

# Coefficient of determination in some atypical situations: use in chemical correlation analysis

Otto Exner<sup>1\*</sup> and Karel Zvára<sup>2</sup>

<sup>1</sup>Institute of Organic Chemistry and Biochemistry, Academy of Sciences of the Czech Republic, 16610 Praha 6, Czech Republic

<sup>2</sup>Department of Probability and Mathematical Statistics, Faculty of Mathematics and Physics, Charles University, Sokolovská 83, 186 00 Praha 8, Czech Republic

Received 1 April 1998; revised 4 June 1998; accepted 24 June 1998

**ABSTRACT:** The popular statistics used in correlation analysis, the correlation coefficient  $R$  and coefficient of determination  $R^2$ , belong to the model of common linear regression with an exact explanatory variable and freely fitted intercept. When they are used in other models, the common equation must be modified and the interpretation revised. The resulting equations are discussed for the model with a fixed intercept. For the models with errors in either variable and with both modifications, new equations are suggested. Under these conditions,  $R$  may be defined differently, always in at least two ways; the different interpretations are discussed. Some properties of these definitions are shown, particularly the sensitivity to coordinate transformations. These properties and interpretations differ from those of  $R$  in the common linear regression, hence the values in different models must not be directly compared. Applications in chemistry and some cases of misuse in the literature are mentioned. Copyright © 1999 John Wiley & Sons, Ltd.

**KEYWORDS:** correlation coefficient; linear regression; no-intercept regression; functional model

## INTRODUCTION

In the correlation analysis of chemical data, simple or multiple regression is the most common mathematical tool.<sup>1</sup> The adherence of the model to the data is often expressed by the correlation coefficient,  $R$ . In connection with regression, the more exact term is the coefficient of determination,  $R^2$ . Two notes have appeared<sup>2</sup> suggesting the suitable statistics which should be published in chemical papers to evaluate the fit and to compare various regressions; the coefficient of determination was always included. Jaffé's classical review<sup>3</sup> gives, in addition to  $R$ , the standard deviation from the regression line,  $SD$ , and the number of data points,  $N$ . These statistics can be considered as standard and are also in the output of all programs. Some papers are excessively brief when they give just  $R$  as the only characteristic (e.g. Ref. 4), and sometimes even  $N$  is not explicitly given.<sup>5</sup> On the other hand, some cautionary notes<sup>6</sup> have drawn attention to the asymmetric distribution of  $R$  and its strong dependence on the number of degrees of freedom. Great

caution is necessary when interpreting  $R$  derived with only a few degrees of freedom, and in this case even the adjusted coefficient of determination,<sup>7</sup>  $R_a^2$ , will not help. Several times, the use of  $R$  was completely rejected in connection with regression,<sup>8</sup> for not convincing reasons. In fact,  $R$  is closely related to one specific model of linear regression whose conditions are not fulfilled in all applications.<sup>9</sup> Alternative models have also been used in chemistry but are not common. Regression with a fixed origin (the no-intercept model) has been advocated,<sup>8b,10</sup> mainly in cases when this origin has a particular physical meaning, in correlation analysis, for instance, with hydrogen as substituent. In other cases, regressions are of importance with experimental errors on both coordinates.<sup>11</sup> Sometimes,  $R$  was calculated even in these cases, although already its definition may be doubtful. Summarizing, applications of  $R$  in chemistry are not always satisfactory and clear instructions for its use in various cases have not been presented.

In this paper, we deal with the extension of  $R^2$  to these non-typical regressions. There are questions of what modifications of the mathematical formula are needed and to what extent the interpretation can be maintained that is common in the normal regression. The problem has been tackled by Kvålseth<sup>9</sup> in a more general way. He collected several general expressions for  $R^2$ , postulated their desirable properties and applied them also to the no-intercept model (denoted B below). In this paper, these

\*Correspondence to: O. Exner, Institute of Organic Chemistry and Biochemistry, Academy of Sciences of the Czech Republic, 16610 Praha 6, Czech Republic.

