

Provided for non-commercial research and educational use.
Not for reproduction, distribution or commercial use.

Serdica

Bulgariacae mathematicae
publicationes

Сердика

Българско математическо
списание

The attached copy is furnished for non-commercial research and education use only.
Authors are permitted to post this version of the article to their personal websites or institutional repositories and to share with other researchers in the form of electronic reprints.

Other uses, including reproduction and distribution, or selling or licensing copies, or posting to third party websites are prohibited.

For further information on
Serdica Bulgaricae Mathematicae Publicationes
and its new series Serdica Mathematical Journal
visit the website of the journal <http://www.math.bas.bg/~serdica>
or contact: Editorial Office
Serdica Mathematical Journal
Institute of Mathematics and Informatics
Bulgarian Academy of Sciences
Telephone: (+359-2)9792818, FAX:(+359-2)971-36-49
e-mail: serdica@math.bas.bg

ESTIMATION OF MULTIPLE MISSING VALUES IN FIXED EFFECTS EXPERIMENTAL DESIGNS

SHAWNM ABDUL-KADIR

This paper deals with the estimation of multiple missing values for some unifactor experimental designs. The method of minimizing the residual sum of squares is used and an example is given.

Introduction. The statistical analysis of unifactor experimental designs is complicated very often by the lacking of some observations. When some values are missing, the usual method of computing the various sums of squares cannot be used unless the missing values are estimated through the existing data.

Nowadays a number of computer algorithms are known for estimating missing values in experimental designs. Some of them use iterative methods (Pearce & Jeffers [3], Preece [4]), others are non-iterative (see Draper [2], Rubin [5, 6]). Usually such algorithms are based primarily on subroutines used for complete data set analysis. In most of the methods the vector of estimated residuals is obtained as a part of the computer analysis.

This paper describes a procedure for estimating multiple missing values in Cross-Over design and Graeco-Latin square design by minimizing the residual sum of squares and shows that for a (3×3) Graeco-Latin square design the estimates of the missing values are always zero, i. e. $\{\theta_h = 0\}_{h=1}^t$. Also an example is given.

1. Cross-Over design. The model is as follows:

$$(1.1) \quad Y_{ijk} = \mu + T_i + R_j + C_k + \varepsilon_{ijk},$$

$$1 \leq i \leq t, \quad 1 \leq j \leq r, \quad 1 \leq k \leq c, \quad r=t, \quad r \leq c.$$

Here (Y_{ijk}) denotes the value observed in the i -th treatment, j -th row and k -th column; $Y_{ijk} = 0$ stands for the missing value in the cell (i, j, k) . Now let us estimate the missing values for a Cross-Over design by minimizing the residual sum of squares (MSSR).

The following notations shall be used.

Let $\{\theta_h\}_{h=1}^m$ be the unknown missing values, and let A_1, A_2 and A_3 be a treatment, a row and a column factor respectively. The (i, j, k) cell in the design matrix is the only one comprising 1 and the remaining values are zero.

Let us define the dummy variable corresponding to θ_h , as follows:

$$a_{hi} = 1 \text{ if } \theta_h \text{ is in treatment } i$$

$$0 \text{ otherwise,}$$

$$b_{hj} = 1 \text{ if } \theta_h \text{ is in row } j$$

$$0 \text{ otherwise,}$$

$$c_{hk} = 1 \text{ if } \theta_h \text{ is in column } k$$

$$0 \text{ otherwise.}$$

The following three dummy variables will be defined by the equation (1.2) if θ_g and θ_h , $h \neq g$ are missing.

$$(1.2) \quad \begin{aligned} \phi_{gh}(A_1) &= 1 \text{ if } \theta_g \text{ and } \theta_h \text{ are from one and the same treatment} \\ & \quad 0 \text{ otherwise,} \\ \phi_{gh}(A_2) &= 1 \text{ if } \theta_g \text{ and } \theta_h \text{ appear in the same row} \\ & \quad 0 \text{ otherwise,} \\ \phi_{gh}(A_3) &= 1 \text{ if } \theta_g \text{ and } \theta_h \text{ are in the same column} \\ & \quad 0 \text{ otherwise,} \end{aligned}$$

Also let

$$(1.3) \quad \begin{aligned} T_h(A_1) & \text{ be the total sum of all observations for the treatment where } \theta_h \\ & \text{ appears,} \\ T_h(A_2) & \text{ be the total sum of all observations for the row for which } \theta_h \text{ appears,} \\ T_h(A_3) & \text{ be the total sum of all observations for the column for which } \theta_h \\ & \text{ appears.} \end{aligned}$$

Theorem 1.1. For the Cross-Over design given by (1.1), the missing values $\{\theta_h\}_{h=1}^m$ are obtained by solving the following system of equations:

$$(1.4) \quad \begin{aligned} (c-1)(t-2)\theta_h + \sum_{g \neq h} \theta_g \{2-t[\sum \phi_{gh}(A_1) + \sum \phi_{gh}(A_2)] - c\sum \phi_{gh}(A_2)\} \\ = c\sum T_h(A_2) + t[\sum T_h(A_1) - \sum T_h(A_3)] - 2G \end{aligned}$$

for $1 \leq h \leq m$, when $(t-1)(r-2) > m$.

