

# ON THE EFFICIENT DESIGN OF STATISTICAL INVESTIGATIONS

BY ABRAHAM WALD

*Columbia University.*

**1. Introduction.** A theory of efficient design of statistical investigations has been developed by R. A. Fisher<sup>1</sup> and his followers mainly in connection with agricultural experimentation. However, the same methods can be applied to other fields also. All statistical designs treated in the aforementioned theory refer to problems of testing linear hypotheses. By testing a linear hypothesis we mean the following problem: Let  $y_1, \dots, y_N$  be  $N$  independently and normally distributed variates with a common variance  $\sigma^2$ . It is assumed that the expected value of  $y_\alpha$  is given by

$$(1) \quad E(y_\alpha) = \beta_1 x_{1\alpha} + \beta_2 x_{2\alpha} + \dots + \beta_p x_{p\alpha} \quad (\alpha = 1, \dots, N)$$

where the quantities  $x_{i\alpha}$  ( $i = 1, \dots, p; \alpha = 1, \dots, N$ ) are known constants and  $\beta_1, \dots, \beta_p$  are unknown constants. The coefficients  $\beta_1, \dots, \beta_p$  are called the population regression coefficients of  $y$  on  $x_1, x_2, \dots$ , and  $x_p$ , respectively. The hypothesis that the unknown regression coefficients  $\beta_1, \dots, \beta_p$  satisfy a set of linear equations

$$(2) \quad g_{i1}\beta_1 + \dots + g_{ip}\beta_p = g_i \quad (i = 1, \dots, r; r \leq p),$$

is called a linear hypothesis. The problem under consideration is that of testing the hypothesis (2) on the basis of the observed values  $y_1, \dots, y_N$ .

In many cases the experimenter has a certain amount of freedom in the choice of the values  $x_{i\alpha}$ . The efficiency of the test is greatly affected by the values of  $x_{i\alpha}$ . The statistical investigation is efficiently designed if the values  $x_{i\alpha}$  are chosen so that the sensitivity of the test is maximized. Let us illustrate this by a simple example. Suppose that  $x$  and  $y$  have a bivariate normal distribution and we want to test the hypothesis that the regression coefficient  $\beta$  of  $y$  on  $x$  has a particular value  $\beta_0$ . Suppose, furthermore, that the test has to be carried out on the basis of  $N$  pairs of observations  $(x_1, y_1), \dots, (x_N, y_N)$ , where the experiments are performed in such a way that  $x_1, \dots, x_N$  are not random variables but have predetermined fixed values. It is known that the variance of the least square estimate  $b$  of  $\beta$  is inversely proportional to  $\sum_{\alpha=1}^N (x_\alpha - \bar{x})^2$  where  $\bar{x} = (x_1 + \dots + x_N)/N$ . Hence, if we can freely choose the values  $x_1, \dots, x_N$  in a certain domain  $D$ , the greatest sensitivity of the test will be achieved by choosing  $x_1, \dots, x_N$  so that  $\sum (x_\alpha - \bar{x})^2$  becomes a maximum.

In the next section we will introduce a measure of the efficiency of the design

<sup>1</sup> See for instance R. A. FISHER, *The Design of Experiments*, Oliver and Boyd, London, 1935.

of a statistical investigation for testing a linear hypothesis. In sections 3 and 4 it will be shown that some well known experimental designs, used widely in agricultural experimentation, are most efficient in the sense of the definition given in section 2.

**2. A measure of the efficiency of the design of a statistical investigation for testing a linear hypothesis.** The hypothesis (2) can be reduced by a suitable linear transformation to the canonical form

$$(3) \quad \beta_1 = \beta_2 = \cdots = \beta_r = 0, \quad (r \leq p).$$

Hence, we can restrict ourselves without loss of generality to the consideration of the hypothesis (3).

