

Regression Algebra

Robert DeSerio

Abstract

Linear algebra is used to derive the main formulas of linear regression—those determining the fit parameters and their covariance matrix in terms of the data and their covariance matrix. It is also used to demonstrate various assertions, notably that the expectation value of the chi-square statistic χ^2 is equal to the degrees of freedom—the number of data points less the number of fit parameters. The principles and formulas involved are also applied to nonlinear fitting functions and the required approximations and formula modifications are given. Relative to their best-fit values, the change in χ^2 is shown to be a quadratic function of the change in fit parameters. The quadratic is shown to demonstrate the “ $\Delta\chi^2 = 1$ ” rule—that if any fit parameter is offset one standard deviation from its best-fit value, the minimum χ^2 will increase by one. All unspecified parameters must be re-fit with the offset parameter held fixed; multiplying the parameter changes that reach the new minimum $\chi^2 + 1$ by the fixed parameter’s one-sigma offset are shown to give the covariances with the fixed parameter.

Statement and Solution

The dependent variable is represented by the set $\{y_i\}$, its N elements y_i , $i = 1\dots N$, or the column vector \mathbf{y} (of length N whose i th element is y_i). Each y_i is assumed to be a random sample from its own Gaussian distribution. The one-sigma uncertainty σ_i in each y_i is assumed known in advance and determines the standard deviation of the distribution from which y_i is a sample. All y_i are assumed to be statistically independent and thus the $N \times N$ covariance matrix $[\sigma_y]$ has non-zero elements only on the diagonal.

$$[\sigma_y]_{ij} = \sigma_i^2 \delta_{ij} \tag{1}$$

Usually not known in advance, the distributions’ means μ_i , $i = 1\dots N$ (column vector $\boldsymbol{\mu}$) are predicated on experimental conditions associated with point i (such as the value x_i of an independent variable) and theoretical considerations. The theory provides a fitting function $F_i(\{\alpha_j\})$ predicting μ_i in terms of a set of M fitting parameters α_j , $j = 1\dots M$, (column vector $\boldsymbol{\alpha}$).

$$\mu_i = F_i(\{\alpha_j\}) \tag{2}$$

In addition to the *true* parameter values $\{\alpha_j\}$ that determine the true means, another parameter set is of paramount importance. The *best-fit* parameter values are denoted a_j , $j = 1 \dots M$ (column vector \mathbf{a}), and are defined by the condition that they minimize the χ^2 for a particular data set $\{y_i\}$. The best-fit parameters then determine the associated best-fit y -values denoted y_i^{fit} , $i = 1 \dots N$ (column vector \mathbf{y}^{fit})

$$y_i^{\text{fit}} = F_i(\{a_j\}) \quad (3)$$

For linear regression techniques to apply, the functions $F_i(\{a_j\})$ are restricted to a form that is linear in the fitting parameters. One such form arises when an independent variable, say x_i , is specified for each y_i and $F_i(\{a_j\})$ can be expressed as a superposition of linearly independent functions $f_j(x)$ of the independent variable.

$$y_i^{\text{fit}} = \sum_{j=1}^M a_j f_j(x_i) \quad (4)$$

For example, for a quadratic fit: $y_i^{\text{fit}} = a_1 + a_2 x_i + a_3 x_i^2$ and the three corresponding functions are $f_1(x) = 1$, $f_2(x) = x$, $f_3(x) = x^2$. The functions $F_i(\{a_j\})$ are linear in the regression coefficients because the coefficients only appear as amplitudes and do not appear in the functions $f_j(x)$. More specifically, the derivative of all $F_i(\{a_j\})$ with respect to every a_j is independent of the complete set $\{a_j\}$.

Explicitly representing the N values $f_j(x_i)$, $i = 1 \dots N$ as a column vector \mathbf{f}_j , the full set of M such vectors becomes an $N \times M$ matrix $[J]$; one column of length N for each of the M functions. This form is typically required for spreadsheet linear regression programs and when constructed from functions of an independent variable x_i , the matrix elements are

$$[J]_{ij} = f_j(x_i) \quad (5)$$

On the other hand, the column constructions need not be so constrained. For example, entries might depend on more than one independent variable. The only requirement for a unique solution is that all columns be linearly independent.

With this restriction on the fitting function, Eq. 4 can then be expressed

$$\mathbf{y}^{\text{fit}} = [J] \mathbf{a} \quad (6)$$

and the relation between the true means and true parameters becomes

$$\boldsymbol{\mu} = [J] \boldsymbol{\alpha} \quad (7)$$

The chi-square,

$$\chi^2 = \sum_{i=1}^N \frac{(y_i - y_i^{\text{fit}})^2}{\sigma_i^2} \quad (8)$$

can be expressed as the following weighted inner product

$$\chi^2 = (\mathbf{y}^T - \mathbf{y}^{\text{fit}T}) [\sigma_y^{-1}] (\mathbf{y} - \mathbf{y}^{\text{fit}}) \quad (9)$$

where \mathbf{y}^T and $\mathbf{y}^{\text{fit}T}$ are the corresponding transposes (row vectors) and $[\sigma_y^{-1}]$ is the weighting matrix—the inverse of $[\sigma_y]$, an $N \times N$ diagonal matrix with elements

$$[\sigma_y^{-1}]_{ij} = \frac{1}{\sigma_i^2} \delta_{ij} \quad (10)$$

Be sure to appreciate the accuracy and elegance of these simple relations. For example, keep track of the multiplication rules and vector and matrix sizes to verify the sums involved and the dimensionality of any intermediate or final results.

Substituting Eq. 6 and its transpose $\mathbf{y}^{\text{fit}T} = \mathbf{a}^T [J^T]$ into Eq. 9 shows how χ^2 would depend on \mathbf{a}

$$\chi^2 = (\mathbf{y}^T - \mathbf{a}^T [J^T]) [\sigma_y^{-1}] (\mathbf{y} - [J] \mathbf{a}) \quad (11)$$

By temporarily treating the fit parameters as variables, the best-fit $\{a_j\}$ will be determined from the fact that the derivatives $d\chi^2/da_j$ are zero for all a_j at the χ^2 minimum.

An expression for these derivatives is obtained by first noting that

$$d\mathbf{a}/da_j = \boldsymbol{\delta}_j \quad (12)$$

where $\boldsymbol{\delta}_j$ is a column vector of length M with the only non-zero element being a unity element in the j th position. Then, simple chain rule differentiation of Eq. 11 gives

$$\frac{d\chi^2}{da_j} = -\boldsymbol{\delta}_j^T [J^T] [\sigma_y^{-1}] (\mathbf{y} - [J] \mathbf{a}) - (\mathbf{y}^T - \mathbf{a}^T [J^T]) [\sigma_y^{-1}] [J] \boldsymbol{\delta}_j \quad (13)$$

And noting that

$$[J] \boldsymbol{\delta}_j = \mathbf{f}_j \quad (14)$$

along with the transpose equation $\boldsymbol{\delta}_j^T [J^T] = \mathbf{f}_j^T$ gives

$$\frac{d\chi^2}{da_j} = -\mathbf{f}_j^T [\sigma_y^{-1}] (\mathbf{y} - [J] \mathbf{a}) - (\mathbf{y}^T - \mathbf{a}^T [J^T]) [\sigma_y^{-1}] \mathbf{f}_j \quad (15)$$

