $p'$ main plot, $q'$ sub plot and $n'$ samples per plot over r replication, p main plot, q sub plot and n samples per plot may then be estimated

$$RI = \frac{1}{\left[\dfrac{q'n'(E_1 - E_2) + qn'(E_2 - E_3) + sqn'E_3 + sqS(n - n')}{sqnr'p'q'n'}\right]} \bigg/ \frac{1}{\left[\dfrac{E_1 - E_3 + sE_3}{spqrn}\right]}$$

$$= \frac{n'p'q'r'(E_1 - E_3 + sE_s)}{pr\left[n'q'(E_1 - E_2) + qn'(E_2 - E_3) + n'qsE_3 + (n - n')Sqs\right]}.$$

$$= \frac{\{n'p'q'r'(E_1 - E_3 + sE_3)\}/n'}{\left[pr\{n'q'(E_1 - E_2) + qn'(E_2 - E_3) + n'sqE_3 + (n - n')Sqs\}\right]/n'}$$

As $n' \to \infty$, $RI = \dfrac{p'q'r'(E_1 - E_3 + sE_3)}{pr\{q'(E_1 - E_2) + q(E_2 - E_3) + sq(E_3 - S)\}}$

and this is the maximum that can be achieved by increasing the number of sub sampling units. It will often be uneconomical to increase the number of sub sampling units i.e. n and it is necessary to consider the cost of the various operations to determine optimum n.

For the cost function, let us consider

C = pre-harvesting cost per plot

$C_n$ = harvesting and post-harvesting cost per sample

With r replications p main plots, q sub plots and n sample size per plot, the cost of the experiment per treatment (sub-sub plot) is $(rpqC + rpqnC_n) = rpq(C + nC_n)$

The information on each treatment mean, assuming the variance components to be known is

$$\frac{pqrn}{\sigma^2{}_n + n(q\sigma^2{}_a + \sigma^2{}_e + \sigma^2{}_b)} \tag{1}$$

We shall choose r, p, q and n to maximize the information in (1) for a given cost per treatment $C_0 = rpq(C + nC_n)$ \hfill (2)

From (2), $r = \dfrac{C_0}{pq(C + nC_n)}$

Putting this value in equation (1), we get the information expressed in terms of costs as

$$\frac{nC_0}{(C+nC_n)\left\{\sigma^2_n + n\left(q\sigma^2_a + \sigma^2_b + \sigma^2_e\right)\right\}} \tag{3}$$

Then maximization of (1) is equivalent to maximization of (3)

The reciprocal of (3) is given by

$$\frac{1}{C_0}\left\{\frac{C\sigma^2 n}{n} + C_n\sigma^2 n + C\left(q\sigma^2_a + \sigma^2_b + \sigma^2_e\right) + nC_n\left(q\sigma^2_a + \sigma^2_b + \sigma^2_e\right)\right\} \tag{4}$$

Then the maximization of (3) is equivalent to the minimization (4)

Differentiating (4) with respect to n and equating to zero, we get

$$-\frac{C\sigma^2_n}{n^2} + C_n\left(q\sigma^2_a + \sigma^2_b + \sigma^2_e\right) = 0,$$

or
$$n^2 = \frac{C\sigma^2_n}{C_n\left(q\sigma^2_a + \sigma^2_b + \sigma^2_e\right)},$$

which yields the optimum sample size as $n_{(opt)} = \sqrt{\dfrac{C\sigma^2_n}{C_n\left(q\sigma^2_a + \sigma^2_b + \sigma^2_e\right)}}$

The optimum number of replication is obtained substituting the optimum value $n_{(opt)}$ of n in equation (2) $r_{(opt)} = \dfrac{C_0}{pq\left[C + \sqrt{\dfrac{CC_n\sigma^2_n}{\left(q\sigma^2_a + \sigma^2_b + \sigma^2_e\right)}}\right]}$

## Concluding remarks

The optimum sample size as well as optimum number of replications depend on unknown error variances $\sigma^2_n, \sigma^2_e, \sigma^2_b \text{ and } \sigma^2_a$. Unbiased estimators of the error variances are given from Table 1. Using such unbiased estimators of $\sigma^2_n, \sigma^2_e, \sigma^2_b \text{ and } \sigma^2_a$, estimates optimum sample size $(\hat{n})$ and optimum number of replications $(\hat{r})$ can be obtained. The costs C and $C_n$ can be obtained from similar experiments already conducted in real fields.

## References

Kempthorne, O. 1952. Design and analysis of experiment, John Wiley and Sons, Inc, New York, USA.

Islam M. S., K. Sen and K. Rahim. 2000. Estimation of optimum sample size and number of replications in RCBD with unequal observation per cell. *Dhaka University Journal of Science* **48**(1): 89-94.

Islam, M. S. 2001. Statistical development of field plot technique for agronomic experiments. Unpublished Ph. D. thesis, Dhaka University, Dhaka, Bangladesh.

Federer W. T. 1963. Experimental design. Theory and application, Oxford and IBH Publishing Co. New Delhi, India.