Denote  $\sum_{\alpha=1}^N x_{i\alpha} x_{j\alpha}$  by  $a_{ij}$  and let the matrix  $\|c_{ij}\|$  be the inverse of the matrix  $\|a_{ij}\|$  ( $i, j = 1, \dots, p$ ). Denote by  $b_i$  the least square estimate of  $\beta_i$  ( $i = 1, \dots, p$ ). It is known that the estimates  $b_1, \dots, b_p$  have a joint normal distribution with mean values  $\beta_1, \dots, \beta_p$ , respectively. It is furthermore known that the covariance of  $b_i$  and  $b_j$  is equal to  $c_{ij}\sigma^2$ . The statistic used for testing the hypothesis (3) is given by

$$(4) \quad F = \frac{N - p}{r} \frac{\sum_{i=1}^r \sum_{m=1}^r a_{im}^* b_i b_m}{\sum_{\alpha=1}^N (y_\alpha - b_1 x_{1\alpha} - \cdots - b_p x_{p\alpha})^2}$$

where  $\|a_{im}^*\|$  is the inverse of  $\|c_{lm}\|$  ( $l, m = 1, \dots, r$ ). The statistic  $F$  has the  $F$ -distribution with  $r$  and  $N - p$  degrees of freedom. The critical region for testing the hypothesis (3) is given by the inequality

$$(5) \quad F \geq F_0,$$

where the constant  $F_0$  is determined so that the probability that  $F \geq F_0$  (calculated under the assumption that (3) holds) is equal to the level of significance we wish to have.

It is known that the power function<sup>2</sup> of the critical region (5) depends only on the single parameter

$$(6) \quad \lambda = \frac{1}{\sigma^2} \sum_{i=1}^r \sum_{m=1}^r a_{im}^* \beta_i \beta_m.$$

Furthermore this power function is a monotonically increasing function of  $\lambda$ . The coefficients  $a_{im}^*$  are functions of the quantities  $x_{i\alpha}$  ( $i = 1, \dots, p$ ;  $\alpha = 1, \dots, N$ ). The choice of the values  $x_{i\alpha}$  ( $i = 1, \dots, p$ ;  $\alpha = 1, \dots, N$ ) is the better the greater the corresponding value of  $\lambda$ . If  $r = 1$ , the expression  $\lambda$

<sup>2</sup> See for instance P. C. TANG, "The power function of the analysis of variance tests," *Stat. Res. Mem.*, Vol. II, 1938.

reduces to  $\frac{1}{\sigma^2} a_{11}^* \beta_1^2$ . Hence, if  $r = 1$ , we maximize  $\lambda$  by maximizing  $a_{11}^*$ . Since  $a_{11}^* = 1/c_{11}$ , we maximize  $\lambda$  by minimizing  $c_{11}$ . Thus, if  $r = 1$ , we can say that we obtain the most powerful test by minimizing  $c_{11}$ , i.e. by minimizing the variance of  $b_1$ . If  $r > 1$ , the difficulty arises that no set of values  $x_{i\alpha}$  ( $i = 1, \dots, p; \alpha = 1, \dots, N$ ) can be found for which  $\lambda$  becomes a maximum irrespective of the values of the unknown parameters  $\beta_1, \dots, \beta_r$ . Hence, if  $r > 1$ , we have to be satisfied with some compromise solution. For this purpose let us consider the unit sphere

$$(7) \quad \beta_1^2 + \dots + \beta_r^2 = 1,$$

in the space of the parameters  $\beta_1, \dots, \beta_r$ . It is known that the smallest root in  $\rho$  of the determinantal equation

$$(8) \quad \begin{vmatrix} a_{11}^* - \rho & a_{12}^* & \dots & a_{1r}^* \\ a_{21}^* & a_{22}^* - \rho & \dots & a_{2r}^* \\ \vdots & \vdots & \ddots & \vdots \\ a_{r1}^* & a_{r2}^* & \dots & a_{rr}^* - \rho \end{vmatrix} = 0,$$

is equal to the minimum value of  $\sigma^2 \lambda$  on the unit sphere (7). Similarly the greatest root of (8) is equal to the maximum value of  $\sigma^2 \lambda$  on the sphere (7). The compromise solution of maximizing the smallest root of (8) seems to be a very reasonable one. However, for the sake of certain mathematical simplifications, we propose to maximize the product of the  $r$  roots of (8). Since the product of the roots of (8) is equal to the determinant

$$(9) \quad \begin{vmatrix} a_{11}^* & \dots & a_{1r}^* \\ \vdots & \ddots & \vdots \\ a_{r1}^* & \dots & a_{rr}^* \end{vmatrix},$$