As they should be, both terms on the right are scalars and are in fact the same scalars, simply formed with expressions that are transposes of one another. Choosing the left-hand expression gives

$$\frac{d\chi^2}{da_j} = -2\mathbf{f}_j^T [\sigma_y^{-1}] (\mathbf{y} - [J] \mathbf{a}) \quad (16)$$

The M equations for the chi-square minimum, i.e., $d\chi^2/da_j = 0$ for each a_j , are then obtained by setting the right side of Eq. 16 equal to zero for each j . These equations can be rewritten

$$\mathbf{f}_j^T [\sigma_y^{-1}] [J] \mathbf{a} = \mathbf{f}_j^T [\sigma_y^{-1}] \mathbf{y} \quad (17)$$

for $j = 1 \dots M$. These equations are just the M components of the vector equation

$$[J^T] [\sigma_y^{-1}] [J] \mathbf{a} = [J^T] [\sigma_y^{-1}] \mathbf{y} \quad (18)$$

One cannot simply cancel the matrix product $[J^T][\sigma_y^{-1}]$ on both sides. $[J^T][\sigma_y^{-1}]$ is not square and does not possess a unique inverse. However, defining

$$[X] = [J^T][\sigma_y^{-1}][J] \quad (19)$$

which is an $M \times M$ matrix and does have a unique inverse $[X^{-1}]$ (if all columns of $[J]$ are linearly independent), Eq. 18 then becomes

$$[X] \mathbf{a} = [J^T][\sigma_y^{-1}] \mathbf{y} \quad (20)$$

Multiplying both sides by $[X^{-1}]$ on the left, gives the solution

$$\mathbf{a} = [X^{-1}][J^T][\sigma_y^{-1}] \mathbf{y} \quad (21)$$

Defining the combination

$$[J^\dagger] = [X^{-1}][J^T][\sigma_y^{-1}] \quad (22)$$

gives the final result

$$\mathbf{a} = [J^\dagger] \mathbf{y} \quad (23)$$

Equation 23 says that the M regression coefficients \mathbf{a} are obtained by a simple multiplication between an $M \times N$ matrix $[J^\dagger]$ and the N -component vector \mathbf{y} . Moreover, the recipe for $[J^\dagger]$ has only two ingredients: the regression functions $[J]$ and the data uncertainties $[\sigma_y]$. It does not depend on the data \mathbf{y} . Consequently, assuming the regression functions and data uncertainties can be determined before an experiment is performed, so too can $[J^\dagger]$. This principal is behind such constructs as Savitsky-Golay filters for smoothing and differentiating data uniformly acquired in time.

$[J^\dagger]$ is called the weighted Moore-Penrose pseudoinverse of $[J]$. One inverse-like property between them is demonstrated by Eqs. 6 and 23: $[J] \mathbf{a} = \mathbf{y}^{\text{fit}}$ and $[J^\dagger] \mathbf{y} = \mathbf{a}$. One matrix undoes what the other does to the extent that $\mathbf{y} \approx \mathbf{y}^{\text{fit}}$.

Expectation values

Expectation values provide important information about linear regression results. If the data set $\{y_i\}$ could be sampled over and over again, what should be expected for the returned parameters' means, variances and correlations? How big a χ^2 should be expected?

Parameter expectation values

For example, taking the expectation value of both sides of Eq. 23

$$\langle \mathbf{a} \rangle = \langle [J^\dagger] \mathbf{y} \rangle \quad (24)$$

would be the start of a proof to demonstrate that the parameters will average to their predicted mean, $\langle a_j \rangle = \alpha_j$.

Expectation values are with respect to the random variables y_i . Because $[X]$, $[\sigma_y]$, $[J]$ and $[J^\dagger]$ (and their inverses or transposes) do not depend on the y_i , they can be factored out when evaluating the expectation value on the right side of Eq. 24. Furthermore, each y_i is from a distribution of mean μ_i . Thus, $\langle y_i \rangle = \mu_i$ and

$$\langle \mathbf{y} \rangle = \boldsymbol{\mu} \quad (25)$$

Factoring out the constant $[J^\dagger]$ from the right side of Eq. 24, then substituting Eqs. 22 and 25 for $[J^\dagger]$ and $\langle \mathbf{y} \rangle$, followed by substitutions using Eqs. 7 and 19 for $\boldsymbol{\mu}$ and $[X]$ gives

$$\begin{aligned} \langle \mathbf{a} \rangle &= [J^\dagger] \langle \mathbf{y} \rangle \\ &= [X^{-1}] [J^T] [\sigma_y^{-1}] \boldsymbol{\mu} \\ &= [X^{-1}] [J^T] [\sigma_y^{-1}] [J] \boldsymbol{\alpha} \\ &= [X^{-1}] [X] \boldsymbol{\alpha} \\ &= \boldsymbol{\alpha} \end{aligned} \quad (26)$$

where a cancellation of the identity matrix $[X^{-1}] [X]$ completes the proof.

Taking expectation values of Eq. 23 now demonstrates

$$\boldsymbol{\alpha} = [J^\dagger] \boldsymbol{\mu} \quad (27)$$

And finishing the cycle by substituting Eq. 7 for $\boldsymbol{\mu}$ illustrates another inverse-like property of the pseudoinverse $[J^\dagger]$:

$$\boldsymbol{\alpha} = [J^\dagger] [J] \boldsymbol{\alpha} \quad (28)$$

which implies that $[J^\dagger] [J]$ is an $M \times M$ identity matrix.

$$[J^\dagger] [J] = [I] \quad (29)$$

A solution for $[J^\dagger]$ satisfying only Eq. 29 is not unique. The pseudoinverse is unique because $[J^\dagger] \mathbf{y}$ gives the best-fit parameters \mathbf{a} that minimize the χ^2 . Equation 29 is a result, not a requirement.

Parameter variances and covariances

Determining the $M \times M$ covariance matrix for the fitting parameters starts with the defining equations for its elements

$$[\sigma_a]_{jk} = \langle (a_j - \alpha_j)(a_k - \alpha_k) \rangle \quad (30)$$

where $\langle a_j \rangle = \alpha_j$ has been used. By definition, these are just the elements of the expectation value of the matrix obtained from the outer product of the column vector $\mathbf{a} - \boldsymbol{\alpha}$ and its corresponding row vector $\mathbf{a}^T - \boldsymbol{\alpha}^T$.

$$[\sigma_a] = \langle (\mathbf{a} - \boldsymbol{\alpha}) (\mathbf{a}^T - \boldsymbol{\alpha}^T) \rangle \quad (31)$$

Substituting Eq. 23 and its transpose $\mathbf{a}^T = \mathbf{y}^T [J^{\dagger T}]$ in Eq. 31 gives

$$[\sigma_a] = \langle ([J^\dagger] \mathbf{y} - \boldsymbol{\alpha}) (\mathbf{y}^T [J^{\dagger T}] - \boldsymbol{\alpha}^T) \rangle \quad (32)$$

