

Error Analysis Using the Variance–Covariance Matrix

Carl Salter

Department of Chemistry, Moravian College, Bethlehem, PA 18018; csalter@cs.moravian.edu

Modern microcomputer software makes least-squares (LS) fits quite tractable (1–5). An important task in any fit is the propagation of error (6, 7) from the measurements into the fit parameters. Harris showed how to use Solver in Excel to perform a linear LS fit and implemented a jackknife method for estimating the errors in the fit parameters (8). Using Harris's example and others requiring nonlinear least squares, de Levie showed that an Excel macro could be used to compute errors in the fit parameters obtained from Solver (9). Zielinski and Allendoerfer reviewed the theory of nonlinear LS and used Mathcad to analyze kinetics data, obtaining estimates of rate constants and their errors (10). More generally, they noted that the errors in the LS parameters are obtained from the diagonal elements of the "error matrix", which is the inverse of the curvature matrix for the sum-of-squared errors (SSE) surface in parameter space.

But how is the error computed when the quantity of interest is not a fit parameter, but some derived property such as an extrapolated value of the function, or its x intercept, or the area under the curve? Such quantities are computed from two or more fit parameters and, because the parameters from the fit are correlated, the usual propagation of error equation does not apply. The errors in these *derived quantities from the fit* must be computed from a more elaborate and more general version of the error propagation equation that includes terms involving the off-diagonal elements of the error matrix. In a fit to a straight line there is one unique off-diagonal element, and its size indicates the degree of correlation between the slope and the y intercept. This element is commonly called the *covariance*.

The x -intercept is a good example of a derived quantity from a LS fit, and in two common chemistry experiments the x -intercept of the fit line is the key quantity: the Charles's law determination of absolute zero from a pressure-vs-temperature fit (11, 12), and the determination of an analyte concentration by standard additions (13, 14). The x intercept is obtained from the slope m and the y -intercept b as $x_{\text{int}} = -b/m$; because m and b are correlated, the covariance must be used to estimate the error in x_{int} . Bruce and Gill showed that neglecting covariance in standard additions analyses can lead to an underestimated standard error in the analyte concentration, and they noted that one analytical textbook has an error equation that neglects covariance (14). Meyer illustrated the effect of neglecting covariance in extrapolating a linearized Clausius–Clapeyron fit, showing that the usual propagation of error equation can overestimate the error of extrapolation by more than an order of magnitude (15). Note that in these two examples neglecting covariance leads to opposite effects on the calculated error.

These two papers notwithstanding, the topic of covariance in LS error analysis has largely fallen through the cracks. Meyer notes that regression packages in spreadsheets do not supply covariances, and this may be one reason why the topic

has received scant attention. But another may be the unfamiliarity of the general expression for propagated error in functions of correlated variables, *which takes on a particularly simple form in matrix notation*. This matrix equation for propagated error cannot be found in the standard chemical education references. The paper on nonlinear LS fitting by Wentworth (16a), which follows almost exactly the analysis of Deming (17), contains the correct formula for determining error in derived quantities, but doesn't write it as a matrix equation. His follow-up paper gives only a simple error-band example of the equation (16b). The second edition of Harris's *Quantitative Chemical Analysis* textbook also contains the formula for the error in derived quantities of fits, but again not as a matrix equation (18). The comprehensive treatment of LS fitting in *Experiments in Physical Chemistry* fails to mention errors in derived quantities (19). Boqué, Rius, and Massart mentioned the problem of covariance briefly with regard to straight line fits (4). Even Meyer avoids matrix notation and doesn't describe how to compute the needed covariance term.

The purpose of this paper is to present the matrix equation for propagated error, and in particular to show the relationship between the *variance–covariance matrix* \mathbf{V} and the error matrix of ref 10. The use of \mathbf{V} will be illustrated for computing the errors of functions of LS parameters in the commonly occurring case of a fit to a straight line, $y = b + mx$. For this fit the covariance can easily be computed and added to a spreadsheet regression. Since spreadsheets do provide for matrix multiplication, \mathbf{V} can be used to propagate error in any function of the fit parameters.