we have to maximize the determinant (9). The value of the determinant  $|c_{lm}|$  ( $l, m = 1, \dots, r$ ) is the reciprocal of that of (9). Hence we maximize (9) by minimizing the determinant  $|c_{lm}|$ . The generalized variance of the set of variates  $b_1, \dots, b_r$  is equal to the product of  $\sigma^{2r}$  and the determinant  $|c_{lm}|$ . Thus, our result can be expressed as follows: The optimum choice of the values of  $x_{i\alpha}$  is that for which the generalized variance of the variates  $b_1, \dots, b_r$  becomes a minimum.

Any set of  $pN$  values  $x_{i\alpha}$  ( $i = 1, \dots, p; \alpha = 1, \dots, N$ ) can be represented by a point in the  $pN$ -dimensional Cartesian space. Denote by  $D$  the set of all points in the  $pN$ -dimensional space which we are free to choose. If  $N$  is fixed and if any point of  $D$  can be equally well chosen, the following two definitions seem to be appropriate:

**DEFINITION 1.** Denote by  $c$  the minimum value of the determinant  $|c_{lm}|$  ( $l, m = 1, \dots, r$ ) in the domain  $D$ . Then the ratio  $c/|c_{lm}|$  is called the efficiency of the design of the statistical investigation for testing the hypothesis (3).

**DEFINITION 2.** The design of the statistical investigation for testing the hypothesis (3) is said to be most efficient if its efficiency is equal to 1.

**3. Efficiency of the Latin square design.** A widely used and important design in agricultural experimentation is the so-called Latin square. Suppose we wish to find out by experimentation whether there is any significant difference among the yields of  $m$  different varieties  $v_1, \dots, v_m$ . For this purpose the experimental area is subdivided into  $m^2$  plots lying in  $m$  rows and  $m$  columns and each plot is assigned to one of the varieties  $v_1, \dots, v_m$ . If each variety appears exactly once in each row and exactly once in each column, we have a Latin square arrangement. Denote by  $y_{ijk}$  the yield of the variety  $v_k$  on the plot which lies in the  $i$ -th row and  $j$ -th column. The subscript  $k$  is, of course, a single valued function of the subscripts  $i$  and  $j$ , since to each plot only one variety is assigned. The following assumptions are made: the variates  $y_{ijk}$  are independently and normally distributed with a common variance  $\sigma^2$  and the expected value of  $y_{ijk}$  is given by

$$(10) \quad E(y_{ijk}) = \mu_i + \nu_j + \rho_k.$$

The parameters  $\sigma^2$ ,  $\mu_i$ ,  $\nu_j$  and  $\rho_k$  are unknown. The hypothesis to be tested is the hypothesis that variety has no effect on yield, i.e.

$$(11) \quad \rho_1 = \rho_2 = \dots = \rho_k.$$

We associate the positive integer  $\alpha(i, j) = (i - 1)m + j$  with the plot which lies in the  $i$ -th row and  $j$ -th column. ( $i, j = 1, \dots, m$ ). It is clear that for any positive integer  $\alpha \leq m^2$  there exists exactly one plot, i.e. exactly one pair of values  $i$  and  $j$ , such that  $\alpha = \alpha(i, j)$ . In the following discussions the symbol  $y_\alpha$  ( $\alpha = 1, \dots, m^2$ ) will denote the yield  $y_{ijk}$  where the indices  $i$  and  $j$  are determined so that  $\alpha(i, j) = \alpha$ . The plot in the  $i$ -th row and  $j$ -th column will be called the  $\alpha$ -th plot where  $\alpha = \alpha(i, j)$ .