Expanding the product and factoring constant matrices from expectation values where appropriate

$$[\sigma_a] = [J^\dagger] \langle \mathbf{y} \mathbf{y}^T \rangle [J^{\dagger T}] - [J^\dagger] \langle \mathbf{y} \rangle \boldsymbol{\alpha}^T - \boldsymbol{\alpha} \langle \mathbf{y}^T \rangle [J^{\dagger T}] + \boldsymbol{\alpha} \boldsymbol{\alpha}^T \quad (33)$$

Substituting Eq. 25 and its transpose $\langle \mathbf{y}^T \rangle = \boldsymbol{\mu}^T$ in the middle two terms above and then substituting Eq. 27 and its transpose $\boldsymbol{\mu}^T [J^{\dagger T}] = \boldsymbol{\alpha}^T$ in these terms leaves

$$\begin{aligned} [\sigma_a] &= [J^\dagger] \langle \mathbf{y} \mathbf{y}^T \rangle [J^{\dagger T}] - [J^\dagger] \boldsymbol{\mu} \boldsymbol{\alpha}^T - \boldsymbol{\alpha} \boldsymbol{\mu}^T [J^{\dagger T}] + \boldsymbol{\alpha} \boldsymbol{\alpha}^T \\ &= [J^\dagger] \langle \mathbf{y} \mathbf{y}^T \rangle [J^{\dagger T}] - \boldsymbol{\alpha} \boldsymbol{\alpha}^T - \boldsymbol{\alpha} \boldsymbol{\alpha}^T + \boldsymbol{\alpha} \boldsymbol{\alpha}^T \\ &= [J^\dagger] \langle \mathbf{y} \mathbf{y}^T \rangle [J^{\dagger T}] - \boldsymbol{\alpha} \boldsymbol{\alpha}^T \end{aligned} \quad (34)$$

To take the expectation value in the first term on the right note that the matrix $\mathbf{y} \mathbf{y}^T$ is an $N \times N$ symmetric matrix with elements $[\mathbf{y} \mathbf{y}^T]_{ij} = y_i y_j$. Further noting that $\langle y_i y_j \rangle = \mu_i \mu_j + \sigma_i^2 \delta_{ij}$ gives

$$\langle \mathbf{y} \mathbf{y}^T \rangle = \boldsymbol{\mu} \boldsymbol{\mu}^T + [\sigma_y] \quad (35)$$

and Eq. 34 becomes

$$[\sigma_a] = [J^\dagger] \boldsymbol{\mu} \boldsymbol{\mu}^T [J^{\dagger T}] + [J^\dagger] [\sigma_y] [J^{\dagger T}] - \boldsymbol{\alpha} \boldsymbol{\alpha}^T \quad (36)$$

Using Eq. 27 and its transpose $\boldsymbol{\mu}^T [J^{\dagger T}] = \boldsymbol{\alpha}^T$ shows that the first and last term cancel giving

$$[\sigma_a] = [J^\dagger] [\sigma_y] [J^{\dagger T}] \quad (37)$$

Substituting Eq. 22 and its transpose $[J^{\dagger T}] = [\sigma_y^{-1}] [J] [X^{-1}]$ gives

$$\begin{aligned} [\sigma_a] &= [X^{-1}] [J^T] [\sigma_y^{-1}] [\sigma_y] [\sigma_y^{-1}] [J] [X^{-1}] \\ &= [X^{-1}] [J^T] [\sigma_y^{-1}] [J] [X^{-1}] \\ &= [X^{-1}] [X] [X^{-1}] \\ &= [X^{-1}] \end{aligned} \quad (38)$$

Equation 38 gives the recipe for the parameter variances and covariances. Find the inverse of the $M \times M$ matrix $[X]$ whose two ingredients are the same as before: the regression functions $[J]$ and the data uncertainties $[\sigma_y]$. It does not depend on \mathbf{y} . Consequently, to the extent that the regression functions and the data uncertainties can be predetermined, so too can the parameter uncertainties. This fact can be useful in experimental design.

The expectation value of χ^2

Distinct from the “best-fit” χ^2 given by Eq. 8 and based on $\{y_i^{\text{fit}}\}$, a “true” χ^2 can be defined based on the true means $\{\mu_i\}$ and expressed

$$\chi^2 = \sum_{i=1}^N \frac{(y_i - \mu_i)^2}{\sigma_i^2} \quad (39)$$

It is easy to show that the expectation values of this χ^2 evaluates to N . By the definition of sample variance, $\sigma_i^2 = \langle (y_i - \mu_i)^2 \rangle$, each of the N terms in Eq. 39 has an expectation value of unity and thus the sum evaluates to N .

Because the true parameters and true means are not generally known, the true χ^2 cannot generally be calculated. Exceptions occur in simulations such as those in the spreadsheet that accompanies this paper. Obviously, once the best-fit parameters are determined, $\{y_i^{\text{fit}}\}$ and the best-fit χ^2 of Eq. 8 can then be calculated. However, the best-fit χ^2 is clearly a different statistic from the true χ^2 . Most importantly, the best-fit parameters are chosen to minimize the chi-square for a particular data set. Using any other parameter set—including the true parameters—can only make the chi-square increase. Consequently, the best-fit chi-square will **always** be smaller than the true chi-square for that particular data set. How much smaller (on average) is the question to be answered next.

Before evaluating the expectation value of the best-fit χ^2 , let’s again evaluate the expectation value of the true χ^2 —this time proceeding from the matrix algebra form of Eq. 39

$$\chi^2 = (\mathbf{y}^T - \boldsymbol{\mu}^T) [\sigma_y^{-1}] (\mathbf{y} - \boldsymbol{\mu}) \quad (40)$$

Of course, it will give the known result $\langle \chi^2 \rangle = N$, but the process will demonstrate techniques that will be used again for evaluating the expectation value of the best-fit χ^2 .

The evaluation proceeds as follows:

$$\langle \chi^2 \rangle = \langle (\mathbf{y}^T - \boldsymbol{\mu}^T) [\sigma_y^{-1}] (\mathbf{y} - \boldsymbol{\mu}) \rangle \quad (41)$$

Expanding the product, factoring out constants, and using Eq. 25 and its transpose gives

$$\begin{aligned} \langle \chi^2 \rangle &= \langle \mathbf{y}^T [\sigma_y^{-1}] \mathbf{y} \rangle - \langle \mathbf{y}^T \rangle [\sigma_y^{-1}] \boldsymbol{\mu} - \boldsymbol{\mu}^T [\sigma_y^{-1}] \langle \mathbf{y} \rangle + \boldsymbol{\mu}^T [\sigma_y^{-1}] \boldsymbol{\mu} \\ &= \langle \mathbf{y}^T [\sigma_y^{-1}] \mathbf{y} \rangle - \boldsymbol{\mu}^T [\sigma_y^{-1}] \boldsymbol{\mu} - \boldsymbol{\mu}^T [\sigma_y^{-1}] \boldsymbol{\mu} + \boldsymbol{\mu}^T [\sigma_y^{-1}] \boldsymbol{\mu} \\ &= \langle \mathbf{y}^T [\sigma_y^{-1}] \mathbf{y} \rangle - \boldsymbol{\mu}^T [\sigma_y^{-1}] \boldsymbol{\mu} \end{aligned} \quad (42)$$

The expectation value of an N by N matrix $[Q]$ in the form $\langle \mathbf{y}^T [Q] \mathbf{y} \rangle$ can be simplified by examining the implicit multiplications of such a vector-matrix-vector product.

