

## APPARENT LINEAR RELATIONSHIP, COMPENSATION LAW AND OTHERS. PART I

J. NORWISZ and Z. SMIESZEK

*Institute for Non-Ferrous Metals, Gliwice (Poland)*

Z. KOLENDA

*University of Mining and Metallurgy, Cracow (Poland)*

(Received 1 May 1989)

### ABSTRACT

The graphical representation of the error in a one-dimensional variable is a sector, similar to an ellipse for a two-dimensional variable. If the slope  $a$  and the constant  $b$  of the straight line  $y = ax + b$  are estimated by the least-squares method, the geometrical properties of the error ellipse are determined by values of the independent variables, e.g. the tangent to the main axis is equal to minus the average value of the independent variable. If this same experiment is repeated, the estimated values of the straight line coefficients will fall inside the error ellipse, giving points lying on a straight line in a system with  $a$  and  $b$  axes. This phenomenon is more obvious when the experiments have a bigger experimental error.

A similar effect will occur when different experiments are really the same, e.g. when the experimental error is larger than the differences between experiments.

### INTRODUCTION

The estimation of linear equation coefficients is a procedure most frequently used for elaboration of the experimental results. The least-squares method used here is widely known and accepted [1]. According to general opinion, being “strictly scientific” it prevents mistakes. However, it is intriguing how frequently data indicating the relationship between the linear equation coefficients can be found in the literature.

For example, a linear relationship has been observed between the natural logarithm of the reaction constant and the activation energy, for many groups of experimental results. This relationship is presented either as the so-called linear compensation law [2] or is hidden in the tabulated data.

There are several publications devoted to justification [3,4] of the observed relationship; there are also several publications [5,6] denying its existence and attempting to prove that the relationship results from inaccuracy in the measurements.

The aim of this work is a re-analysis of the problem of the estimation of linear equation coefficients using the least-squares method, and to demonstrate the existence of an apparent relationship between these coefficients.

As there is a lack of understanding of the meaning of the error of the estimated values, there is a subsequent overestimation of the accuracy of the measurements. The linear compensation law, and others which are similar, most frequently result from this.

Theoretical considerations have been the basis for a re-calculation of sets of experimental results selected from recent publications, in order to determine whether a relationship exists between the linear equation coefficients, and whether it results from the estimation error or is a reflection of actual differences between the experiments.

#### THE COVARIANCE MATRIX OF THE LINEAR EQUATION COEFFICIENTS

According to the method of least squares, the best values of the straight line coefficients

$$y = ax + b \quad (1)$$

describing a given set of experimental results, are calculated from the equations

$$\tilde{a} = \frac{\sum yx + \sum y \sum x/n}{\sum x^2 + (\sum x)^2/n} \quad (2)$$

and

$$\tilde{b} = \sum y/n - \tilde{a} \sum x/n \quad (3)$$

where

$$\sum x = \sum_{i=1}^{i=n} x_i, \text{ etc.} \quad (4)$$

$a$  and  $b$  are the coefficients of the linear equation and  $n$  is the number of experiments

It has been assumed that the values measured,  $y_1, y_2, \dots, y_n$ , are statistically independent and have an error of normal distribution. The values  $x_1, x_2, \dots, x_n$  are assumed to be correctly determined.

The values of the linear equation coefficients, calculated using eqns. (2) and (3), are the estimate of the two-dimensional normal random variable with the covariance matrix

$$\mathbf{M}(a, b) = \begin{bmatrix} s^2(a) & \text{cov}(a, b) \\ \text{cov}(a, b) & s^2(b) \end{bmatrix} \quad (5)$$

The  $\mathbf{M}(a, b)$  matrix elements are calculated according to the so-called error propagation principle

$$\mathbf{M}(a, b) = \mathbf{P}^T \cdot \mathbf{M} \cdot \mathbf{P} \quad (6)$$

where

$$\mathbf{P}^T = \begin{bmatrix} \partial a / \partial y_1 & \partial a / \partial y_2 & \dots & \partial a / \partial y_n \\ \partial b / \partial y_1 & \partial b / \partial y_2 & \dots & \partial b / \partial y_n \end{bmatrix} \quad (7)$$

$$\mathbf{M} = \text{diag}[\mu, \mu, \dots, \mu] \quad (8)$$

where  $\mu$  is the standard deviation.