We define the symbols  $t_{i\alpha}$ ,  $u_{j\alpha}$ ,  $z_{k\alpha}$  ( $i, j, k = 1, \dots, m$ ;  $\alpha = 1, \dots, m^2$ ), as follows:  $t_{i\alpha} = 1$  if the  $\alpha$ -th plot lies in the  $i$ -th row, and  $t_{i\alpha} = 0$  otherwise. Similarly  $u_{j\alpha} = 1$  if the  $\alpha$ -th plot lies in the  $j$ -th column, and  $u_{j\alpha} = 0$  otherwise. Finally  $z_{k\alpha} = 1$  if the  $k$ -th variety is assigned to the  $\alpha$ -th plot, and  $z_{k\alpha} = 0$  otherwise. Then equation (10) can be written as

$$(12) \quad E(y_\alpha) = \mu_1 t_{1\alpha} + \dots + \mu_m t_{m\alpha} + \nu_1 u_{1\alpha} + \dots + \nu_m u_{m\alpha} + \rho_1 z_{1\alpha} + \dots + \rho_m z_{m\alpha}.$$

Denote the arithmetic means  $\frac{1}{m^2} \sum_{\alpha=1}^{m^2} t_{i\alpha}$ ,  $\frac{1}{m^2} \sum_{\alpha=1}^{m^2} u_{i\alpha}$ , and  $\frac{1}{m^2} \sum_{\alpha=1}^{m^2} z_{i\alpha}$  by  $\bar{t}_i$ ,  $\bar{u}_i$  and  $\bar{z}_i$  respectively. Let  $t'_{i\alpha} = t_{i\alpha} - \bar{t}_i$ ,  $u'_{i\alpha} = u_{i\alpha} - \bar{u}_i$ ,  $z'_{i\alpha} = z_{i\alpha} - \bar{z}_i$ ,  $\mu'_i = \mu_i - \mu_m$ ,  $\nu'_i = \nu_i - \nu_m$  and  $\rho'_i = \rho_i - \rho_m$  for  $i = 1, \dots, m - 1$ . Let furthermore  $w_\alpha = 1$  for  $\alpha = 1, \dots, m^2$ . Then we have

$$(13) \quad \begin{cases} t_{i\alpha} = t'_{i\alpha} + \bar{t}_i w_\alpha; & u_{i\alpha} = u'_{i\alpha} + \bar{u}_i w_\alpha; & z_{i\alpha} = z'_{i\alpha} + \bar{z}_i w_\alpha; \\ & (i = 1, \dots, m - 1) \\ t_{m\alpha} = (1 - \bar{t}_1 - \dots - \bar{t}_{m-1})w_\alpha - t'_{1\alpha} - \dots - t'_{m-1,\alpha}, \\ u_{m\alpha} = (1 - \bar{u}_1 - \dots - \bar{u}_{m-1})w_\alpha - u'_{1\alpha} - \dots - u'_{m-1,\alpha}, \\ z_{m\alpha} = (1 - \bar{z}_1 - \dots - \bar{z}_{m-1})w_\alpha - z'_{1\alpha} - \dots - z'_{m-1,\alpha}. \end{cases}$$

From (12) and (13) we obtain

$$(14) \quad E(y_\alpha) = \xi w_\alpha + \sum_{i=1}^{m-1} \mu'_i t'_{i\alpha} + \sum_{i=1}^{m-1} \nu'_i u'_{i\alpha} + \sum_{i=1}^{m-1} \rho'_i z'_{i\alpha}$$

where

$$\xi = \sum_{i=1}^{m-1} \mu'_i \bar{t}_i + \sum_{i=1}^{m-1} \nu'_i \bar{u}_i + \sum_{i=1}^{m-1} \rho'_i \bar{z}_i + \mu_m + \nu_m + \rho_m.$$

The hypothesis (11) can be written as

$$(15) \quad \rho'_1 = \rho'_2 = \dots = \rho'_{m-1} = 0.$$

This is a linear hypothesis in canonical form as given in (3). The values  $z'_{i\alpha}$  ( $i = 1, \dots, m-1$ ;  $\alpha = 1, \dots, m^2$ ) depend on the way in which the varieties  $v_1, \dots, v_m$  are assigned to the  $m^2$  plots. We will show that we obtain a most efficient design if we distribute the varieties over the  $m^2$  plots in a Latin square arrangement, i.e. if each variety appears exactly once in each row and exactly once in each column.