$$\mathbf{y}^T [Q] \mathbf{y} = \sum_{i=1}^N \sum_{j=1}^N y_i [Q]_{ij} y_j \quad (43)$$

Noting that $\langle y_i y_j \rangle = \mu_i \mu_j + \sigma_i^2 \delta_{ij}$ gives

$$\begin{aligned} \langle \mathbf{y}^T [Q] \mathbf{y} \rangle &= \sum_{i=1}^N \sum_{j=1}^N \langle y_i [Q]_{ij} y_j \rangle \\ &= \sum_{i=1}^N \sum_{j=1}^N [Q]_{ij} \langle y_i y_j \rangle \\ &= \sum_{i=1}^N \sum_{j=1}^N [Q]_{ij} (\mu_i \mu_j + \sigma_i^2 \delta_{ij}) \\ &= \sum_{i=1}^N \sum_{j=1}^N \mu_i [Q]_{ij} \mu_j + \sum_{i=1}^N [Q]_{ii} \sigma_i^2 \end{aligned} \quad (44)$$

or

$$\langle \mathbf{y}^T [Q] \mathbf{y} \rangle = \boldsymbol{\mu}^T [Q] \boldsymbol{\mu} + \text{Tr}([Q] [\sigma_y]) \quad (45)$$

Where the trace of any $N \times N$ square matrix $[B]$ is the sum of its diagonal elements

$$\text{Tr}([B]) = \sum_{i=1}^N [B]_{ii} \quad (46)$$

Applying Eq. 45 to the first term on the right side of Eq. 42 gives

$$\begin{aligned} \langle \mathbf{y}^T [\sigma_y^{-1}] \mathbf{y} \rangle &= \boldsymbol{\mu}^T [\sigma_y^{-1}] \boldsymbol{\mu} + \text{Tr}[\sigma_y^{-1}] [\sigma_y] \\ &= \boldsymbol{\mu}^T [\sigma_y^{-1}] \boldsymbol{\mu} + \text{Tr}[I] \\ &= \boldsymbol{\mu}^T [\sigma_y^{-1}] \boldsymbol{\mu} + N \end{aligned} \quad (47)$$

where the identity matrix $[I]$ above is of the $N \times N$ variety as it is formed by the product of the $N \times N$ matrix $[\sigma_y]$ and its inverse. In the last step, the trace of this identity matrix is taken.

Substituting Eq. 47 into Eq. 42 then gives the expected result

$$\langle \chi^2 \rangle = N \quad (48)$$

Finding the expectation value of the best-fit χ^2 will be based on Eq. 9 and begins

$$\langle \chi^2 \rangle = \langle (\mathbf{y}^T - \mathbf{y}^{\text{fit}T}) [\sigma_y^{-1}] (\mathbf{y} - \mathbf{y}^{\text{fit}}) \rangle \quad (49)$$

The \mathbf{y}^{fit} are replaced by Eq. 6 and its transpose $\mathbf{y}^{\text{fit}T} = \mathbf{a}^T [J^T]$ giving

$$\langle \chi^2 \rangle = \langle (\mathbf{y}^T - \mathbf{a}^T [J^T]) [\sigma_y^{-1}] (\mathbf{y} - [J] \mathbf{a}) \rangle \quad (50)$$

Substituting Eq. 23 and its transpose $\mathbf{a}^T = \mathbf{y}^T [J^{\dagger T}]$ then gives

$$\begin{aligned} \langle \chi^2 \rangle &= \langle (\mathbf{y}^T - \mathbf{y}^T [J^{\dagger T}] [J^T]) [\sigma_y^{-1}] (\mathbf{y} - [J] [J^\dagger] \mathbf{y}) \rangle \\ &= \langle \mathbf{y}^T (1 - [J^{\dagger T}] [J^T]) [\sigma_y^{-1}] (1 - [J] [J^\dagger]) \mathbf{y} \rangle \end{aligned} \quad (51)$$

Expanding the product gives

$$\begin{aligned} \langle \chi^2 \rangle &= \\ &\langle \mathbf{y}^T ([\sigma_y^{-1}] - [\sigma_y^{-1}] [J] [J^\dagger] - [J^{\dagger T}] [J^T] [\sigma_y^{-1}] + [J^{\dagger T}] [J^T] [\sigma_y^{-1}] [J] [J^\dagger]) \mathbf{y} \rangle \end{aligned} \quad (52)$$

Substituting Eq. 22 and its transpose $[J^{\dagger T}] = [\sigma_y^{-1}] [J] [X^{-1}]$ for $[J^\dagger]$ and its transpose in the last three matrix products, and then in the final term substituting $[X]$ for the combination $[J^T] [\sigma_y^{-1}] [J]$ (Eq. 19) and then canceling the product $[X] [X^{-1}]$ shows that, aside from their signs, the last three matrix products are all the same and given by $[\sigma_y^{-1}] [J] [X^{-1}] [J^T] [\sigma_y^{-1}]$.

Taking account of the signs and using the form of the first matrix product, the resulting expression is thus

$$\langle \chi^2 \rangle = \langle \mathbf{y}^T [\sigma_y^{-1}] \mathbf{y} \rangle - \langle \mathbf{y}^T [\sigma_y^{-1}] [J] [J^\dagger] \mathbf{y} \rangle \quad (53)$$

Applying Eq. 45 to the first term in Eq. 53 gives

$$\begin{aligned} \langle \mathbf{y}^T [\sigma_y^{-1}] \mathbf{y} \rangle &= \boldsymbol{\mu}^T [\sigma_y^{-1}] \boldsymbol{\mu} + \text{Tr} [\sigma_y^{-1}] [\sigma_y] \\ &= \boldsymbol{\mu}^T [\sigma_y^{-1}] \boldsymbol{\mu} + \text{Tr} [I] \\ &= \boldsymbol{\mu}^T [\sigma_y^{-1}] \boldsymbol{\mu} + N \end{aligned} \quad (54)$$

Applying Eq. 45 to the second term in Eq. 53 gives

$$\langle \mathbf{y}^T [\sigma_y^{-1}] [J] [J^\dagger] \mathbf{y} \rangle = \boldsymbol{\mu}^T [\sigma_y^{-1}] [J] [J^\dagger] \boldsymbol{\mu} + \text{Tr} [\sigma_y^{-1}] [J] [J^\dagger] [\sigma_y] \quad (55)$$

Substituting Eq. 27 and then Eq 7 provides the simplification of the first term on the right.

$$\begin{aligned} \boldsymbol{\mu}^T [\sigma_y^{-1}] [J] [J^\dagger] \boldsymbol{\mu} &= \boldsymbol{\mu}^T [\sigma_y^{-1}] [J] \boldsymbol{\alpha} \\ &= \boldsymbol{\mu} [\sigma_y^{-1}] \boldsymbol{\mu} \end{aligned} \quad (56)$$

Simplification of the second term in Eq. 55 begins by substituting Eq. 22 for $[J^\dagger]$ and then proceeds as follows

$$\text{Tr} [\sigma_y^{-1}] [J] [J^\dagger] [\sigma_y] = \text{Tr} [\sigma_y^{-1}] [J] [X^{-1}] [J^T] [\sigma_y^{-1}] [\sigma_y]$$

$$\begin{aligned}
&= \text{Tr} [\sigma_y^{-1}] [J] [X^{-1}] [J^T] \\
&= \text{Tr} [J^{\dagger T}] [J^T] \\
&= \text{Tr} [J] [J^\dagger] \\
&= \text{Tr} [J^\dagger] [J] \\
&= \text{Tr} [I] \\
&= M
\end{aligned} \tag{57}$$