Contract/grant sponsor: Grant Agency of the Czech Republic; Contract/grant number: 203/96/1658; Contract/grant number: Copernicus IRP-10053.

**Table 1.** Some properties of various regression models

Parameter	Model			
	A	B	C	D
Equation of dependence	$y = \alpha + \beta x$	$y = \beta x$	$y = \alpha + \beta x$	$y = \beta x$
Criterion to be minimized	$\Sigma (y - a - bx)^2$	$\Sigma (y - bx)^2$	$\Sigma (y - a - bx)^2 / (1 + b^2)$	$\Sigma (y - bx)^2 / (1 + b^2)$
Estimate of $\beta$	$b_3$ , Eqn 3	$b_5$ , Eqn 5	$b_9$ , Eqn 9	$b_{12}$ , Eqn 12
Definition based on $R_1^2$	$R_4^2$ , Eqn 4	$R_6^2$ , Eqn 6	$R_{10}^2$ , Eqn 10	$R_{13}^2$ , Eqn 13
Invariance to multiplication by a positive constant	Yes	Yes	Yes if both the same	Yes if both the same
Invariance to shift	Yes	No	Yes	No
Invariance to rotation	No	No	Yes	Yes
Definition based on $R_2^2$	$R_4^2$ , Eqn 4	$R_7^2$ , Eqn 7	$R_4^2$ , Eqn 4	$R_7^2$ , Eqn 7
Invariance to multiplication by a positive constant	Yes	Yes	Yes	Yes
Invariance to shift	Yes	No	Yes	No
Invariance to rotation	No	No	No	No

considerations will be extended also to the functional models (denoted C and D) with particular attention to the applications in the chemical correlation analysis: several misunderstandings in the literature will be pointed out.

## DEFINITIONS OF R

Of the various possible definitions<sup>7,9</sup> we have chosen the following two:

$$R_1^2 = 1 - \frac{RSS}{\sum (y - \bar{y})^2} \quad (1)$$

$$R_2^2 = \frac{(\sum (y - \bar{y})(\hat{y} - \bar{\hat{y}}))^2}{\sum (y - \bar{y})^2 \sum (\hat{y} - \bar{\hat{y}})^2} \quad (2)$$

They are general and can be applied in various types of regressions and even outside regression, i.e. always when real values of  $y$  and those predicted by any theory,  $\hat{y}$ , are to be compared.  $R_1^2$  in Eqn 1 may be understood as the amount of information gained from the regression.<sup>9</sup> The numerator  $RSS$  is the residual sum of squares, i.e. of squared differences between the values of  $y$  and the predicted values  $\hat{y}$ . It may be denoted also as the information entropy after the regression. The denominator is the entropy before any regression has been attempted, i.e. when only the mean value  $\bar{y}$  and the standard deviation from this mean were known.  $R_2$  in Eqn 2 can be viewed as cosine of the angle of two vectors in the  $N$ -dimensional sample space: the vector of real values  $y$  and the vector of predicted values  $\hat{y}$ ;  $\bar{y}$  and  $\bar{\hat{y}}$  are the mean values of  $y$  and  $\hat{y}$ , respectively. These two equations are considered here as basic definitions and will be applied in regressions of four types which differ by the explanatory variable (exact or loaded with error) and by

the intercept (freely fitted or fixed). Kvålseth<sup>9</sup> gives some other general equations but some of them differ very little. We found the two definitions in Eqns 1 and 2 to be sufficient for differentiating and characterizing our models.

