

A PROBLEM IN LEAST SQUARES

BY JAN K. WIŚNIEWSKI

§1. We are dealing with two variables, the observed values of which are denoted x and y respectively. The pairs of observations are divided into r groups, numbering n_1, n_2, \dots, n_r pairs. Suppose in each group we determine a regression equation of the following shape:

$$y_i = a_i + b_i x + \dots + m_i x^s \tag{1}$$

where y_i denotes the value of the "dependent" variable obtained from the regression equation, while y without any subscript denotes its observed value. The r regression equations of type (1) are not assumed independent; on the contrary, we postulate that

$$\sum_1^r y_i = a_0 + b_0 x + \dots + m_0 x^s \tag{2}$$

be fulfilled identically in x ; a_0, b_0, \dots, m_0 being predetermined numbers. This leads to the following conditions:

$$\sum_1^r a_i = a_0 \quad \sum_1^r b_i = b_0 \quad \dots \quad \sum_1^r m_i = m_0. \tag{3}$$

The magnitude to be minimized under the theory of least squares is now

$$Z = \sum_1^{r-1} \sum_i [y - (a_i + b_i x + \dots + m_i x^s)]^2 + \sum_r \left\{ y - \left[\left(a_0 - \sum_1^{r-1} a_i \right) + \left(b_0 - \sum_1^{r-1} b_i \right) x + \dots + \left(m_0 - \sum_1^{r-1} m_i \right) x^s \right] \right\}^2. \tag{4}$$

The normal equations derived from (4) are of the following shape:

.....

$$n_j a_j + n_r \sum_1^{r-1} a_i + b_j \sum_i x + \left(\sum_1^{r-1} b_i \right) (\sum_r x) + \dots + m_j \sum_i x^s + \left(\sum_1^{r-1} m_i \right) (\sum_r x^s) = \sum_i y - \sum_r y + n_r a_0 + b_0 \sum_r x + \dots + m_0 \sum_r x^s \tag{5}$$

.....

$$\begin{aligned}
 & a_j \sum_i x + \left(\sum_1^{r-1} a_i \right) (\sum_r x) + b_j \sum_i x^2 + \left(\sum_1^{r-1} b_i \right) (\sum_r x^2) \\
 & + \cdots m_j \sum_i x^{s+1} + \left(\sum_1^{r-1} m_i \right) (\sum_r x^{s+1}) = \sum_i xy - \sum_r xy + a_0 \sum_r x \\
 & \qquad \qquad \qquad + b_0 \sum_r x^2 + \cdots m_0 \sum_r x^{s+1} \\
 & \dots\dots\dots \\
 & \dots\dots\dots \tag{5}
 \end{aligned}$$

$$\begin{aligned}
 & a_j \sum_i x^s + \left(\sum_1^r a_i \right) (\sum_r x^s) + b_j \sum_i x^{s+1} + \left(\sum_1^{r-1} b_i \right) (\sum_r x^{s+1}) \\
 & + \cdots m_j \sum_i x^{2s} + \left(\sum_1^{r-1} m_i \right) (\sum_r x^{2s}) = \sum_i x^s y - \sum_r x^s y \\
 & \qquad \qquad \qquad + a_0 \sum_r x^s + b_0 \sum_r x^{s+1} + \cdots m_0 \sum_r x^{2s} \\
 & \dots\dots\dots
 \end{aligned}$$

\sum_i meaning a summation extended over the i -th group. As (1) is of the s -th degree, we have $(s + 1)(r - 1)$ parameters to determine and as many equations, the problem thus being in theory solved.* As to the numerical solution, Doolittle's method or any other may be applied. We do not enter at present the question, how much labor would the actual solution require.

Examples. Allen and Bowley in their book on "Family Expenditure" (London, 1935) assume the expenditure on some defined item f to be a linear function of the total expenditure e

$$f = ke + c. \tag{6}$$

Evidently $\sum k = 1, \sum c = 0$ (cfr. pp. 10-11). Another example I give in a paper on seasonal variation, which appeared in "Economic Studies" III (Kraków). Actual values y of a time series are assumed to be linear functions of certain "normal" values x

$$y = a + bx \tag{7}$$

a and b changing from month to month but constant from year to year. Then $\sum a = 0, \sum b = 12$.

§2. Methods of solution in special cases. The generally recognized methods of solving normal equations become extremely laborious as the product $(s + 1)(r - 1)$ grows large. As a matter of fact, the amount of computer's work is approximately proportional to the cube of the number of parameters to determine. Therefore short cuts seem to be indispensable. A most elegant one is at our disposal in the special case¹ when the values of x in the several groups

* The remaining $s + 1$ parameters $a_r, b_r, \dots m$ are, of course, found from (3).

¹ This seems to be realized in Allen and Bowley's work.

§3. If this condition is not fulfilled, we can, indeed, replace the power series in x by orthogonal polynomials $X_{h,i}$, the second subscript being appended in order to show that the values of the X polynomials are no more identical for the several groups; these polynomials are now orthogonalized separately within each group. But we are no more able to predetermine the values of A_0, B_0, \dots, M_0 , as they depend on each other; this will be made clear a little later. Therefore we have to resort to an approximation: the values of the parameters will not be found from simultaneous equations, but successively, step by step, beginning with those corresponding to the highest degree of the independent variable.

The values of a_0, b_0, \dots, m_0 are given. It is evident that $m_0 = M_0$. The j -th normal equation is now:

$$M_j \sum_i X_{s,j}^2 - M_0 \sum_r X_{s,r}^2 + \left(\sum_1^{r-1} M_i \right) (\sum_r X_{s,r}^2) = \sum_j X_{s,j} y - \sum_r X_{s,r} y. \quad (12)$$

We see at once that

$$M_i = \frac{M_j \sum_i X_{s,j}^2 + \sum_i X_{s,i} y - \sum_j X_{s,j} y}{\sum_i X_{s,i}^2}. \quad (13)$$

Inserting this into /12/ we get

$$M_j = \frac{\sum_j X_{s,j} y}{\sum_j X_{s,j}^2} - \frac{1}{\sum_i X_{s,i}^2} \cdot \frac{\sum_1^r \frac{\sum_i X_{s,i} y}{\sum_i X_{s,i}^2} - M_0}{\sum_1^r \sum_i X_{s,i}^2}. \quad (14)$$

The second member of the right hand side of /14/ is again a correction term, the necessary amount of correction being distributed in inverse proportion to $\sum_i X_{s,i}^2$. Now we determine the value of L_0 , this coefficient corresponding to $s - 1$, the second highest degree of x , and calculate the several L 's from equations strictly analogous to (14) thus accomplishing the second step of our work, and so on, down to the A 's. L_0 is found from the following equation:

$$L_0 = l_0 - \sum_1^r [\alpha_{s-1}^s(i) \cdot M_i]. \quad (15)$$

To α_{s-1}^s is now appended a bracketed i , this to stress its variation from group to group. We see from (15) that before the several M 's are calculated we are not in a position to determine L_0 . On the other hand, if α_{s-1}^s is the same for all groups, the second member of the right hand side of (15) simply reduces to $\alpha_{s-1}^s \cdot m_0$ and L_0 can be determined in advance, i.e. before calculating the M 's. This is the case treated first (in §2). In any case, if no definite correlation is to be expected between $\alpha_{s-1}^s(i)$ and M_i , the approximative method developed here should give very nearly correct results. The writer applied this method of solution to the simple problem of seasonal variation mentioned in §1 and found the results very satisfactory.