Let  $q_{1\alpha} = w_\alpha$ ,  $q_{i+1,\alpha} = t'_{i\alpha}$  ( $i = 1, \dots, m-1$ ),  $q_{m+j,\alpha} = u'_{j\alpha}$  ( $j = 1, \dots, m-1$ ) and  $q_{2m-1+k,\alpha} = z'_{k\alpha}$  ( $k = 1, \dots, m-1$ ). Denote  $\sum_{\alpha=1}^{m^2} q_{i\alpha} q_{j\alpha}$  by  $a_{ij}$  ( $i, j = 1, 2, \dots, 3m-2$ ) and let the matrix  $\|c_{ij}\|$  be the inverse of the matrix  $\|a_{ij}\|$  ( $i, j = 1, \dots, 3m-2$ ). Let us denote by  $\Delta$  the determinant  $|a_{ij}|$  ( $i, j = 1, \dots, 3m-2$ ), by  $\Delta_1$  the determinant  $|a_{ij}|$  ( $i, j = 1, \dots, 2m-1$ ), by  $\Delta_2$  the determinant  $|a_{ij}|$  ( $i, j = 2m, \dots, 3m-2$ ) and  $\Delta'_2$  the determinant  $|c_{ij}|$  ( $i, j = 2m, \dots, 3m-2$ ). We have to show that for the Latin square arrangement  $\Delta'_2$  becomes a minimum. From a known theorem<sup>3</sup> about determinants it follows that

$$(16) \quad \Delta'_2 = \Delta_1 / \Delta.$$

Hence, we have merely to show that  $\Delta / \Delta_1$  becomes a maximum for the Latin square arrangement. Denote by  $\bar{\Delta}$ ,  $\bar{\Delta}_1$  and  $\bar{\Delta}_2$  the values taken by  $\Delta$ ,  $\Delta_1$  and  $\Delta_2$ , respectively, in the case of a Latin square arrangement. Since, for the Latin square arrangement, as is known,

$$\sum_{\alpha=1}^{m^2} z'_{k\alpha} u'_{j\alpha} = \sum_{\alpha=1}^{m^2} z'_{k\alpha} t'_{i\alpha} = \sum_{\alpha=1}^{m^2} z'_{k\alpha} w_\alpha = 0 \quad (i, j, k = 1, \dots, m-1)$$

we have

$$(17) \quad \frac{\bar{\Delta}}{\bar{\Delta}_1} = \bar{\Delta}_2.$$

Since the matrix  $\|a_{ij}\|$  ( $i, j = 1, \dots, 3m-2$ ) is positive definite we have

$$(18) \quad \frac{\Delta}{\Delta_1} \leq \Delta_2.$$

<sup>3</sup> See M. BÔCHER, *Introduction to Higher Algebra*, 1931, pp. 31.

Because of (17) and (18) the Latin square design is proved to be most efficient if we show that  $\Delta_2 \leq \bar{\Delta}_2$ .

Denote by  $\Delta_2^*$  the  $m$ -rowed determinant  $|a_{ij}|$  ( $i, j = 1, 2m, 2m+1, \dots, 3m-2$ ). Since  $a_{ij} = 0$  for  $j \neq 1$ , we have

$$(19) \quad \Delta_2^* = a_{11}\Delta_2 = m^2\Delta_2.$$

Denote  $\sum_{\alpha=1}^{m^2} z_{i\alpha}z_{j\alpha}$  by  $b_{ij}$  ( $i, j = 1, \dots, m$ ). Then

$$(20) \quad \begin{cases} b_{ij} = 0, & \text{for } i \neq j \\ \text{and } b_{ii} = N_i, \end{cases}$$

where  $N_i$  denotes the number of plots to which the variety  $v_i$  has been assigned. Because of (20) we have

$$(21) \quad \begin{vmatrix} b_{11} & \cdot & \cdot & \cdot & b_{1m} \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ b_{m1} & & & & b_{mm} \end{vmatrix} = N_1 N_2 \cdots N_m.$$

According to (13) we have

$$(22) \quad \begin{aligned} z'_{i\alpha} + \bar{z}_i w_\alpha &= z_{i\alpha}, & (i = 1, \dots, m-1) \\ -z'_{1\alpha} - \dots - z'_{m-1,\alpha} + w_\alpha(1 - \bar{z}_1 - \dots - \bar{z}_{m-1}) &= z_{m\alpha}. \end{aligned}$$