In the second line $[\sigma_y^{-1}] [\sigma_y] = [I]$ is used and in the third line the transpose of Eq. 22 is used. In the fourth line the equality between the trace of a square matrix ($[J] [J^\dagger]$) and its transpose ($[J^{\dagger T}] [J^T]$) is used. But note that $[J] [J^\dagger]$ appearing in line four is an $N \times N$ matrix whereas $[J^\dagger] [J]$ appearing in line five is an $M \times M$ matrix. Consequently, the two matrices cannot be equal. However, their traces are equal and only that fact is used in line five. The substitution is valid because for any two matrices $[A]$ of shape $N \times M$ and $[B]$ of shape $M \times N$, both $[A] [B]$ and $[B] [A]$ exist and the rules for multiplication and trace give

$$\begin{aligned}
\text{Tr} [A] [B] &= \sum_{i=1}^N \sum_{j=1}^M [A]_{ij} [B]_{ji} \\
&= \sum_{j=1}^M \sum_{i=1}^N [B]_{ji} [A]_{ij} \\
&= \text{Tr} [B] [A]
\end{aligned} \tag{58}$$

In the sixth line of Eq. 57, Eq. 29 giving $[J^\dagger] [J]$ as the $M \times M$ identity matrix is used and in line seven its trace is taken. Finally, using Eq. 56 and Eq. 57 in Eq. 55 and combining this with Eq. 54 in Eq. 53 then gives

$$\langle \chi^2 \rangle = N - M \tag{59}$$

Note the familiar result that the expectation value of χ^2 is its number of degrees of freedom. The true χ^2 based on Eq. 39 has no adjustable parameters and the number of degrees of freedom is the number of data points. The best-fit χ^2 , on the other hand, is minimized by adjusting the M fitting parameters \mathbf{a} , thereby reducing the number of degrees of freedom by M . Note that the best-fit χ^2 is not just smaller than the true χ^2 on average. It is, in fact, always smaller. It is not always smaller by the amount M —the reduction can be smaller or larger than this. The reduction is, in fact, a χ^2 random variable with M degrees of freedom.

Relationship to non-linear regression

Fitting data to functions that are nonlinear in the fitting parameters can be complicated and problematic. There are many software packages that do the job well and no attempt will be

made here to describe the many algorithms that can be used. The intent of this section is to give the conditions under which the linear regression results are still applicable and to show the required modifications.

As an example, consider a nonlinear function $F(\{\alpha_j\}, x)$ of a single independent variable x and a set of M fitting parameters $\{\alpha_j\}$. Taking the $\{\alpha_j\}$ as the true parameters, the function evaluated at the data set's values for x_i then gives the true means μ_i from which the data set's N measured y_i values are samples.

$$\mu_i = F(\{\alpha_j\}, x_i) \quad (60)$$

As with linear regression, the best-fit parameters will be denoted $\{a_j\}$ and are defined as those values that minimize the χ^2 . The $\{a_j\}$ then determine the best-fit y -values.

$$y_i^{\text{fit}} = F(\{a_j\}, x_i) \quad (61)$$

Most linear regression results, with modifications to be discussed shortly, will be sound as long as the fitting function is well approximated by a first-order Taylor series about any set of parameters $\{\beta_j\}$ in the vicinity of the best-fit values $\{a_j\}$ —say, within a range of a few parameter uncertainties σ_j for each parameter. The $\{\beta_j\}$ should not be considered crude initial guesses. Rather, they should be considered “almost there” values—say, the next-to-last parameter set in the typical iterative approach to finding the best-fit $\{a_j\}$. Then $F(\{\beta_j\}, x_i)$ would be the fitted y -values associated with that almost-there solution.

The first-order Taylor series expansion that the fitting function would need to closely follow would be written

$$y_i^{\text{fit}}(\{a_j\}, x_i) = F(\{\beta_j\}, x_i) + \sum_{j=1}^M \frac{\partial F(\{\beta_j\}, x_i)}{\partial \beta_j} (a_j - \beta_j) \quad (62)$$

First, we define

$$\Delta y_i = y_i - F(\{\beta_j\}, x_i) \quad (63)$$

Δy_i will play the role of a modified measured y -value in a linear regression model for a nonlinear fit. This modified data is the deviation of the raw data y_i from the almost-there fit values given by $F(\{\beta_j\}, x_i)$. Next, we define

$$\Delta a_j = a_j - \beta_j \quad (64)$$

Δa_j will play the role of a modified fitting parameter and is the deviation of each a_j from its almost-there value β_j . Finally, we define the $N \times M$ set of coefficients

$$[J]_{ij} = \frac{\partial F(\{\beta_j\}, x_i)}{\partial \beta_j} \quad (65)$$

Considered as an $N \times M$ matrix of values, each column of the matrix will play the role of a modified fitting function $f_j(x_i)$ and the whole matrix $[J]$ will play the same role it played in the linear regression formulation.

Note how this formulation is exact in the case that $F(\{\beta_j\}, x)$ is linear in the fitting parameters. Consequently, all results in the following sections will also apply to a linear regression.

With these three definitions, the χ^2 —now to be considered a function of the fitting parameters $\{a_j\}$ —becomes

$$\begin{aligned}
 \chi^2 &= \sum_{i=1}^N \frac{(y_i - y^{\text{fit}}(\{a_j\}, x_i))^2}{\sigma_i^2} \\
 &= \sum_{i=1}^N \frac{(y_i - F(\{\beta_j\}, x_i) - \sum_{j=1}^M \frac{dF(\{\beta_j\}, x_i)}{d\beta_j} (a_j - \beta_j))^2}{\sigma_i^2} \\
 &= \sum_{i=1}^N \frac{(\Delta y_i - \sum_{j=1}^M [J]_{ij} \Delta a_j)^2}{\sigma_i^2}
 \end{aligned} \tag{66}$$

where Eq. 62 was substituted for $y^{\text{fit}}(\{a_j\}, x_i)$ in the second line and Eqs. 63-65 were substituted in the third line.

Equation 66 can be rewritten

$$\chi^2 = (\Delta \mathbf{y}^T - \Delta \mathbf{a}^T [J^T]) [\sigma_y^{-1}] (\Delta \mathbf{y} - [J] \Delta \mathbf{a}) \tag{67}$$

where $\Delta \mathbf{y}$ is the column vector (of length N) representing the Δy_i and $\Delta \mathbf{a}$ is the column vector (of length M) representing the Δa_j . $\Delta \mathbf{y}^T$ and $\Delta \mathbf{a}^T$ and $[J^T]$ are the transposed quantities.

Equation 63 shows that the elements of $\Delta \mathbf{y}$ are simply offset from the elements of \mathbf{y} describing the raw, measured y_i . Consequently, variances and covariances of the Δy_i are the same as those of the y_i and, consequently, $[\sigma_y]$ and its inverse $[\sigma_y^{-1}]$ do not need modification.