## R IN VARIOUS REGRESSION MODELS

### (A) Common linear regression

This common type of linear regression with an exact explanatory variable and freely fitted intercept is used in almost all chemical applications. The regression line is a common straight line and the least-squares condition concerns the sum of squares of vertical deviations from this line (see Table 1 where the equations are summarized). The slope of the regression line is estimated as  $b_3$  (we use always the same subscript as is the number of the equation):

$$b_3 = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2} \quad (3)$$

The two definitions of  $R^2$  in Eqns 1 and 2 yield in this case the same popular equation known from textbooks:

$$R_4^2 = \frac{[\sum (x - \bar{x})(y - \bar{y})]^2}{\sum (x - \bar{x})^2 \sum (y - \bar{y})^2} \quad (4)$$

Some properties of this characteristic are given in Table 1. Note that  $R_4^2$  according to Eqn 4 is invariant to any linear transformation of either variable, i.e. to both shift of the intercept and change of scaling. However, it is affected by rotation of coordinates. This is all well

known. The only problem in this type of regression may be whether it has been applied correctly in a given case. Particularly the following applications may be accepted as correct: (a) the variable  $x$  is (approximately) free of error, the variable  $y$  is loaded with an experimental error and the deviations from the regression are caused only by these errors; (b) both variables are (approximately) exact but  $y$  contains an unknown component which does not depend on  $x$  and possesses a regular distribution; the deviations are caused by the latter factor; (c) both variables are exact and the deviations are caused by their inexact dependence; there are physical reasons why  $y$  can be predicted from  $x$  and not vice versa (for instance, the reaction rate can depend on temperature, not conversely). Many applications in chemistry meet these presumptions only roughly, or not at all.

### (B) No-intercept model

This model assumes an exact explanatory variable and fixed intercept. The intercept may lie at the point (0,0). The regression line goes through this point and the least-squares criterion is the same as in the preceding model (see Table 1). The estimate  $b_5$  of the unknown slope  $\beta$  is given by

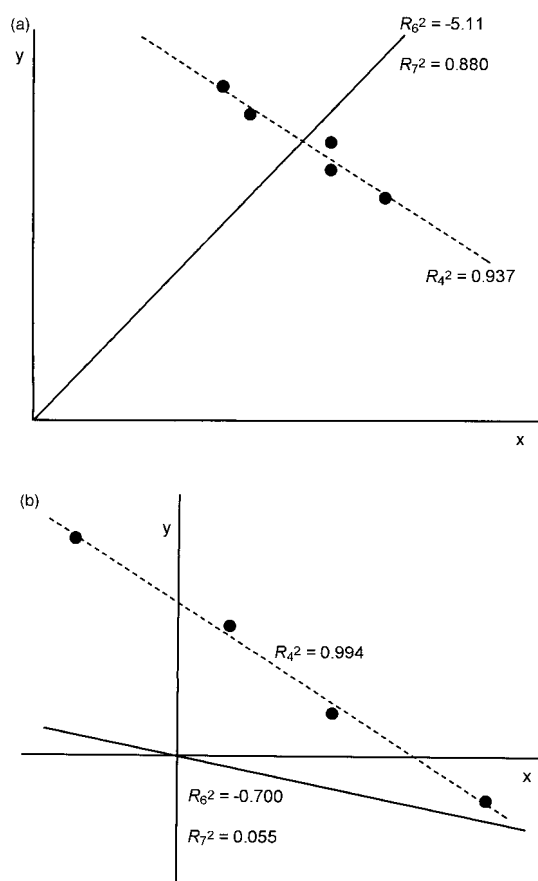
$$b_5 = \frac{\sum xy}{\sum x^2} \quad (5)$$

The two definitions of  $R^2$  in Eqns 1 and 2 in this case yield different results, Eqns 6 and 7, respectively, but the derivation is not self-evident. Here we retained the meaning of  $\bar{y}$  in Eqn 6. In Eqn 7, which was recommended as the only possibility,<sup>9</sup> both  $\bar{y}$  and  $\bar{x}$  were taken as zero, i.e. at the fixed origin.