The Variance–Covariance Matrix for the Linear LS Fit to a Straight Line

A set of n ordered pairs of data (x_i, y_i) can be fitted to the equation $y = b + mx$. Assuming that all the experimental error is contained in the y values, minimizing the SSE in the y_i 's leads to the following matrix equation:

$$\begin{pmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{pmatrix} \begin{pmatrix} b \\ m \end{pmatrix} = \begin{pmatrix} \sum y_i \\ \sum x_i y_i \end{pmatrix}$$

or $\mathbf{A} \quad \mathbf{p} = \mathbf{y}$

where the summations are over the n data pairs. The equation can be solved for the parameter vector \mathbf{p} by inverting the matrix \mathbf{A} and multiplying on the left and right by the inverse,

$$\begin{pmatrix} b \\ m \end{pmatrix} = \mathbf{p} = \mathbf{A}^{-1} \mathbf{y} = \frac{1}{D} \begin{pmatrix} \sum x_i^2 & -\sum x_i \\ -\sum x_i & n \end{pmatrix} \mathbf{y}$$

$$D = \det(\mathbf{A}) = n \sum x_i^2 - \left(\sum x_i \right)^2$$

\mathbf{A}^{-1} is equivalent to the error matrix discussed by Zielinski and Allendoerfer for nonlinear fits (10), of which linear fits are a special subset.

For any LS fit, the estimated fit variance s_y^2 is SSE/v , where the number of degrees of freedom is $v = n - p$ and p is the number of fit parameters. The fit variance is the key indicator of the goodness of the fit. A good fit should have a variance roughly equal to the variance of the y -measuring instrument (if the assumption that y is the only source of error is correct). For the simple case of a fit to slope and y intercept, the fit variance is

$$s_y^2 = \frac{\sum (y_i - mx_i - b)^2}{n - 2}$$

and the variance–covariance matrix \mathbf{V} is

$$\mathbf{V} = s_y^2 \mathbf{A}^{-1} = \frac{s_y^2}{D} \begin{pmatrix} \sum x_i^2 & -\sum x_i \\ -\sum x_i & n \end{pmatrix} \quad (1)$$

Equation 1 is equivalent to eq 41 of ref 16a. The diagonal elements of \mathbf{V} are the usual variances for the y intercept and the slope: $s_b^2 = \mathbf{V}_{11} = s_y^2 \sum x_i^2 / D$, and $s_m^2 = \mathbf{V}_{22} = s_y^2 n / D$. The covariance is the off-diagonal element $\mathbf{V}_{12} = \mathbf{V}_{21} = -s_y^2 \sum x_i / D$.¹ The covariance contains the sum of x_i^3 , a term not present in s_b^2 or s_m^2 . Notice that the covariance will be zero if the mean of the x_i 's is zero, which means that a suitable choice of x_i 's (experimental design) will produce a slope and y intercept that are independent.

Errors in Functions of the Fit Parameters

If F is some function of the slope and y intercept $F(b, m)$, the variance in F can be computed from

$$s_F^2 = \mathbf{d}_F^t \mathbf{V} \mathbf{d}_F \quad (2)$$

where

$$\mathbf{d}_F = \begin{pmatrix} \frac{\partial F}{\partial b} \\ \frac{\partial F}{\partial m} \end{pmatrix}$$

and \mathbf{d}_F^t is its transpose (a row vector). Note that in \mathbf{d}_F the derivative with respect to the b always comes first because b was taken to be the first parameter. Equation 2 is the matrix equivalent of eq 43 in ref 16a, the latter being given in algebraic form for the specific case of three parameters. *Equation 2 is the key equation in LS error analysis.*

Since $\mathbf{V} = s_y^2 \mathbf{A}^{-1}$, it follows that $s_F^2 = s_y^2 \mathbf{d}_F^t \mathbf{A}^{-1} \mathbf{d}_F$, which shows that the variance in F depends on the curvature of the SSE surface and on the overall quality of the fit.

Analytical Examples of the Use of the Variance–Covariance Matrix

Let's employ the formalism of eq 2 to obtain error expressions for some derived quantities from a straight-line fit.