The χ^2 of Eq. 67 is now in a form analogous to Eq. 11 and the values for the Δa_j that minimize it can be determined by following the linear regression formulation already presented, i.e., by taking derivatives with respect to each Δa_j and setting them all equal to zero at the minimum. We need not go through that formulation again and simply use the results. The solution for the best-fit Δa_j analogous to Eq. 23 becomes

$$\Delta \mathbf{a} = [J^\dagger] \Delta \mathbf{y} \tag{68}$$

where the pseudoinverse $[J^\dagger]$ is still given by Eq. 22 with the $[X]$ matrix still given by Eq. 19, but with the elements of $[J]$ now given by Eq. 65.

What was a direct solution in the linear regression treatment has become a solution involving the changes from the almost-there solution in the nonlinear case. The solution gives the changes $\Delta \mathbf{a}$ (from the almost-there parameter values) in terms of raw data deviations $\Delta \mathbf{y}$ (from the almost-there fit values).

Eq. 64 is then used to find the best-fit $\{a_j\}$ from the solution values $\{\Delta a_j\}$ and the starting parameter set $\{\beta_j\}$ from which the solution is defined.

$$a_j = \beta_j + \Delta a_j \tag{69}$$

An important result implied from the linear regression treatment is that $[X^{-1}]$ is the covariance matrix for the $\{\Delta a_j\}$. The simple constant offset relation between a_j and Δa_j expressed by Eq. 69 implies that the set $\{a_j\}$ and $\{\Delta a_j\}$ will have the same variances and covariances and thus, $[X^{-1}]$ is also the covariance matrix for the parameters $\{a_j\}$ themselves.

Note that any set of parameters $\{\beta_j\}$ sufficiently close to the best-fit values $\{a_j\}$ (along with all the derivatives in the matrix $[J]$ evaluated at $\{\beta_j\}$) contain all the information necessary to find the $\{a_j\}$ and its covariance matrix. Further note that the best-fit $\{a_j\}$ would also be such a set. Of course, this set must give $\Delta \mathbf{a} = 0$, but it would still give non-zero values for $\{\Delta y_i\}$ and a well-defined, non-zero covariance matrix $[X^{-1}]$. Consequently, the difference between a covariance matrix evaluated from an almost-there parameter set $\{\beta_j\}$ and one evaluated from the best-fit parameter set $\{a_j\}$ should be small. Nonetheless, $[X^{-1}]$ is normally reevaluated one last time based on the derivative matrix $[J]$ evaluated at $\{a_j\}$ rather than $\{\beta_j\}$.

One might try iterating Eqs. 68 and 69 (with the ending values of $\{a_j\}$ used as the starting values $\{\beta_j\}$ for the next iteration) as an algorithm to find the best-fit parameters. However, this algorithm has known problems if the starting parameters are not close enough to the best-fit $\{a_j\}$. More sophisticated nonlinear fitting routines use algorithms better suited for reliably finding a solution from more distant initial guesses. Nonetheless, this approach does work for limited applications and is demonstrated in the Excel spreadsheet called *Matrix Algebra and Regression* on the lab web site.

Dependence of χ^2 on the fitting parameters

It is worthwhile to understand how χ^2 would vary about its minimum as the parameters vary about their best-fit values $\{a_j\}$. The χ^2 of Eq. 67 will be the starting point, but it is to be understood now to use $\{a_j\}$ as the almost-there parameter set and the Δa_j are to be interpreted as giving the deviations from that set. We will not try to prove again that $\Delta \mathbf{a} = 0$ at the χ^2 minimum, but rather ask for the form of χ^2 away from the minimum as a function of $\Delta \mathbf{a}$.

Expanding the right side of Eq. 67 gives

$$\begin{aligned} \chi^2 = & \Delta \mathbf{y}^T [\sigma_y^{-1}] \Delta \mathbf{y} \\ & - \Delta \mathbf{y}^T [\sigma_y^{-1}] [J] \Delta \mathbf{a} - \Delta \mathbf{a}^T [J^T] [\sigma_y^{-1}] \Delta \mathbf{y} + \Delta \mathbf{a}^T [J^T] [\sigma_y^{-1}] [J] \Delta \mathbf{a} \end{aligned} \quad (70)$$

Since all deviations are now from the best-fit values, $\Delta y_i = y_i - y_i^{\text{fit}}$ and the first term is just the best-fit χ^2 , which we now call χ_0^2 . The equivalent of Eq. 20 for the nonlinear regression formulation is

$$[X] \Delta \mathbf{a} = [J^T] [\sigma_y^{-1}] \Delta \mathbf{y} \quad (71)$$

Using this and its transpose to eliminate the $\Delta \mathbf{y}$ and its transpose in the second and third term and substituting Eq. 19 in the fourth shows that, aside from their signs, the final three

terms are equivalent and evaluate to $\Delta \mathbf{a}^T [X] \Delta \mathbf{a}$. After subtracting χ_0^2 from both sides, the final result is

$$\Delta \chi^2 = \Delta \mathbf{a}^T [X] \Delta \mathbf{a} \quad (72)$$

where $\Delta \chi^2 = \chi^2 - \chi_0^2$ is the deviation of χ^2 from its best-fit value.

$\Delta \chi^2 = 1$ Rule

Recall that $[X^{-1}]$ is the covariance matrix with diagonal elements giving the parameter variances σ_j^2 and off-diagonal elements giving the parameter covariances σ_{jk} . Moreover $[X]$ and $[X^{-1}]$ are inverses of one another and both have a self-transpose symmetry, that is, $[X]_{kj} = [X]_{jk}$ and $[X^{-1}]_{kj} = [X^{-1}]_{jk}$.

In the very rare case where all parameter covariances are zero, the covariance matrix $[X^{-1}]$ is diagonal (zero for all off-diagonal elements). Consequently, its inverse $[X]$ is also diagonal with elements $[X]_{jj} = 1/\sigma_j^2$ and the χ^2 variation given by Eq. 72 becomes

$$\Delta \chi^2 = \sum_{j=1}^M \frac{(\Delta a_j)^2}{\sigma_j^2} \quad (73)$$

This equation gives $\Delta \chi^2 = 1$ when any single parameter is moved one standard deviation σ_k from its best-fit value. Changing any other parameters independently increases the χ^2 in a like manner—quadratically, in $\Delta a_j/\sigma_j$.

While the independent increases in χ^2 given above for uncorrelated parameters is consistent with the $\Delta \chi^2 = 1$ rule, the rule requires a caveat for the general—and vastly more likely—case of correlated parameters in which case a χ^2 increase from one parameter change can be partially canceled by appropriate changes in the other parameters. Offsetting any parameter one standard deviation from its best-fit value (keeping the others at their best-fit values) generally gives a $\Delta \chi^2 \geq 1$. However, the $\Delta \chi^2$ value always decreases back to one if the other parameters are then re-optimized for the smallest possible χ^2 holding the offset parameter fixed one-sigma from its best-fit value. This second χ^2 minimization is the caveat to the $\Delta \chi^2 = 1$ rule.