$$R_6^2 = \frac{(\sum xy)^2 - (\sum x^2)(\sum y)^2/N}{\sum x^2 \sum (y - \bar{y})^2} \quad (6)$$

$$R_7^2 = \frac{(\sum xy)^2}{\sum x^2 \sum y^2} \quad (7)$$

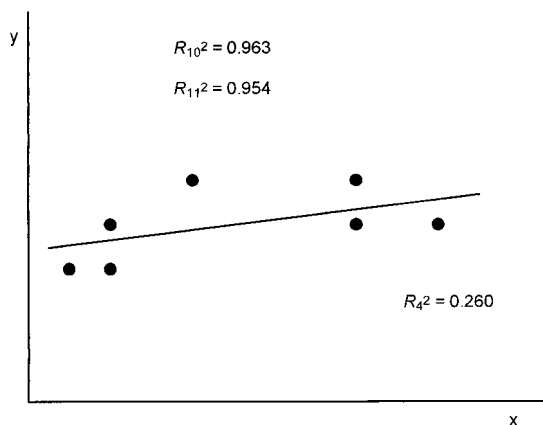
Following the definitions, the meanings of  $R_6^2$  and  $R_7^2$  are different. Whereas  $R_6^2$  characterizes the adherence to the model and the success of predicting within the data set,  $R_7^2$  describes the fit viewed from the point (0,0), with respect to the whole size of  $y$ . In extreme cases, particularly when all data points are distant from the origin,  $R_6^2$  and  $R_7^2$  may differ sharply and  $R_6^2$  may be even negative.<sup>9</sup> Such an extreme case is shown in Fig. 1(a). A negative  $R_6^2$  means that the model is completely wrong in this case; a better prediction would be to predict for all points the mean value of  $y$ :  $\hat{y}_i = \bar{y}$ . On the other hand,  $R_7^2$  tells us that the prediction is not so bad with respect to the whole value of  $y$ : from the point (0,0) we see the predicted and real values in similar directions.



**Figure 1.** Schematic representation of linear regression with the fixed intercept (model B) in two extreme cases of bad fitting; scaling on both axes is arbitrary. Regression lines (full lines) are controlled by the hard condition of the fixed intercept: very low values of the coefficient of determination  $R_6^2$  indicate that the model is wrong and unable to predict the values of  $y$  from  $x$ . In Fig. (b) the prediction is bad even with respect to the whole value of  $y$  (low  $R_7^2$ ) but in Fig. (a) it is not so bad (reasonable  $R_7^2$ ; the points are seen from the intercept in the right direction). When the constraint of the fixed intercept is left, the model is changed to model A with a much better fit (broken lines and coefficients  $R_4^2$ )

Another extreme case, Fig. 1(b), shows dramatically the consequences of the hard condition of a fixed intercept. All points are situated on the same side of the regression line: the fit is bad according to both  $R_6^2$  and  $R_7^2$ . A characteristic feature of both  $R_6^2$  and  $R_7^2$  is also their sensitivity to shift (Table 1). When the data points are far from the intercept,  $R_6^2$  and  $R_7^2$  may differ considerably from each other and also from  $R_4^2$  in the model A. On the other hand, for  $\bar{x} = \bar{y} = 0$ ,  $R_4^2 = R_6^2 = R_7^2$ .

Figure 1 also shows that the models A and B can yield divergent results. Therefore, the common correlation coefficient  $R_4^2$  must not be introduced into the no-intercept model. An often used statistic<sup>10</sup> of this model is  $f$  [Eqn 8]. In this equation,  $SD$  is the standard deviation