Extrapolation and Interpolation $y(x')$

The fit function can be used to estimate y for any value of x ($= x'$). The function F in this case is the fit function itself $F = y_{\text{ext}} = b + mx'$, and the derivative matrix is $\mathbf{d}_{y_{\text{ext}}}^t = (1 \ x')$. Using eq 2 the variance is

$$s_{y_{\text{ext}}}^2 = (1 \ x') \mathbf{V} \begin{pmatrix} 1 \\ x' \end{pmatrix} = \frac{s_y^2}{D} (1 \ x') \begin{pmatrix} \sum x_i^2 & -\sum x_i \\ -\sum x_i & n \end{pmatrix} \begin{pmatrix} 1 \\ x' \end{pmatrix} = \frac{s_y^2}{D} (nx'^2 - 2x' \sum x_i + \sum x_i^2) \quad (3)$$

Equation 3 is the vertical error band equation. The positive and negative square roots represent the range of ± 1 standard deviation about the fit function. When the error band is multiplied by an appropriate value from the Student's t table it becomes the confidence band. It is easy to show that the minimum error occurs when x' equals the mean of the x_i 's, and the band flares out at both ends beyond the observed range of x .

Extrapolation and Interpolation $x(y')$

The x value corresponding to a particular value y' is given by $x_{\text{ext}} = (y' - b)/m$. The derivative matrix for this function is $\mathbf{d}_{x_{\text{ext}}}^t = (-1/m \ (b - y')/m^2)$. The variance of x_{ext} is

$$s_{x_{\text{ext}}}^2 = \frac{s_y^2}{Dm^2} \left(\frac{n(y' - b)^2}{m^2} + \frac{2(b - y')}{m} \sum x_i + \sum x_i^2 \right)$$

Substituting into this expression the equation $x_{\text{ext}} = (y' - b)/m$ leads to

$$s_{x_{\text{ext}}}^2 = \frac{s_y^2}{Dm^2} (nx_{\text{ext}}^2 - 2x_{\text{ext}} \sum x_i + \sum x_i^2) \quad (4)$$

Equation 4, the horizontal error band equation, is remarkably similar to eq 3. Once again, the variance is minimized when x_{ext} is the mean of the x_i 's. Equation 4 is obtained directly from eq 2 when $\mathbf{d}_{x_{\text{ext}}}^t$ is expressed as $-1/m(1 \ x_{\text{ext}})$.

A special case of eq 4 is the error in the x -intercept, where $y' = 0$. The function is $x_{\text{int}} = -b/m$, and $\mathbf{d}_{x_{\text{int}}}^t$ is $-1/m(1 \ x_{\text{int}})$. Inserting this derivative matrix into eq 2 yields

$$s_{x_{\text{int}}}^2 = \frac{s_y^2}{Dm^2} (nx_{\text{int}}^2 - 2x_{\text{int}} \sum x_i + \sum x_i^2) \quad (5)$$

Cast in slightly different form, this is eq 16 of ref 14, derived from eq 13 of ref 14, which Bruce and Gill call the "extrapolation method" of determining the error in the analyte concentration by standard additions.

Calibration Error

Often a regression line is used as a calibration curve: y is measured for an unknown sample so that x of the unknown can be determined: $x_{\text{cc}} = (y_{\text{meas}} - b)/m$. The measurement of y adds an additional source of error. The variance of x_{cc} is the sum of the variance derived from y_{meas} and the variance from the LS fit:

$$s_{x_{cc}}^2 = \left(\frac{\partial x_{cc}}{\partial y_{meas}} \right)^2 s_{y_{meas}}^2 + \mathbf{d}_{cc}^t \mathbf{V} \mathbf{d}_{cc} \quad (6)$$

where \mathbf{d}_{cc}^t equals $\mathbf{d}_{x_{ext}}^t$ given above. The derivative $\partial x_{cc}/\partial y_{meas}$ is $1/m$; if the unknown is measured n_{unk} times, then $s_{y_{meas}}^2 = s_y^2/n_{unk}$. Substitution into eq 6 yields

$$s_{x_{cc}}^2 = \left(\frac{1}{m} \right)^2 \frac{s_y^2}{n_{unk}} + \frac{s_y^2}{Dm^2} (nx_{cc}^2 - 2x_{cc} \sum x_i + \sum x_i^2) \quad (7)$$

$$s_{x_{cc}}^2 = \frac{s_y^2}{m^2} \left(\frac{1}{n_{unk}} + \frac{1}{D} (nx_{cc}^2 - 2x_{cc} \sum x_i + \sum x_i^2) \right) \quad (8)$$