By appropriate differentiation of eqns. (2) and (3), using eqns. (6) and (7), the following relations are obtained

$$s^2(a) = \mu / \left[ \sum x^2 - (\sum x)^2 / n \right] \quad (9)$$

$$s^2(b) = \mu \left\{ 1/n + \bar{x}^2 / \left[ \sum x^2 - (\sum x)^2 / n \right] \right\} \quad (10)$$

$$\text{cov}(ab) = \mu \left\{ -\bar{x} / \left[ \sum x^2 - (\sum x)^2 / n \right] \right\} \quad (11)$$

The best  $\mu$  value is given by the expression

$$\mu \approx s^2(y_0) = \sum_{i=1}^n (y_i - \tilde{a}x_i - \tilde{b})^2 / (n - 2) \quad (12)$$

where  $\bar{x}$  denotes the mean value of  $x_i$ .

The values of the  $\mathbf{M}(a, b)$  covariance matrix elements depend (with  $\mu$  factor accuracy) only on the value of the so-called independent experimental variables  $x_1, x_2, \dots, x_n$ .

## THE RANDOM VARIABLE DISTRIBUTION

In the case of the two-dimensional random variable  $\mathbf{l}$ , its density distribution is given by

$$p(\tilde{a}, \tilde{b}) = \frac{\{\det[\mathbf{M}(a, b)]^{-1}\}^{\frac{1}{2}}}{2\pi} \exp\left[-\frac{1}{2}(\mathbf{m} - \mathbf{l})^T \mathbf{M}^{-1}(a, b)(\mathbf{m} - \mathbf{l})\right] \quad (13a)$$

where

$$\mathbf{m}^T = [m(a), m(b)], \quad (13b)$$

and

$$\mathbf{l}^T = [\tilde{a}, \tilde{b}]. \quad (13c)$$

$m(a)$  and  $m(b)$  are the expected values of the random variables  $a$  and  $b$ .

There is a certain probability  $\alpha$  that the given estimate of the  $\tilde{a}$ ,  $\tilde{b}$  pair of values will be found within the ellipse

$$p(\tilde{a}, \tilde{b}) = P(\alpha) \quad (14)$$

Its equation is determined by the relationship

$$R^2(\alpha) = (\mathbf{m} - 1)^T \mathbf{M}^{-1}(a, b)(\mathbf{m} - 1) \quad (15)$$

where  $R^2(\alpha)$  has the chi-square distribution of  $n$  degrees of freedom. The covariance ellipse determines the area of the estimated error of the  $\tilde{a}$ ,  $\tilde{b}$  pair, in the same way as the confidence interval in the case of a one-dimensional random variable.

### GEOMETRICAL PROPERTIES OF THE COVARIANCE ELLIPSE

The covariance ellipse matrix and the direction of the expected values determine explicitly the geometric properties of the covariance ellipse. Utilizing simple analytical geometric properties [7], the following may be stated.

(1) Ellipses calculated for the different  $R^2(\alpha)$  values are alike and have the same centre of gravity (Fig. 1)

$$|m_1 0| / |m'_1 0| = |m_2 0| / |m'_2 0| = R(\alpha) / R'(\alpha) \quad (16)$$

(2) The tangent  $\theta$  of the ellipse axes is determined by the equations

$$\tan(2\theta) = -2\text{cov}(a, b) / [s^2(b) - s^2(a)] = T \quad (17)$$

$$\tan(\theta_{1,2}) = -1/T \pm (1/T^2 + 1)^{1/2} \quad (18)$$

(3) The ellipse axes are perpendicular to each other

$$\tan(\theta_1) \tan(\theta_2) = 1 \quad (19)$$

(4) The diameters of the axes in the  $\theta_1$  and  $\theta_2$  directions are given by

$$S_{1,2} = 2R(\alpha) \left\{ \frac{[\text{cov}^2(a, b) - s^2(a)s^2(b)][1 + \tan^2(\theta_{1,2})]}{2\text{cov}(a, b) \tan(\theta_{1,2}) - s(a)^2 \tan^2(\theta_{1,2}) - s^2(b)} \right\}^{1/2} \quad (20)$$

(5) The ellipse is inscribed within a rectangle

$$B_a = 2R(\alpha)s(a), \quad B_b = 2R(\alpha)s(b) \quad (21)$$

(6) The ellipse area is proportional to  $R^2(\alpha)$  and to the square root of the covariance matrix determinant

$$\begin{aligned} F &= \pi R^2(\alpha) [s^2(a)s^2(b) - \text{cov}^2(a, b)]^{1/2} \\ &= \pi R^2(\alpha) \{ \det[\mathbf{M}(a, b)] \}^{1/2} \end{aligned} \quad (22)$$

