

THE PROCESSING OF EXPERIMENTAL RESULTS BY THE METHOD OF LEAST SQUARES

S. P. Malysenko

Inzhenerno-Fizicheskii Zhurnal, Vol. 14, No. 2, pp. 309-313, 1968

UDC 512.897

It is demonstrated that the processing of experimental results by the method of least squares should be accomplished in various cases by minimizing the sum of the squares of the relative deviations in the estimates of the regression coefficients, rather than those of the absolute deviations, as is usually the case.

In processing experimental results we frequently use the method of least squares to smooth the experimental data, to determine the relationship between a certain measured quantity and another, etc. This method is a special case of the method of maximum probability,\* provided the observations are distributed in normal fashion [1].

Here, if it is only the dependent variables  $y_{i\lambda}$  that are subject to random measurement errors, and if  $y_{i\lambda}$  are not correlated, the problem of finding the equation of the curve which best represents the experimental data for  $y_{i\lambda}(x_i)$  (the regression curve) reduces to a minimization of the quadratic form

$$M = \sum_{i=1}^n \omega_i \left[ y_i - \sum_{k=1}^m a_k P_k(x_i) \right]^2 \quad (1)$$

with respect to the estimates of  $a_k$  for the coefficient of regression  $\alpha_k$ .

We assume that the theoretical relationship  $y = \eta(x)$  can be expanded into a system of linear independent functions  $P_k(x)$  and that it is well approximated by a sum of no more than  $m < n$  functions  $P_k(x)$  over the entire  $x$  interval:

$$\eta(x) = \sum_{k=1}^m \alpha_k P_k(x),$$

$M$  contains the scattering of the observations about the empirical curve  $\hat{\eta}(x) = \sum_{k=1}^m a_k P_k(x)$ , which is an estimate for the theoretical curve reduced to  $m$  terms.

Minimization of (1) with respect to  $a_k$  leads to a system of linear equations for  $a_k$

$$\sum_{k'=1}^m G_{kk'} a_{k'} = Y_k, \quad (2)$$

where

$$G_{kk'} = \sum_{i=1}^n P_k(x_i) \omega_i P_{k'}(x_i), \quad (3)$$

and

$$Y_k = \sum_{i=1}^n P_k(x_i) \omega_i y_i. \quad (4)$$

The solution for system (2) has the form

$$a_k = \sum_{k'=1}^m G_{kk'}^{-1} Y_{k'}, \quad (5)$$

where  $G^{-1}$  is the inverse matrix of the coefficients for the system satisfying the equation

$$\sum_{k'=1}^m G_{kk'} G_{k'l}^{-1} = \delta_{kl} \quad (\delta_{kk} = 1, \delta_{kl} = 0 \text{ when } k \neq l).$$

If the orthogonal polynomials  $K_k(x)$  are taken as  $P_k(x)$  so that

$$G_{kk'} = \sum_{i=1}^n K_k(x_i) \omega_i K_{k'}(x_i) = \delta_{kk'},$$

then

$$a_k = \sum_{i=1}^n K_k(x_i) \omega_i y_i. \quad (6)$$

For the estimates of the variance of the regression coefficients we can use the corresponding diagonal elements  $\sigma_{kk}^2$  of the matrix

$$\sigma_{kk'}^2 = s^2 G_{kk'}^{-1},$$

where

$$s^2 = \frac{M}{n-m}. \quad (7)$$

For the estimates of  $a_k$  to be effective, we have to choose  $m$  so that (7) is minimum [1]:

$$s^2(m) = \min.$$

The results (from the processing of the experimental data) and the possibility of interpreting these data depend significantly on the weights ascribed to the individual measurements.

The weight of the measurement can be evaluated on the basis of the observational data for the  $i$ -th group from the mean, and it is a function of  $x$ .

If  $h_i$  is a measure of the accuracy found in the  $i$ -th group of measurements, then  $\omega_i \sim h_i^2$  and  $h_i \sim 1/\Delta y_i$ , where  $\Delta y_i$  is the absolute measurement error, and the last ratio is satisfied if  $\Delta y_i$  is understood to be the probable, the mean, the mean square, or some other characteristic error [8].