To see how the rule arises,¹ the parameter vector $\Delta \mathbf{a}$ is partitioned into a vector $\Delta \mathbf{a}_1$ of length K and another vector $\Delta \mathbf{a}_2$ of length $L = M - K$.

$$\Delta \mathbf{a} = \begin{pmatrix} \Delta \mathbf{a}_1 \\ \Delta \mathbf{a}_2 \end{pmatrix} \quad (74)$$

¹This derivation is from R. A. Arndt and M. H. MacGregor, Nucleon-Nucleon Phase Shift Analysis by Chi-Squared Minimization in *Methods of Mathematical Physics*, vol. 6, pp. 253-296, Academic Press, New York (1966).

The parameters are ordered so those in the second group will be offset from their best-fit values to see how the χ^2 changes if the parameters in the first group are automatically re-optimized for any variations in that second group. At the end of the derivation, the second group will be taken to consist of a single parameter to prove the $\Delta\chi^2 = 1$ rule.

The matrix $[X]$ is then partitioned into an $K \times K$ matrix A , an $L \times L$ matrix C , a $K \times L$ matrix B and its $L \times K$ transpose B^T

$$[X] = \left(\begin{array}{c|c} [A] & [B] \\ \hline [B^T] & [C] \end{array} \right) \quad (75)$$

Note that the self-transpose symmetry of $[X]$ implies that same symmetry for the square matrices $[A]$ and $[C]$ and is why the two off-diagonal matrices $[B]$ and $[B^T]$ are transposes of one another.

Equation 72 then becomes

$$\begin{aligned} \Delta\chi^2 &= \Delta\mathbf{a}^T [X] \Delta\mathbf{a} \\ &= \left(\Delta\mathbf{a}_1^T \mid \Delta\mathbf{a}_2^T \right) \left(\begin{array}{c|c} [A] & [B] \\ \hline [B^T] & [C] \end{array} \right) \begin{pmatrix} \Delta\mathbf{a}_1 \\ \Delta\mathbf{a}_2 \end{pmatrix} \\ &= \left(\Delta\mathbf{a}_1^T \mid \Delta\mathbf{a}_2^T \right) \left(\begin{array}{c} [A] \Delta\mathbf{a}_1 + [B] \Delta\mathbf{a}_2 \\ [B^T] \Delta\mathbf{a}_1 + [C] \Delta\mathbf{a}_2 \end{array} \right) \\ &= \Delta\mathbf{a}_1^T [A] \Delta\mathbf{a}_1 + \Delta\mathbf{a}_1^T [B] \Delta\mathbf{a}_2 + \Delta\mathbf{a}_2^T [B^T] \Delta\mathbf{a}_1 + \Delta\mathbf{a}_2^T [C] \Delta\mathbf{a}_2 \end{aligned} \quad (76)$$

Now we consider $\Delta\mathbf{a}_2$ to be fixed and adjust $\Delta\mathbf{a}_1$ to minimize $\Delta\chi^2$ for this fixed $\Delta\mathbf{a}_2$. The minimization proceeds as usual; finding this value of $\Delta\mathbf{a}_1$ requires solving the set of K equations obtained by setting to zero the derivative of $\Delta\chi^2$ with respect to each Δa_j , $j = 1 \dots K$ in $\Delta\mathbf{a}_1$. The equation for index j becomes

$$\begin{aligned} 0 &= \frac{d\Delta\chi^2}{d\Delta a_j} \\ &= \boldsymbol{\delta}_j^T [A] \Delta\mathbf{a}_1 + \Delta\mathbf{a}_1^T [A] \boldsymbol{\delta}_j + \boldsymbol{\delta}_j^T [B] \Delta\mathbf{a}_2 + \Delta\mathbf{a}_2^T [B^T] \boldsymbol{\delta}_j \end{aligned} \quad (77)$$

where $\boldsymbol{\delta}_j = d\Delta\mathbf{a}_1/d\Delta a_j$ is the unit column-vector with a single 1 in the j th row, and $\boldsymbol{\delta}_j^T$ is its transposed row vector. All terms are scalars and the first two and second two are identical—the expressions are simply transposes of one another. Dividing both sides by two, Eq. 77 can then be taken to be

$$0 = \boldsymbol{\delta}_j^T [A] \Delta\mathbf{a}_1 + \boldsymbol{\delta}_j^T [B] \Delta\mathbf{a}_2 \quad (78)$$

and the full set of all K of these equations is the vector equation

$$0 = [A] \Delta \mathbf{a}_1 + [B] \Delta \mathbf{a}_2 \quad (79)$$

which has the solution

$$\Delta \mathbf{a}_1 = -[A^{-1}] [B] \Delta \mathbf{a}_2 \quad (80)$$

where $[A^{-1}]$ is the inverse of $[A]$ and is also a self-transpose matrix. Inserting this equation and its transpose in Eq. 76 then gives

$$\begin{aligned} \Delta \chi^2 = & \Delta \mathbf{a}_2^T [B^T] [A^{-1}] [A] [A^{-1}] [B] \Delta \mathbf{a}_2 - \Delta \mathbf{a}_2^T [B^T] [A^{-1}] [B] \Delta \mathbf{a}_2 \\ & - \Delta \mathbf{a}_2^T [B^T] [A^{-1}] [B] \Delta \mathbf{a}_2 + \Delta \mathbf{a}_2^T [C] \Delta \mathbf{a}_2 \end{aligned} \quad (81)$$

After canceling an $[A^{-1}] [A] = [I]$ in the first term, the first three terms are all identical (aside from their signs) and the equation can be rewritten

$$\Delta \chi^2 = \Delta \mathbf{a}_2^T \left([C] - [B^T] [A^{-1}] [B] \right) \Delta \mathbf{a}_2 \quad (82)$$

Note that this equation gives $\Delta \chi^2$ as a function of the parameters in a subset of its parameter space without specifying the values for the other parameters. Whenever this is done the re-optimization of all unspecified parameters for a minimum χ^2 is implied.

Now we must consider the inverse matrix $[X^{-1}]$. It can be put in the same submatrix form

$$[X^{-1}] = \left(\begin{array}{c|c} [A'] & [B'] \\ \hline [B'^T] & [C'] \end{array} \right) \quad (83)$$

The equation $[X^{-1}] [X] = [I]$ then gives

$$\begin{aligned} [I] &= \left(\begin{array}{c|c} [A'] & [B'] \\ \hline [B'^T] & [C'] \end{array} \right) \left(\begin{array}{c|c} [A] & [B] \\ \hline [B^T] & [C] \end{array} \right) \\ &= \left(\begin{array}{c|c} [A'] [A] + [B'] [B^T] & [A'] [B] + [B'] [C] \\ \hline [B'^T] [A] + [C'] [B^T] & [B'^T] [B] + [C'] [C] \end{array} \right) \end{aligned} \quad (84)$$