Proof. The marginal totals are, as follows:

$$(1.5) \quad \begin{aligned} T_i &= \sum_j Y_{i..} + \sum_{h=1}^m \theta_h a_{hi}, \quad i=1, \dots, t, \\ R_j &= \sum_i Y_{.j.} + \sum_{h=1}^m \theta_h b_{hj}, \quad j=1, \dots, r, \\ C_k &= \sum_i Y_{..k} + \sum_{h=1}^m \theta_h c_{hk}, \quad k=1, \dots, c, \\ G &= \sum_i \sum_j \sum_k Y_{ijk} + \sum_h \theta_h. \end{aligned}$$

Here $Y_{i..}$, for example, means averaging over the indices replaced by points while the rest index remains fixed.

The usual sum of squares for residual is

$$\text{SSE} = \text{SSG} - \text{SST} - \text{SSR} - \text{SSC},$$

where

$$\begin{aligned} \text{total sum of squares} &= \text{SSG} = \sum_i \sum_j \sum_k Y_{ijk}^2 - \frac{Y_{\dots}^2}{tc} \\ \text{sum of squares for treatment} &= \text{SST} = \sum_i \frac{Y_{i..}^2}{r} - \frac{Y_{\dots}^2}{tc} \\ \text{sum of squares for row} &= \text{SSR} = \sum_j \frac{Y_{.j.}^2}{c} - \frac{Y_{\dots}^2}{tc} \\ \text{sum of squares for column} &= \text{SSC} = \sum_k \frac{Y_{..k}^2}{t} - \frac{Y_{\dots}^2}{tc} \end{aligned}$$

The values of θ_h which minimize SSE are to be found from the equations:

$$\frac{\partial \text{SSE}}{\partial \theta_h} = 2\theta_h - \frac{2}{r} \sum_i T_i a_{hi} - \frac{2}{c} \sum_j R_j b_{hj} - \frac{2}{t} \sum_k C_k c_{hk} + \frac{4Y \dots}{tc}$$

for any $h, 1 \leq h \leq m$.

Replacing T_i, R_j, C_k and G by (1.5), we obtain the following equations:

$$(1.6) \quad (c-1)(t-2)\theta_h + \sum_{g \neq h} \theta_g [2-t(\sum a_{hi} a_{gi} + \sum c_{hk} c_{gk}) - c b_{hj} b_{gj}] \\ = t \sum b_{hj} \sum Y_{..k} + c (\sum a_{hi} \sum Y_{i..} + \sum c_{hk} \sum Y_{.j.}) - 2Y \dots$$

for fixed $h, 1 \leq h \leq m$.

The desired equation (1.4) is obtained by using relations (1.2) and (1.3) in (1.6), i. e.

$$(c-1)(t-2)\theta_h + \sum_{g \neq h} \theta_g \{2-t [\sum \emptyset_{gh}(A_1) + \sum \emptyset_{gh}(A_3)] - c \sum \emptyset_{gh}(A_2)\} \\ = c \sum T_h(A_2) + t [\sum T_h(A_1) + \sum T_h(A_3)] - 2G$$

for $(c-1)(t-2) > m$.

2. Graeco-Latin square design. The mathematical model for Graeco-Latin square design is given, as follows:

$$(2.1) \quad Y_{ijkl} = \mu + T_i + R_j + C_k + G_l + \varepsilon_{ijkl},$$

for $1 \leq i \leq t, 1 \leq j \leq r, 1 \leq k \leq c, 1 \leq l \leq s$, and $t=r=c=s, r > 3$.

Usually $Y_{ijkl} = 0$ for the missing value in cell (i, j, k, l) .