**Figure 2.** Schematic representation of linear regression with errors in both variables, model C, in an extreme case with a small slope; scaling is arbitrary but equal on both axes. Adherence to the regression line is not bad according to  $R_{10}^2$  (or  $R_{11}^2$  with the same significance), but a prediction of  $y$  from  $x$  would be ineffective (low value of  $R_4^2$ )

from the regression line, which is connected to  $R_7^2$  [see Eqn 8, right part]. A similar parameter  $\psi$ , suggested many years ago by Exner,<sup>12</sup> is related in a similar way to  $R_4^2$  and belongs to the model A but its main significance is outside the regression models. Including  $f$  in model A, for instance,<sup>5b,13</sup> is a misunderstanding, similar to using  $R_4^2$  in the no-intercept model.<sup>14</sup>

$$f = \frac{SD}{\sqrt{\sum y^2/N}} = \sqrt{\frac{N}{N-1}} (1 - R_7^2) \quad (8)$$

In practical applications, the no-intercept model is warranted particularly in the following situations: (a) the point (0,0) has an absolute significance (for instance, null dose–no response); (b) one point has been obtained experimentally with an extraordinary precision; (c) each value  $y_i$  has been determined against a reference  $y_i^0$ , measured simultaneously;<sup>11</sup> the actual variable is then  $(y - y^0)_i$ . These conditions are not met in the most popular chemical applications;<sup>8b,10</sup> there the point (0,0) may sometimes be known with an improved precision and sometimes not; in no case has its higher precision been statistically proven. In most of these applications, however,  $\bar{y}$  is not far from the origin, so that the difference between models A and B is not dramatic.

### (C) Functional model

This model assumes both variables loaded with error and freely fitted intercept. The whole model is symmetrical with respect to the two variables, which means that their interchange does not affect the results. The errors in both

variables can be assumed with the same variance and independent. If the variances are not the same, one variable may be rescaled. The least-squares condition (see Table 1) refers to the squared perpendicular distances from the regression line. It gives the estimate of the slope  $\beta$  according to the equation

$$b_9 = \frac{\sum (y - \bar{y})^2 - \sum (x - \bar{x})^2}{2 \sum (x - \bar{x})(y - \bar{y})} + \frac{\sqrt{[\sum (y - \bar{y})^2 - \sum (x - \bar{x})^2]^2 + 4[\sum (x - \bar{x})(y - \bar{y})]^2}}{2 \sum (x - \bar{x})(y - \bar{y})} \quad (9)$$

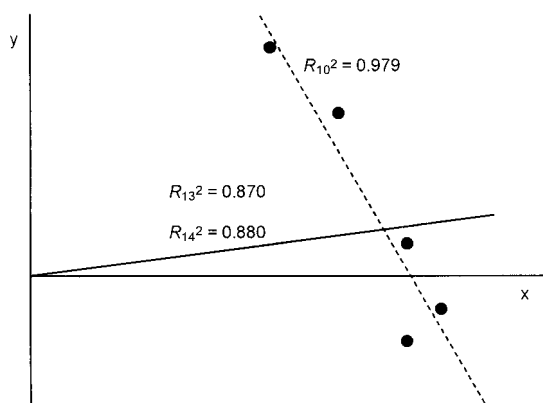
When applying Eqn 1, one must give a new meaning to the denominator, which generally means the information entropy before the regression. We suggest here the following definition:

$$R_{10}^2 = 1 - \frac{RSS}{RSS_{\max}} = \frac{(b_9^2 + 1) \sum (x - \bar{x})(y - \bar{y})}{b_9 \sum (y - \bar{y})^2 + \sum (x - \bar{x})(y - \bar{y})} \quad (10)$$

$RSS$  means the minimum sum of squares when the regression line has the optimum slope and  $RSS_{\max}$  means the residual sum of squares when the line has the worst possible slope (perpendicular to the optimum slope). An alternative definition of the denominator leads to Eqn 11, which corresponds to that used in a previous paper.<sup>11</sup> The properties of  $R_{10}^2$  and  $R_{11}^2$  are similar, but we prefer the former, which can be understood in a simpler way. The main defect of  $R_{11}^2$  is that it is always greater than 0.5. In the region of close correlations their values are near to each other, generally  $R_{11}^2 > R_{10}^2 > R_4^2$ .