The absolute error  $\Delta y_i$  can be expressed in terms of the relative  $\varphi_i$  as  $\Delta y_i = \varphi_i y_i$ , so that

$$\omega_i = \xi \frac{1}{\varphi_i^2 y_i^2}, \quad (8)$$

\* We will use the terminology adopted in [1].

where  $\xi$  is independent of  $x$  and can be omitted in the following.

With consideration of (8), we can write relationships (1)–(7) as follows:

$$M_1 = \sum_{i=1}^n \omega_{1i} \left[ 1 - \frac{1}{y_i} \sum_{k=1}^m a_{1k} P_k(x_i) \right]^2, \quad (1a)$$

$$\sum_{k'=1}^m G_{1kk'} a_{1k'} = Y_{1k}, \quad (2a)$$

$$G_{1kk'} = \sum_{i=1}^n \omega_{1i} \frac{1}{y_i^2} P_k(x_i) P_{k'}(x_i), \quad (3a)$$

$$Y_{1k} = \sum_{i=1}^n \omega_{1i} \frac{1}{y_i} P_k(x_i), \quad (4a)$$

$$a_{1k} = \sum_{k'=1}^m G_{1kk'}^{-1} Y_{1k'}, \quad (5a)$$

$$a_{1k} = \sum_{i=1}^n \omega_{1i} \frac{1}{y_i} K_k(x_i), \quad (6a)$$

$$\sigma_{1kk'}^2 = s_1^2 G_{1kk'}^{-1} \text{ and } s_1^2 = \frac{M_1}{n-m}, \quad (7a)$$

where

$$\omega_{1i} = \frac{1}{\Psi_i^2}.$$

In this notation the determination of the regression curve reduces to the minimization of the sum of the weighted squares of the relative deviations over the estimate of  $a_{1k}$ .

In processing an extensive amount of experimental material, we can save considerable work and time by assuming the measurements to be equal in accuracy. This is the usual approach in processing the results of thermophysical experiments: the derivation of the equations of state, the compilation of tables for the thermophysical properties of matter, the determination of virial coefficients, etc. [2–7]. In this case the processing is carried out on the basis of Eqs. (1)–(7) with  $\omega_{1i} = 1$ . This would be correct if the absolute measurement error remained approximately constant for each series of experiments over the entire interval.\* This is occasionally the case, but frequently an experiment—particularly, a thermophysical experiment—is formulated so that it is  $\varphi_1(x_i)$  rather than  $\Delta y_1$  that is kept approximately constant in the  $x$  interval. In this case, the processing of the experimental data should be carried out in accordance with Eqs. (1a)–(7a), assuming that  $\omega_{1i} = 1$ .

\*We know that a change in the weights of the individual observations by a factor of 2 to 3, as a rule, has virtually no effect on the magnitudes of the regression coefficients [8]. Therefore,  $\Delta y_1$  need not be rigorously constant in the  $x$  interval, but may vary slightly. Not only does such an approximation not contradict theory, but it is justified by its very nature.

This circumstance is not taken into consideration in the processing of a thermophysical experiment [2–7], and thus either leads to errors or to an unjustifiably great increase in the expenditure of labor and time, and frequently to both, since in the processing of data for which we can assume  $\varphi_1 \approx \text{const}$  on the basis of Eqs. (1)–(7) with  $\omega_{1i} = 1$ , unjustifiably small weights are ascribed to the low values of  $y_i$ , and the small  $y_i$  have virtually no effect on the magnitudes of the estimates for the regression coefficients. This becomes particularly apparent in processing data in which  $y$  varies markedly in the  $x$  interval ( $p$ ,  $\rho$ , and  $T$  are functions given in the interval  $\rho = 0 - \rho = x$ , and  $p_S$  and  $T_S$  are related).

Thus the authors of [4] write: "In order to obtain approximately identical errors for the approximation functions, the values of the argument must be specified nonuniformly, increasing the number of points in the region of the small values of the functions. In this case it is impossible to prepare tables of orthogonal polynomials in advance..." This difficulty could be avoided by carrying out the procedure according to (1a)–(7a). It is appropriate at this point to note that in the approximation of any tables by means of equations it is generally necessary for the approximation function to retain a certain number of true signs over the entire interval of the argument. If this approximation is carried out by the method of least squares, the calculation has to be carried out on the basis of (1a)–(7a) with  $\omega_{1i} = 1$ .