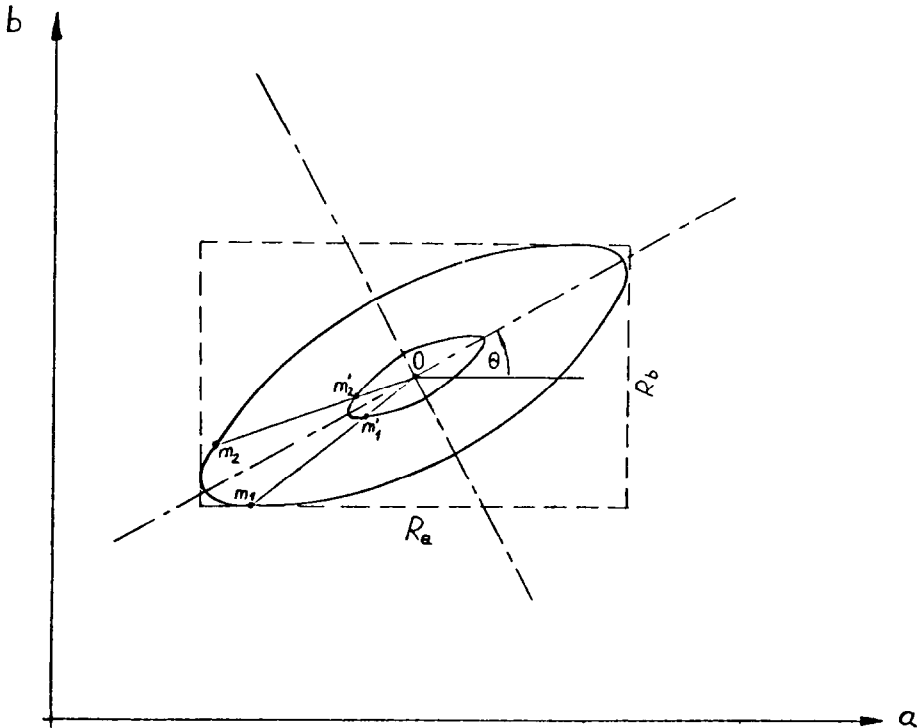


Fig. 1. Ellipses calculated for different  $R^2(\alpha)$  values.

(7) Similar ellipses, of the same covariance matrix, may be described by the same correlation coefficient

$$r(a, b) = \text{cov}(a, b) / [s(a)s(b)] \quad (23)$$

$$\mathbf{M}(a, b) = [s(a)s(b)] \begin{bmatrix} r(a, b) & 1 \\ 1 & r(a, b) \end{bmatrix} \begin{bmatrix} s(a) \\ s(b) \end{bmatrix} \quad (24)$$

#### THE APPARENT LINEAR RELATIONSHIP BETWEEN THE LINEAR EQUATION COEFFICIENTS

Because the covariance matrix of the linear equation coefficients is determined with an accuracy of the  $\mu$  factor, when the independent variable values are known, e.g. from experimental results, many of the ellipse properties mentioned above are known before the experiment and do not depend on the dependent variables  $y_1, y_2, \dots, y_n$ .

Only the scale factor, the  $\mu$  value and subsequently the properties (5) and (6) remain unknown. The most relevant properties, such as the ellipse axes

slope angle and the diameter ratio are known. These values may be calculated from eqns. (18) and (22), using eqn. (9)  $\div$  eqn. (11).

Repeating the experiment results in a subsequent pair of estimation values  $\tilde{a}$ ,  $\tilde{b}$  which should occur within the given covariance ellipse along a large ellipse axis. This phenomenon will be more obvious with greater experimental error, i.e. large  $\mu$  values, as in this case the relatively narrow covariance ellipse will be inscribed within a relatively large rectangle and filled with points, characterizing the results of successive estimations. Where there are no lines fixing the rectangle sides, the image of the points will lie along the line

$$\tilde{b} = \alpha \tilde{a} + \beta \quad (25)$$

The dimension of  $\beta$  is the same as the dimension of the dependent variable  $y$ ; the dimension of  $\alpha$  equals the reciprocal dimension of the independent  $x$ .