$$R_{11}^2 = 1 - \frac{RSS}{\sum (x - \bar{x})^2 + \sum (y - \bar{y})^2} = \frac{\sum (x - \bar{x})^2 + b_9 \sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2 + \sum (y - \bar{y})^2} \quad (11)$$

The second definition of  $R^2$  [Eqn 2] yields the same result as in model A, viz. Eqn 4. Again, the meanings of  $R_{10}^2$  and  $R_4^2$  in this model are different.  $R_{10}^2$  characterizes the adherence of the points to a straight line, irrespective of its slope  $b$ . Its invariance to rotation of coordinates is the most striking feature (Table 1); we may imagine that it characterizes a property of an array of points irrespective of its position in the plane. On the other hand,  $R_4^2$  in this model describes the effectiveness of predicting  $y$  from  $x$  or, with respect to the model symmetry, also  $x$  from  $y$ . A great difference between the two concepts is to be



**Figure 3.** Schematic representation of linear regression with the fixed intercept and errors in both variables, model D, in an extreme case of bad fitting; scaling is arbitrary but equal on both axes. Adherence to the line is apparently fairly good (according to  $R^2_{13}$  or  $R^2_{14}$ ), but this is caused by the hard condition of a fixed intercept and by the artificial comparison with a perpendicular line involved in Eqn 13. When the constraint of the fixed intercept is left, the model is changed to model C (broken line and coefficient  $R^2_{10}$ )

expected when the regression slope is far from unity. Such an extreme case is shown in Fig. 2. Adherence to the regression line is not so bad ( $R^2_{10}$ ) but an efficient prediction is impossible ( $R^2_4$ ). Even in more realistic examples,  $R^2_{10}$  is generally greater than  $R^2_4$  for the same data set, and the two values must not be directly compared.

In practice, this model is applicable particularly in the following cases: (a) the two variables are given in the same units and their experimental errors are equal; (b) the two variables are given in different units and scaled in such a way that the experimental errors are represented by equal given lines (this standardization seems to us more reasonable than that based on the standard deviations<sup>15</sup>  $\sigma_x$  and  $\sigma_y$ ). The case when the variables are in the same units and their precisions differ considerably may be better approximated by model A.

#### (D) Functional no-intercept model

In this model, both variables are loaded with error and the intercept is fixed. Again, the errors in either coordinate are equally dispersed and independent, and the fixed intercept is at the point (0,0). The regression line passes through the origin as in model B, the least-squares condition relates to the perpendicular distances as in model C (see Table 1). The estimate of the slope  $\beta$  is given by Eqn 12 and the coefficient of determination by Eqn 13; both equations follow from Eqns 9 and 10 for  $\bar{x} = \bar{y} = 0$ .

$$b_{12} = \frac{\sum y^2 - \sum x^2 + \sqrt{(\sum y^2 - \sum x^2)^2 + 4(\sum xy)^2}}{2 \sum xy} \quad (12)$$

$$R^2_{13} = \frac{(b_{12}^2 + 1) \sum xy}{b_{12} \sum y^2 + \sum xy} \quad (13)$$

The other possibility is  $R^2_{14}$  in Eqn 14, used previously,<sup>11</sup> which follows in the same way from Eqn 11. The difference between  $R^2_{13}$  and  $R^2_{14}$  is slight, as it was between  $R^2_{10}$  and  $R^2_{11}$

$$R^2_{14} = 1 - \frac{RSS}{\sum x^2 + \sum y^2} = \frac{\sum x^2 + b_{12} \sum xy}{\sum x^2 + \sum y^2} \quad (14)$$