In the work of Michels and his co-workers—devoted to the calculation of the virial coefficients of  $\text{CH}_2$ ,  $\text{H}_2$ , and  $\text{D}_2$  on the basis of  $p\rho T$  data—the calculation is carried out by the method of least squares on the basis of Eqs. (1)–(7) with  $c\omega_{1i} = 1$  [2, 3], whereas the calculation should have been carried out on the basis of (1a)–(7a) with  $\omega_{1i} = 1$ . Since in the case of low density the points remained virtually without consideration in the calculation [2, 3], to obtain stable values for the estimates of the virial coefficients it became necessary to use data for higher densities and to employ polynomials of higher degrees than was required for this purpose. In this connection, the authors of [2, 3] had to carry out the calculation on an electronic digital computer, doubling the number of significant places. All of this resulted in an unjustified and substantial increase in the expenditure of work and time, including the amount of machine time. It is clear that the authors of [2, 3] also failed to evaluate the variance of the coefficients reliably.

Moreover, in connection with the fact that number of experimental points on each isotherm is limited, in the case of low  $T$  the condition  $n \gg m$  with respect to  $\text{H}_2$  and  $\text{D}_2$  [3] was disrupted, making the method of least squares inapplicable to the processing of these experimental results. As a consequence of this, at temperatures  $T \leq 138.16^\circ \text{K}$ , even the third virial coefficient for  $\text{H}_2$  was not determined reliably. On the basis of the data in [3], for  $T \leq 138.16^\circ \text{K}$ , it is independent of  $T$ . However, with proper processing of these varied data on the basis of (1a)–(7a) this phenomenon is not observed. The deliberate violation of

the condition  $n \gg m$  in combination with the correct processing—which is what we did—for the 98.16° K isotherm led to values of  $C(T)$  close to those obtained in [3].

It should be noted that the correct evaluation of the weight is particularly significant if the estimates of the regression coefficients are derived after processing of a theoretical interpretation or if it is the intention to determine the derivatives of the measured quantity with respect to  $x$  from the regression curve or from the parameters.

Since the final verification of the equations of state derived in [4–7] was carried out by comparison against experimental data pertaining to the relative deviations, the errors characteristic of processing procedures using (1)–(7) with  $\omega_i = 1$  might apparently have become evident only in the determination of the derivatives of the thermal quantities and in the interpretation of the regression coefficients.

#### NOTATION

$y_{i\lambda}$  is the experimental value of the  $y$  function at point  $x_i$ ;  $y_i$  is the empirical mean of  $i$ -th group of observation;  $x_i$  is the independent variable;  $\omega_i$  is the weight of  $y_i$ ;  $n$  is the number of measurements at various  $x$ ;  $T$  is the temperature;  $\rho$  is the density;  $p$  is the pressure;  $\Delta y_i$  is the absolute error in  $y_i$ ;  $\varphi_i$  is the relative error in  $y_i$ .

#### REFERENCES

1. N. P. Klepikov and S. N. Sokolov, Analysis and Planning of Experiment by the Method of Maximum Probability [in Russian], Nauka, 1964.
2. A. Michels, J. C. Abels, C. A. Ten Seldam, and W. De Graaff, *Physica*, **26**, 381, 1960.
3. A. Michels, W. De Graaff, and C. A. Ten Seldam, *Physica*, **26**, 393, 1960.
4. A. A. Vasserman, Ya. Z. Kazavchinskii, and V. A. Rabinovich, Thermophysical Properties of Air and Its Components [in Russian], Nauka, 1966.
5. M. P. Vukalovich and V. V. Altunin, Thermophysical Properties of Carbon Dioxide [in Russian], Atomizdat, 1965.
6. A. Stein candidate's dissertation, "Methoden zum Aufstellen von Zustandsgleichungen für reinefluide Stoffe," Braunschweig, 1965.
7. M. P. Vukalovich, Thermodynamic Properties of Water and Water Vapor [in Russian], Mashgiz, Moscow, 1958.
8. A. Worthing and J. Heffner, Methods for the Processing of Experimental Data [Russian translation], IL, 1949.

31 March 1967

High Temperature Institute  
AS USSR, Moscow