The determinant of these equations is given by

$$(23) \quad \lambda = \begin{vmatrix} 1 & 0 & 0 & \cdots & 0 & 0 & \bar{z}_1 \\ 0 & 1 & 0 & \cdots & 0 & 0 & \bar{z}_2 \\ \cdot & \cdot & \cdot & \cdots & \cdot & \cdot & \cdot \\ 0 & 0 & 0 & \cdots & 0 & 1 & \bar{z}_{m-1} \\ -1 & -1 & -1 & \cdots & -1 & -1 & \delta \end{vmatrix}$$

where  $\delta = 1 - \bar{z}_1 - \bar{z}_2 - \dots - \bar{z}_{m-1}$ . It is easy to verify that

$$(24) \quad \lambda = 1.$$

From (21), (22) and (24) it follows that

$$(25) \quad \Delta_2^* = N_1 N_2 \cdots N_m.$$

Hence, from (19) we obtain

$$(26) \quad \Delta_2 = N_1 N_2 \cdots N_m / m^2.$$

In the case of a Latin square design we have  $N_1 = N_2 = \dots = N_m = m$ . Hence

$$(27) \quad \bar{\Delta}_2 = m^{m-2}.$$

Because of the condition  $N_1 + N_2 + \dots + N_m = m^2$ , the right hand side of (26) becomes a maximum when  $N_1 = N_2 = \dots = N_m = m$ . Thus  $\Delta_2 \leq \bar{\Delta}_2$  and consequently the Latin square design is proved to be most efficient.

**4. Efficiency of Graeco-Latin and higher squares.** Consider  $m$  varieties  $v_1, \dots, v_m$  and  $m$  treatments  $q_1, \dots, q_m$ . Suppose that we wish to find out by experimentation whether the yield is affected by varieties or treatments. For this purpose the experimental area is subdivided into  $m^2$  plots lying in  $m$  rows and  $m$  columns and to each plot one of the varieties and one of the treatments is assigned. We call this arrangement a Graeco-Latin square if the following conditions are fulfilled: 1) each variety appears exactly once in each row and exactly once in each column; 2) each treatment appears exactly once in each row and exactly once in each column; 3) each variety is combined with each of the treatments exactly once.

The following general abstract scheme includes the Latin square and Graeco-Latin square as special cases: Consider an  $r$ -way classification with  $m$  classes in each classification. Denote by  $y_{a_1 a_2 \dots a_r}$  the value of a certain characteristic of an individual who is classified in the  $a_1$ -class of the first classification, in the  $a_2$ -class of the second classification,  $\dots$ , and in the  $a_r$ -class of the  $r$ -th classification. Suppose that  $m^2$  observations are made for the purpose of investigating the effect of the classes on the value of the characteristic under consideration. We will say that we have a generalized Latin square design if the following condition is fulfilled: *Let  $i, j, m'$  and  $m''$  be an arbitrary set of four positive integers for which  $i \neq j, i \leq r, j \leq r, m' \leq m$  and  $m'' \leq m$ . Then among the  $m^2$  individuals observed there exists exactly one individual who belongs to the  $m'$ -class of the  $i$ -th classification and  $m''$ -class of the  $j$ -th classification.*

It is clear that if  $r = 3$  the above scheme is a Latin square. If  $r = 4$  we have a Graeco-Latin square.

Assume that the observations  $y_{a_1 \dots a_r}(a_1, a_2, \dots, a_r = 1, \dots, m)$  are normally and independently distributed with a common variance  $\sigma^2$ . Assume furthermore that the expected value of  $y_{a_1 \dots a_r}$  is given by

$$E(y_{a_1 a_2 \dots a_r}) = \gamma_{1a_1} + \dots + \gamma_{ra_r}.$$

The parameters  $\sigma^2$  and  $\gamma_{ia}$  ( $i = 1, \dots, r; a = 1, \dots, m$ ) are unknown constants. Suppose that we wish to test the hypothesis that

$$(28) \quad \gamma_{i1} = \gamma_{i2} = \dots = \gamma_{im}.$$

It can be shown that if the number of observations is limited to  $m^2$ , we obtain a most efficient design by constructing a generalized Latin square. The proof of this statement is similar to that of the efficiency of the Latin square and is therefore omitted.