The second, less suitable, definition of the coefficient of determination according to Eqn 2 yields the same result<sup>11</sup> as in model B, Eqn 7. The meaning of these two definitions is less easy to understand than in the preceding models since the properties of models B and C are combined. The most striking feature of  $R^2_{13}$  is again its invariance to rotation of coordinates (Table 1). It expresses the adherence to a straight line passing through the origin, using a perpendicular straight line as reference. It is therefore sensitive to the distance from the origin. On the other hand,  $R^2_7$  is oriented to the prediction of  $y$  from  $x$  or vice versa. Differences are expected particularly when the slope  $b$  is far from unity. In an extreme case (Fig. 3),  $R^2_{13}$  is high since the origin is relatively distant, although the adherence to the model is evidently not good.  $R^2_7$  is low since any prediction would be very bad. None of these values bears any relation to the adherence to the better model C (dashed line). This is due to the hard condition of a fixed intercept and shows conclusively how misleading such a hard condition can be when it is not realistic.

The conditions when this model is preferable are not encountered often. In one case,<sup>11</sup>  $x$  and  $y$  were measured in the same units and at the same accuracy (as in model C) against the standards; the actual variables were thus the values  $x - x^0$  and  $y - y^0$  (as in model B).

## CONCLUSIONS

The coefficient of determination can be defined in a different way and no definition can be *a priori* preferred. Only in classical linear regression (model A) do all definitions yield the same equation [Eqn 4], which meets all requirements for this statistic.<sup>9</sup> In the other models, B–D, two or more meaningful definitions of  $R^2$  are possible but their values cannot be compared with each other and the less so between individual models.

The main problem is already in the choice of the proper model. In our opinion, the no-intercept regression (model B) is acceptable only if it is dictated by strict physical

grounds; this is not often the case in chemical applications. In common cases, when the zero intercept is theoretically anticipated but may be also loaded with error, we recommend, in agreement with others,<sup>16</sup> the use of the common regression, model A. When a statistically insignificant value is obtained for the intercept, the anticipation has been confirmed and there is no strong necessity to use the no-intercept model B. On the other hand, the improper use of model B can lead to great mistakes (see Fig. 1). In practice, the consequences will not be so great, particularly in the two following limiting cases. In the first case, the center of gravity of all points  $(\bar{x}, \bar{y})$  is not far from the origin, then  $R_4^2$ ,  $R_6^2$  and  $R_7^2$  are near to each other. In the second case, the correlation is so close that  $R^2$  according to any definition is approaching unity.

The functional model (C or D) is different in principle but even here there is some possibility of replacing it by the common type A. One such case is when the error of one variable is much smaller than that of the other. Then the former variable can be taken as exact with a good approximation. Another such case arises when one variable must be considered as independent and can never be predicted from the other (time, temperature). Sometimes the values of this variable can be taken as exact if they are adjusted to a certain value, i.e. not measured<sup>17</sup> (for instance,  $x$  would denote the temperature for which the thermostat was adjusted, not the actual temperature in the reaction vessel). In the statistical sense, such a value represents the expected value for all possible experiments with the given value of  $x$ , hence the whole experimental error is shifted to the dependent variable  $y$ . When the functional model cannot be replaced by the common type with a reasonable approximation, one has to use one of the two types of  $R$  discussed here. The choice between them should depend on the purpose:  $R_{10}$  or  $R_{13}$  quantify the adherence to a straight line and are independent of rotation whereas  $R_4$  or  $R_7$  express the strength of prediction of  $y$  from  $x$  or vice versa.

In conclusion, the classical model of regression and the common equation for  $R_4^2$  can be used to a broader extent than follows from the strict conditions of their derivation. In the cases when it cannot be used, a particular definition of  $R^2$  must be specified and values of  $R^2$  according to different definitions must not be compared. Let us stress that the coefficient of determination is never satisfactory for judging the fit or the importance of an empirical correlation: it should be always considered in connection with other statistics<sup>2</sup> (e.g. residual standard deviation, standard deviations of the regression coefficients). Of

course, all statistical tools can prove only a disagreement with the model, not agreement.