We define dummy variables for treatments, rows, columns and Greek-Letters, as follows:

$$a_{hi} = 1 \text{ if } \theta_h \text{ is in treatment } i \\ 0 \text{ otherwise,} \\ b_{hj} = 1 \text{ if } \theta_h \text{ is in row } j \\ 0 \text{ otherwise,} \\ c_{hk} = 1 \text{ if } \theta_h \text{ is in column } k \\ 0 \text{ otherwise,} \\ d_{hl} = 1 \text{ if } \theta_h \text{ is in Greek-Letter } l \\ 0 \text{ otherwise.}$$

The following four dummy variables will be defined by using equation (2.2) if θ_g and $\theta_h, h=g$ are missing.

$$(2.2) \quad \emptyset_{gh}(A_1) = 1 \text{ if } \theta_g \text{ and } \theta_h \text{ are from one and the same treatment} \\ 0 \text{ otherwise,} \\ \emptyset_{gh}(A_2) = 1 \text{ if } \theta_g \text{ and } \theta_h \text{ appear in the same row} \\ 0 \text{ otherwise,} \\ \emptyset_{gh}(A_3) = 1 \text{ if } \theta_g \text{ and } \theta_h \text{ are in the same column} \\ 0 \text{ otherwise,} \\ \emptyset_{gh}(A_4) = 1 \text{ if } \theta_g \text{ and } \theta_h \text{ are in the same Greek-Letter} \\ 0 \text{ otherwise,}$$

and A_1, A_2, A_3 and A_4 are treatment, row, column and Greek-Letter factors respectively.

Also let

$T_h(A_1)$ be the total sum of all observations for the treatment where θ_h appears,

- (2.3) $T_h(A_2)$ be the total sum of all observations for the row for which θ_h appears,
 $T_h(A_3)$ be the total sum of all observations for the column for which θ_h appears,
 $T_h(A_4)$ be the total sum of all observations for the Greek-Letter for which θ_h appears.

Theorem 2.1. For the Graeco-Latin square design the missing values $\{\theta_h\}_{h=1}^m$ are obtained by solving the following system of equations:

$$(2.4) \quad (t-1)(r-3)\theta_h + \sum_{g \neq h} \theta_g [3-r \sum_i \phi_{gh}(A_i)] = r \sum_i T_h(A_i) - 3G$$

for $1 \leq h \leq m$ when $(t-1)(r-3) > m$.

Proof. The marginal totals are found, as follows:

$$(2.5) \quad \begin{aligned} T_i &= \sum_j Y_{i..} + \sum_{h=1}^m \theta_h a_{hi}, \quad i=1, \dots, t \\ R_j &= \sum_i Y_{.j.} + \sum_{h=1}^m \theta_h b_{hj}, \quad j=1, \dots, r \\ C_k &= \sum_i Y_{..k.} + \sum_{h=1}^m \theta_h c_{hk}, \quad k=1, \dots, c \\ G_l &= \sum_i Y_{...l} + \sum_{h=1}^m \theta_h d_{hl}, \quad l=1, \dots, s \\ G &= \sum_i \sum_j \sum_k \sum_l Y_{ijkl} + \sum_{h=1}^m \theta_h. \end{aligned}$$

The usual sum of squares for residual is

$$SSE = SSG - SST - SSR - SSC - SSGL,$$

where

$$\begin{aligned} \text{total sum of squares} &= SSG = \sum_i \sum_j \sum_k \sum_l Y_{ijkl}^2 - \frac{Y^2 \dots}{tr} \\ \text{sum of squares for treatment} &= SST = \sum_i \frac{Y_{i..}^2}{r} - \frac{Y^2 \dots}{tr} \\ \text{sum of squares for row} &= SSR = \sum_i \frac{Y_{.j.}^2}{c} - \frac{Y^2 \dots}{tr} \\ \text{sum of squares for column} &= SSC = \sum_k \frac{Y_{..k.}^2}{t} - \frac{Y^2 \dots}{tr} \\ \text{sum of squares for Greek-Letter} &= SSCL = \sum_l \frac{Y_{...l}^2}{s} - \frac{Y^2 \dots}{tr} \end{aligned}$$

The values of θ_h which minimize SSE are to be found from the equations:

$$\frac{\partial SSE}{\partial \theta_h} = 0, \quad h=1, \dots, m,$$

i. e.

$$\frac{\partial \text{SSE}}{\partial \theta_h} = 2\theta_h - \frac{2}{t} \sum_i T_i a_{hi} - \frac{2}{r} \sum_j R_j b_{hj} - \frac{2}{c} \sum_k C_k c_{hk} - \frac{2}{s} \sum_l C_l c_{hl} + \frac{6G}{r^2} = 0.$$

Replacing T_i, R_j, C_k, G_l and G by (2.5), we obtain the following equations:

$$(2.6) \quad (t-1)(r-3)\theta_h + \sum_{g \neq h} \theta_g [3 - r(\sum a_{hi} a_{gi} + \sum b_{hj} b_{gj} + \sum c_{hk} c_{gk} + \sum d_{hl} d_{gl})] \\ = r[\sum a_{hi} \sum Y_{i...} + \sum b_{hj} \sum Y_{.j..} + \sum c_{hk} \sum Y_{...k.} + \sum d_{hl} \sum Y_{...l.}] - 3\sum Y_{...}$$

for any $h, 1 \leq h \leq m$.