Equation 8 might be called the “calibration curve” error equation; it can be found in analytical chemistry (20) and chemometrics (21) textbooks. The term s_y^2 appears twice: once from the measurement(s) on the unknown and again from the LS fit, because the fit variance should be the same as the variance of the y -measuring instrument. The sensitivity of the y -measuring instrument is the slope m ; eq 8 indicates that a more sensitive instrument should produce a smaller error in the determination of x . The smallest variance for determining x is observed for values of x centered in the range of the calibration data. Good experimental design dictates that the calibration curve be constructed in the region of the unknown's y value. No amount of signal-averaging on the unknown sample can eliminate the error introduced by the calibration fit, but of course more calibration points will reduce this error.

As presented here the calibration curve error equation is a variant of the horizontal error band. The confidence band equation described by Burdge et al. (22) is an equivalent

expression based on the vertical error band equation (eq 3, cast in different form, is contained in eq 1 of ref 22). Burdge et al. show how this confidence band has been used by IUPAC to define detection limits. It would be more straightforward to define detection limits using horizontal error bands, but the difference is probably not significant.

Area (x_1, x_2)

The error of the definite integral of the fit line (the area under the curve) demonstrates the generality of this method. The area between the limits x_1 and x_2 is

$$A = \int_{x_1}^{x_2} (b + mx) dx = b(x_2 - x_1) + \frac{m}{2}(x_2^2 - x_1^2)$$

The derivative matrix is

$$\mathbf{d}_A^t = \left((x_2 - x_1) \quad \frac{(x_2^2 - x_1^2)}{2} \right)$$

The variance of the area is

$$s_A^2 = \frac{s_y^2}{D} \left(\frac{n}{4} (x_2^2 - x_1^2)^2 - (x_2^2 - x_1^2)(x_2 - x_1) \sum x_i + (x_2 - x_1)^2 \sum x_i^2 \right) \quad (9)$$

Spreadsheet Example: Standard Additions

Following example 1 in ref 14, we fit absorbance vs standard concentration to $y = b + mx$ and obtain the analyte concentration as $x_{int} = b/m$. The fit and error analysis are illustrated in the Excel spreadsheet in Figure 1. The standard additions data and fit results are in columns A–C; \mathbf{V} and the error analysis are in columns D–H.

	A	B	C	D	E	F	G	H	I
1	concn	A		COVARIANCE and ERROR calculation Using V					
2	0	0.240		sum of x's	55.5	is SUM(A5:A9)			
3	5.55	0.437							
4	11.10	0.621		covar	-8.5E-07	is -E2*C26^2/B19			
5	16.65	0.809							
6	22.20	1.009		V matrix			det mat		
7				1.416E-05	-8.5E-07		-29.0576	is -1/B26	
8				-8.505E-07	7.66E-08		203.656	is -E18/B26	
9									
10							V*det		
11							-0.00058	is MMULT(D7:E8,G7:G8)	
12	SUMMARY OUTPUT			det mat			4.03E-05		
13				-29.05759	203.656				
14	Regression Statistics								
15	Multiple R	0.999903							
16	R Square	0.999806		variance	0.025199	is MMULT(D13:E13,F11:F12)			
17	Adjusted R	0.999741							
18	SE	0.004858		x-int	-7.00869		unknown	concn mg/L	
19	# Obs	5		error	0.158742		7.01	+/- 0.51	
20				conf lim	0.505118				
21				t(95%,3)=	3.182				
22									
23									
24		Coefficient	Std Error						
25	Intercept	0.2412	0.003763		V(1,1) is	C25^2			
26	slope	0.034414	0.000277		V(2,2) is	C26^2			
27									

Figure 1. Excel spreadsheet demonstrating error propagation using the \mathbf{V} matrix from a LS fit of data from example 1 of ref 14 for the method of standard additions. The LS fit is extrapolated to its x intercept (cell E18); \mathbf{V} (cells D7:E8) is used to estimate the error in the extrapolation (cell E19) following eq 2 with $\mathbf{d}_{x_{int}}^t = (-1/m)(1 - x_{int})$ in cells G7:G8. The regression output was created by the Regression tool in the Analysis Tool-Pak of Excel; a small portion of the regression output was cut and pasted into the spreadsheet to complete the error analysis. Important numerical results are outlined in dark boxes; 1-D matrices are in italics.