Note that the relationships developed here could in principle be extended to multiple linear regression and the pertinent equation could be developed. In addition to the standard model A, model B may also have practical importance whereas for models C and D it is difficult to imagine a situation with three equivalent interchangeable variables.

## Acknowledgements

This work was supported by grant 203/96/1658 from the Grant Agency of the Czech Republic (to O. E.) and by grant Copernicus JRP-10053 (to K.Z.)

## REFERENCES

1. O. Exner, *Correlation Analysis of Chemical Data*, Chapt. 8. Plenum Press, New York (1988).
2. P. N. Craig, C. Hansch, L. W. McFarland, Y. C. Martin, W. P. Purcell and R. Zahradnik, *J. Med. Chem.* **14**, 447 (1971); M. Charton, S. Clementi, S. Ehrenson, O. Exner, J. Shorter and S. Wold, *Quant. Struct.-Act. Relat.* **4**, 29 (1985).
3. H. H. Jaffé, *Chem. Rev.* **53**, 191–261 (1953).
4. S. Bradamante and G. A. Pagani, *J. Chem. Soc., Perkin Trans. 2* 1055–1061 (1986); M. S. Westwell, M. S. Searle, D. J. Wales and D. H. Williams, *J. Am. Chem. Soc.* **117**, 5013–5015 (1995).
5. (a) W. F. Reynolds, A. Gomes, W. D. MacIntyre, A. Tanin, G. K. Hamer and I. R. Peat, *Can. J. Chem.* **63**, 2376–2384 (1983); (b) A. Cornélis, S. Lambert, P. Laszlo and P. Schaus, *J. Org. Chem.* **46**, 2130–2134 (1981); (c) N. Inamoto, S. Masuda, A. Terui and K. Tori, *Chem. Lett.* 107–110 (1972).
6. C. K. Hancock, *J. Chem. Educ.* **42**, 608–609 (1965); Analytical Methods Committee, *Analyst*, **113**, 1469–1471 (1988).
7. S. Weisberg, *Applied Linear Regression*. Wiley, New York (1980).
8. (a) K.-J. Appenroth, *Z. Phys. Chem. (Leipzig)* **262**, 374–376 (1981); (b) S. Ehrenson, *J. Org. Chem.* **44**, 1793–1797 (1979); (c) W. H. Davis and W. A. Pryor, *J. Chem. Educ.* **53**, 285–287 (1976).
9. T. O. Kvålseth, *Am. Stat.* **39**, 279–285 (1985).
10. S. Ehrenson, R. T. C. Brownlee and R. W. Taft, *Prog. Phys. Org. Chem.* **10**, 1–80 (1973).
11. M. Decouzon, O. Exner, J.-F. Gal and P.-C. Maria, *J. Phys. Org. Chem.* **7**, 615–624 (1994).
12. O. Exner, *Collect. Czech. Chem. Commun.* **31**, 3222–3251 (1966).
13. S. Datta, A. De, S. P. Bhattacharyya, C. Mehdi, A. K. Chakravarty, J. S. A. Brunskill, S. Fadoujou and K. Fish, *J. Chem. Soc., Perkin Trans. 2* 1599–1605 (1988).
14. C.-R. Kramer, *Z. Phys. Chem. (Leipzig)* **261**, 745–758 (1980); I. Lee, Y. S. Lee, B.-S. Lee and H.-W. Lee, *J. Chem. Soc., Perkin Trans. 2* 1441–1445 (1993).
15. D. York, *Can. J. Phys.* **44**, 1079–1086 (1966).
16. S. Wold and M. Sjöström, *Chem. Scr.* **2**, 49–55 (1972); S. Clementi, F. Fringuelli, P. Linda and G. Savelli, *Gazz. Chim. Ital.* **105**, 281–292 (1975).
17. V. Štěpánek, *Sb. Vys. Sk. Chem-Technol. Praze, Anorg. Chem.* 433–442 (1964).