The desired system of equations (2.4) is obtained by using relations (2.2) and (2.3) in (2.6), i. e.

$$(t-1)(r-3)\theta_h + \sum_{g \neq h} \theta_g [3 - r \sum_i \phi_{gh}(A_i)] = r \sum_i T_h(A_i) - 3G$$

for $1 \leq h \leq m$, when $(t-1)(r-3) > m$.

Corollary. For a (3×3) Graeco-Latin Square design the missing values are always zero, i. e. $\{\theta_h = 0\}_{h=1}^m$.

Proof. The general total is found by the equation

$$\sum_h \sum_i T_h(A_i) = (T'_i + R'_j + C'_k + G'_l) = G.$$

Also, any two missing values, θ_h and $\theta_g, h = g$, are in the same treatment or in the same row or in the same column or in the same Greek-Letter.

For $r=3$ by using (2.4), we obtain the identity

$$0 + 0 = 0,$$

and so any value of θ_h is acceptable, which leads to an uncertain situation.

Thus any missing value θ_h has to be set equal to zero.

3. Example. 1. Let us illustrate the MSSS method for estimating the missing values by giving the following example (see [1, p. 315]):

Let us have a (5×5) Graeco-Latin square for estimating the response of *Hylotrupes* larvae to four nutrients, each at five equally spaced log concentrations; y = average log weight after 103 ± 5 days on the experimental diet; rows are yeast extract, columns are poptone levels, treatments are cholesterol levels and Greek-Letters are riboflavin levels.

Let the treatment observation be denoted by D , and let the two observations $D\varepsilon$ and $D\beta$ in columns 2 and 5 and rows 3 and 5, respectively, be missing, where ε and β are Greek letters.

The missing observations are denoted by θ_1 and θ_2 with $m=2$. Then the dummy variable ϕ_{12} is determined, as follows:

$$\phi_{12}(A_1) = 1, \quad \phi_{12}(A_2) = 0, \quad \phi_{12}(A_3) = 0 \text{ and } \phi_{12}(A_4) = 0.$$

Since $r=t=5$, the corresponding totals are

$$T_1(A_1) = 1.45, \quad T_1(A_2) = 2.63, \quad T_1(A_3) = 1.95 \text{ and } T_1(A_4) = 2.25, \\ T_2(A_1) = 1.45, \quad T_2(A_2) = 2.67, \quad T_2(A_3) = 3.14 \text{ and } T_2(A_4) = 2.11, \\ G = 13.21.$$

The system of linear equations for calculating the missing values is obtained immediately from (2.4):

$$\begin{aligned}8\theta_1 - 2\theta_2 &= 1.77, \\ -2\theta_1 + 8\theta_2 &= 7.22\end{aligned}$$

which yields $\theta_1 = 0.477$, $\theta_2 = 1.022$.

Discussion. The method of minimizing the residual sum of squares has been applied to estimating multiple missing values in Cross-Over design and Graeco-Latin square design.

The solution is obtained directly. The technique may be applied to a number of designs. It is useful for the Randomized Block design where blocks are one of the factors. The formula holds also for a $(r \times r)$ Latin-Square design. For this model the system of m linear equations is the following:

$$(r-1)(r-2)\theta_h + \sum_{g=h} \theta_g [2-r \sum_i \varnothing_{gh}(A_i)] = r \sum_i T_h(A_i) - 2G$$

for $h = 1, \dots, m$ and $(r-1)(r-2) > m$.

Acknowledgements. I wish to thank my research supervisor Professor D. Vandev for our discussions. I am much obliged to Professor L. Boneva and Dr I. Tzankova for their recommendations. The referees' comments are most gratefully acknowledged.

REFERENCES

1. G. I. Bliss. *Statistics in Biology*, New York, 1968.
2. N. R. Draper. Missing values in response surface design. *Technometrics*, **3**, 1961, 389-398.
3. S. C. Pearce, J. N. R. Jeffers. Block designs and missing data. *J. R. Statist.*, **33**, 1971, 131-136.
4. D. A. Preece. Iterative procedures for missing values in experiments. *Technometrics*, **13**, 1971, 743-753.
5. D. B. Rubin. A non-iterative algorithm for least squares estimation of missing values in any analysis of variance design. *Appl. Statist.*, **21**, 1972, 136-141.
6. D. B. Rubin. Noniterative least squares estimates, standard error and F-test for analyses of variance with missing data. *J. R. Statist.*, **33**, 1976, 270-274.

Institute of Mathematics
P. O. Box 373
Sofia, Bulgaria

Received 29. 01. 1991
Revised 15. 06. 1991