The phenomenon of ordering of the results of successive estimations along such "compensation" lines will also occur when successive experiments are carried out with slightly different experimental conditions, i.e. using slightly different values of the independent variables  $x_1, x_2, \dots, x_n$ . The slope angle of the covariance ellipse is approximately equal to  $-\bar{x}$ , which results from eqn. (18) by substituting from eqn. (10) and neglecting the  $1/n$  term in eqn. (10). Thus, it is sufficient that the mean arithmetic values of the independent variable  $x$  are close in repeated experiments. This phenomenon may lead to logical mistakes.

If during the course of successive experiments a factor insignificant for the experiment, e.g. moon phases, changes, the occurrence of relationship (25) may be attributed to the as yet unrevealed connections between these facts, i.e. between the experiment and the state of the celestial body.

Similarly, when the actual differences between two different experiments are small in comparison with the true experimental error, the linear relationship (25) will be observed, though now there are "rational" circumstances to explain this in the form of a simple linear equation. If in addition some ordering of the individual estimation results is observed, then eqn. (25) may subsequently be used to estimate the results of other experiments belonging to the same group. Disclosure of a strong statistical dependency, i.e. a high absolute value of the correlation coefficient  $r(a, b)$  for eqn. (25), only confirms these tendencies.

Nevertheless, the value of this coefficient is determined by the experimental conditions and can be calculated before the experiment from eqn. (23), using eqn. (9)  $\div$  eqn. (11). The correlation coefficient  $r(a, b)$  values are usually large.

The existence of the linear relationship (25) between  $a$  and  $b$  (see eqn. (1)) results from the application of the least-squares method to describe experimental results characterized by a normal distribution and is not a property of the objects tested.

As other methods for estimation of the linear equation (1) coefficients give values similar to those obtained by the least-squares method, then a similar relationship should be expected. The differences between the individual estimates do not reflect the actual differences between the given experiments, but only the experimental error. Therefore, no physical significance should be attributed to a linear equation of the type of eqn. (25). The covariance ellipse inclined at  $\arctan(-\bar{x})$  only determines the error in the estimation of a given set of  $[a, b]$  parameters.

The suggestion that the estimation error image of the straight line coefficients is an ellipse with axes perpendicular to the coordinates  $a, b$  axes of diameter proportional to the standard deviation values  $s(a)$  and  $s(b)$ , or may be a rectangle of sides proportional to these standard deviation values results from an incorrect generalization of the ideas connected with the application of the one-dimensional normal random variable in a multi-dimensional case.

However, it cannot be unequivocally decided whether a relationship of the eqn. (25) type is merely an image of an experimental error or a reflection of the true circumstances.

The value of the direction coefficient calculated from eqn. (18) can, however, be compared with the value calculated for the given set of estimates  $\tilde{a}, \tilde{b}$ .

On the one hand, it cannot be denied that the relationship mentioned does exist and that the value of the direction coefficient of the true relationship is equal or close to the value calculated for the given covariance matrix. On the other hand, a distinct difference between the calculated value for the given covariance matrix and the experimental value may result from both the statistical spread and a large experimental error showing no random error features. It should be added that an estimation of the experimental error using the standard deviation for a single measurement merely defines the spread of experimental values around the mean, does not give the true value, and is usually rather optimistic as regards the quality of the measurements.

## CONCLUSIONS

The considerations presented indicate the necessity of exercising great caution when a relationship between the coefficients of a linear equation of type (25) has been given. This caution should be greater still when the direction coefficient value of eqn. (25) is close to the calculated value based on the data from the covariance matrix. In such a case the experiment should be repeated preferably under different experimental conditions, or the result obtained must be compared with other literature data. Only then can it be stated whether or not the relationship observed is of any significance.

## REFERENCES

- 1 W. Volk, *Applied Statistics for Engineers*, McGraw-Hill, New York, 1969.
- 2 J.M. Thomas and W.J. Thomas, *Introduction to the Principles of Heterogenous Catalysis*, Academic Press, New York, 1967.
- 3 A.I. Lesnikovich and S.V. Levchik, *J. Therm. Anal.*, 30 (1985) 237.
- 4 E. Koch, *Non-Isothermal Reaction Analysis*, Academic Press, London, 1977, p. 58.
- 5 O. Exner, *Coll. Czech. Chem. Commun.*, 29 (1964) 1094.
- 6 J. Norwicz and J. Plewa, *J. Therm. Anal.*, 17 (1979) 549.
- 7 A.C. Jones, *An Introduction to Algebraical Geometry*, Oxford University Press, London, 1937.