Excel's Regression Tool produces a lot of output; only the summary and a portion of the parameter statistics were retained for the analysis. The x -intercept is computed using $-b/m$ in cell E18. Σx , needed for the covariance, is computed in E2. The covariance is computed using $-\Sigma x s_m^2/n$ in cell E4; $V_{12} = V_{21} = E4$. V is completed by adding $V_{11} = s_b^2$ and $V_{22} = s_m^2$. The derivative matrix $-1/m(1 \ x_{\text{int}})$ appears twice in the spreadsheet: once as a column matrix (G7:G8) and again as a row matrix (D13:E13). Because of the way matrix multiplication is implemented in Excel, this duplication is required. The column matrix at F11:F12 is the product of V and the derivative column matrix; this intermediate result must be stored on the spreadsheet. The derivative row matrix and this intermediate are multiplied at E16 to yield the variance. The standard error in cell E19 and the 95% confidence limit in E20 are equal to those reported in ref 14 using the "extrapolation method". However, the extrapolation method does not explicitly contain the covariance, and its use requires the additional calculation of Σy .

As Meyer points out, very often the error in a derived quantity is computed disregarding the covariance. Ignoring the covariance is equivalent to "zeroing-out" the off-diagonal elements of V . This is easily done in the spreadsheet. Placing zeros in cells E7 and D8 yields a standard error of 0.39—the same error reported by Bruce and Gill using the incorrect "algebraic method", which, they point out, leaves out covariance. The matrix formalism makes the covariance explicit, and the spreadsheet makes it easy to assess the effect of the covariance on the error in any derived quantity.

Conclusion

When a linear regression has been performed using a spreadsheet program, the covariance can be calculated using $-\Sigma x s_m^2/n$, where it is assumed that s_m is available from the program—the only additional overhead is the calculation of Σx . V can then be assembled and used to propagate error. The spreadsheet can be used as a template: any regression output can be cut and pasted into the template to produce the corresponding V , and the error in any derived quantity can be computed by entering the appropriate derivative matrix. The matrix formalism is also a powerful way to perform error propagation on a hand-held calculator, many of which now can perform matrix multiplication.

Though we examined only the simple case of fits to a straight line, V exists for any LS fit, whether linear or nonlinear, and the relationship between V and the error matrix is always the same. All unweighted linear LS fits start with the creation of a matrix A , which depends only on the x quantities of the data. A^{-1} is the error matrix, and $V = s_y^2 A^{-1}$. Nonlinear fits must converge to the error matrix, but once convergence is achieved V is again just the error matrix times s_y^2 . (Using the notation in ref 10, $V = \sigma^2 \epsilon$.) In either case, all errors associated with the fit are contained in V and can be extracted by application of eq 2.

For linear or nonlinear fits performed using Solver, de Levie's Excel macro can be used to produce the V matrix (9). The macro estimates the elements of the error matrix numerically; although it does not ordinarily print out the off-diagonal elements, the addition of a few lines of code will produce the required off-diagonal elements of the V matrix. A test of

this macro using the standard additions data produced the same covariance term.

The variance-covariance matrix is central to error propagation. It connects traditional error propagation among independent variables with errors in fit parameters, because it contains both operations as special cases. The standard propagation of error theorem can be considered a special case of eq 2 in which V is diagonal. Errors in fit parameters can be extracted from V using "trivial" derivative matrices. For example, the "slope function" is $F = m$; its derivative matrix (0 1) will extract the variance in the slope. V and eq 2 demonstrate their greatest utility when they are used to generate errors in derived quantities from LS fits. This formalism can handle quite complicated tasks in error analysis in straightforward fashion; for example, Tellinghuisen has demonstrated a sophisticated use of eq 2 to determine the errors in a potential energy curve derived from LS fits of spectral data (23).

In some applications the jackknife or other bootstrapping methods may be useful alternatives to the matrix formalism. Though Harris demonstrated the jackknife only on a linear LS fit (8), the method is readily extended to nonlinear fits as well, and could be used to estimate the error in derived quantities of a fit. The reliability of the jackknife in estimating errors in nonlinear LS fits is an open question.

Utility aside, is there any general pedagogical value in learning about the variance-covariance matrix—is there any other topic where these methods and ideas may overlap? Yes. The Schrödinger equation is a recipe for the wave function of a system, and with the appropriate differential operator one can extract information about the system from the wave function. We usually say that the wave function contains *all* the information we can ever have about the system. Equation 1 is a recipe for V , and eq 2 shows that one can extract information about errors in the LS fit from V . The variance-covariance matrix contains *all* the error information we can ever have about the LS fit. In this matrix formalism of error propagation there lies an eerie analogy to quantum mechanics, which is after all a statistical model of Nature.

Acknowledgment

I thank Joel Tellinghuisen of Vanderbilt University. This work was inspired by his lectures and research on least-squares fitting.

Note

1. Meyer calls the covariance σ_{mb} ; Bruce and Gill call it s_{mb} . In Bruce and Gill's notation the covariance equals $-\bar{x}s_y^2/S_{xx}$, which is a recasting of the expression for V_{12} . Since the covariance is dimensionally the same as the variance, the notation σ_{mb}^2 or s_{mb}^2 would be preferable; however, many authors shun the "square" notation because they fear it would lead readers to believe that the covariance can never be negative.

Literature Cited

1. Copeland, T. G. *J. Chem. Educ.* **1984**, *61*, 778–779.
2. Irvin, J. A.; Quickenden, T. I. *J. Chem. Educ.* **1983**, *60*, 711–712.
3. de Levie, R. *J. Chem. Educ.* **1986**, *63*, 10–15.

4. Boqué, R.; Rius, F. X.; Massart, D. L. *J. Chem. Educ.* **1994**, *71*, 230–232.
5. Machuca-Herrera, J. O. *J. Chem. Educ.* **1997**, *74*, 448–449.
6. Guedens, W. J.; Yperman, J.; Mullens, J.; Van Poucke, L. C.; Pauwels, E. J. *J. Chem. Educ.* **1993**, *70*, 776–779, 838–841.
7. Donato, H.; Metz, C. *J. Chem. Educ.* **1988**, *65*, 867–868.
8. Harris, D. C. *J. Chem. Educ.* **1998**, *75*, 119–121.
9. de Levie, R. *J. Chem. Educ.* **1999**, *76*, 1594–1598.
10. Zielinski, T. J.; Allendoerfer, R. D. *J. Chem. Educ.* **1997**, *74*, 1001–1007.
11. Garrett, D. D.; Banta, M. C.; Arney, B. E. *J. Chem. Educ.* **1991**, *68*, 667–668.
12. Strange, R. S.; Lang, F. T. *J. Chem. Educ.* **1989**, *66*, 1054–1055.
13. Bader, M. *J. Chem. Educ.* **1980**, *57*, 703–704.
14. Bruce, G. R.; Gill, P. S. *J. Chem. Educ.* **1999**, *76*, 805–806.
15. Meyer, R. D. *J. Chem. Educ.* **1997**, *74*, 1339–1340.
16. Wentworth, W. E. (a) *J. Chem. Educ.* **1965**, *42*, 96–103. (b) *J. Chem. Educ.* **1965**, *42*, 162–167.
17. Deming, W. E. *Statistical Adjustment of Data*; Wiley: New York, 1943.
18. Harris, D. C. *Quantitative Chemical Analysis*, 2nd ed.; Freeman: New York, 1987; p 625.
19. Shoemaker, D. P.; Garland, C. W.; Nibler, J. W. *Experiments in Physical Chemistry*, 5th ed.; McGraw-Hill: New York, 1989; Chapter 20.
20. Harris, D. C. *Quantitative Chemical Analysis*, 5th ed.; Freeman: New York, 1999; pp 100–104.
21. Hecht, H. G. *Mathematics in Chemistry*; Prentice Hall: Englewood Cliffs, NJ, 1990; p 271.
22. Burdge, J. R.; MacTaggart, D. L.; Farwell, S. O. *J. Chem. Educ.* **1999**, *76*, 434–439.
23. Tellinghuisen, J. *J. Mol. Spectrosc.* **1990**, *141*, 258–264.