Consequently the $K \times K$ and $L \times L$ matrices in the top left and lower right must be identity matrices and the $L \times K$ and $K \times L$ matrices in the upper right and lower left must be identically zero. That is,

$$[A'] [A] + [B'] [B^T] = [I] \quad (85)$$

$$[A'] [B] = -[B'] [C] \quad (86)$$

$$[B'^T] [A] = - [C'] [B^T] \quad (87)$$

$$[B'^T] [B] + [C'] [C] = [I] \quad (88)$$

Multiplying Eq. 87 by $[A^{-1}]$ on the right gives

$$[B'^T] = - [C'] [B^T] [A^{-1}] \quad (89)$$

or

$$[B'^T] [B] = - [C'] [B^T] [A^{-1}] [B] \quad (90)$$

which, when inserted into Eq. 88 gives

$$[C'] ([C] - [B^T] [A^{-1}] [B]) = [I] \quad (91)$$

Multiplying both sides of Eq. 91 by $[C'^{-1}]$ (the inverse of the $L \times L$ covariance submatrix $[C']$) on the left gives

$$([C] - [B^T] [A^{-1}] [B]) = [C'^{-1}] \quad (92)$$

which can be substituted into Eq. 82 to get

$$\Delta\chi^2 = \Delta\mathbf{a}_2^T [C'^{-1}] \Delta\mathbf{a}_2 \quad (93)$$

The transpose of Eq. 89 is

$$[B'] = - [A^{-1}] [B] [C'] \quad (94)$$

(since $[A^{-1}]$ and $[C']$ are self-transpose matrices). Multiplying on the right by $[C'^{-1}]$ gives

$$[B'] [C'^{-1}] = - [A^{-1}] [B] \quad (95)$$

Substituting this in Eq. 80 then gives

$$\Delta\mathbf{a}_1 = [B'] [C'^{-1}] \Delta\mathbf{a}_2 \quad (96)$$

Now we can consider the case that $\Delta\mathbf{a}_2$ consists of a single parameter, say the very last one, Δa_M . The rest are in $\Delta\mathbf{a}_1$ and consist of the set $\{\Delta a_j\}$, $j = 1 \dots M-1$. From Eq. 83, $[C']$ would then be a 1×1 submatrix of the covariance matrix and would be given by $[C'] = \sigma_M^2$ and thus $[C'^{-1}] = 1/\sigma_M^2$. Solving Eq. 93 for σ_M^2 then gives

$$\sigma_M^2 = \frac{(\Delta a_M)^2}{\Delta\chi^2} \quad (97)$$

Furthermore, the submatrix $[B']$ would be the last column of the covariance matrix (less the last element, σ_M^2). Thus, it is a $1 \times (M-1)$ matrix whose elements are the covariances of the

other variables with a_M ; $[B']_{1j} = \sigma_{jM}$. Each element of Eq. 96 is then $\Delta a_j = \sigma_{jM}\Delta a_M/\sigma_M^2$. Eliminating σ_M using Eq. 97 and solving for σ_{jM} then gives

$$\sigma_{jM} = \frac{\Delta a_j \Delta a_M}{\Delta \chi^2} \quad (98)$$

If we further restrict $\Delta \chi^2 = 1$, these last two equations give $\Delta a_M = \pm \sigma_M$, and $\sigma_{jM} = \Delta a_j \Delta a_M$. This is the complete $\Delta \chi^2 = 1$ rule.

Equations 97 and 98 are a more general form of the $\Delta \chi^2 = 1$ rule that can be used with such programs as Microsoft Excel's Solver that perform minimizations but do not return the parameter variances or covariances. The procedure consists of first finding the best-fit parameters $\{a_j\}$ and the corresponding χ_0^2 via a minimization including all parameters. Next a parameter of interest a_M is changed by an arbitrary amount Δa_M (no larger than a few standard deviations σ_M) and the program is used again to minimize the χ^2 , but this time with the parameter of interest kept fixed, i.e., a_M is not included in the parameter list. The new χ^2 and the new parameter values $\{a'_j\}$ are recorded. $\Delta \chi^2 = \chi^2 - \chi_0^2$ and all parameter changes $\Delta a_j = a'_j - a_j$ are calculated and used with Eqs. 97 and 98 to find the parameter's variance σ_M^2 and its covariances σ_{jM} .

The $\Delta \chi^2 = 1$ rule answers the question: "What are the possible values for parameter a_k such that the χ^2 can still be kept within one of its minimum?" Taking the partitioned set $\{a_2\}$ to include only that one variable, Eq. 93 gave the answer: inside the interval $a_k \pm \sigma_k$. Taking the partitioned set $\{a_2\}$ to consist of two elements, that same equation can also be used to answer the question: "What are the possible pairs of values for parameters j and k such that the χ^2 can be kept within one of its minimum?" The answer turns out to be those values inside an ellipse in the space of a_j and a_k on which $\Delta \chi^2 = 1$.

To find this ellipse from Eq. 93, the inverse $[C'^{-1}]$ of a 2×2 submatrix of the covariance matrix will be required. The covariance submatrix for two parameters, say, a and b is expressed

$$[C'] = \begin{pmatrix} \sigma_a^2 & \sigma_{ab} \\ \sigma_{ab} & \sigma_b^2 \end{pmatrix} \quad (99)$$

and its inverse is given by

$$[C'^{-1}] = \frac{1}{\sigma_a^2 \sigma_b^2 - \sigma_{ab}^2} \begin{pmatrix} \sigma_b^2 & -\sigma_{ab} \\ -\sigma_{ab} & \sigma_a^2 \end{pmatrix} \quad (100)$$

Equation 93 then gives

$$\begin{aligned} \Delta \chi^2 &= \frac{1}{\sigma_a^2 \sigma_b^2 - \sigma_{ab}^2} \begin{pmatrix} \Delta a & \Delta b \end{pmatrix} \begin{pmatrix} \sigma_b^2 & -\sigma_{ab} \\ -\sigma_{ab} & \sigma_a^2 \end{pmatrix} \begin{pmatrix} \Delta a \\ \Delta b \end{pmatrix} \\ &= \frac{1}{\sigma_a^2 \sigma_b^2 - \sigma_{ab}^2} \begin{pmatrix} \Delta a & \Delta b \end{pmatrix} \begin{pmatrix} \sigma_b^2 \Delta a - \sigma_{ab} \Delta b \\ -\sigma_{ab} \Delta a + \sigma_a^2 \Delta b \end{pmatrix} \\ &= \frac{1}{\sigma_a^2 \sigma_b^2 - \sigma_{ab}^2} \left[\sigma_b^2 (\Delta a)^2 - 2\sigma_{ab} \Delta a \Delta b + \sigma_a^2 (\Delta b)^2 \right] \end{aligned} \quad (101)$$

With the correlation coefficient ρ_{ab} defined by

$$\rho_{ab} = \frac{\sigma_{ab}}{\sigma_a \sigma_b} \quad (102)$$

and defining each parameter offset in units of its standard deviation by

$$\delta_a = \frac{\Delta a}{\sigma_a} \quad (103)$$

$$\delta_b = \frac{\Delta b}{\sigma_b} \quad (104)$$

Eq. 101 becomes

$$\Delta\chi^2 = \frac{1}{1 - \rho_{ab}^2} [\delta_a^2 - 2\rho_{ab}\delta_a\delta_b + \delta_b^2] \quad (105)$$

With no correlation, $\rho_{ab} = 0$ and the $\Delta\chi^2 = 1$ contour is given by

$$1 = \delta_a^2 + \delta_b^2 \quad (106)$$

This is the equation for the unit circle.

The effect of non-zero correlation is clearer in the space of the variables δ'_a and δ'_b rotated 45° counterclockwise from those of δ_a and δ_b and defined by

$$\begin{aligned} \delta'_a &= \frac{\delta_a + \delta_b}{\sqrt{2}} \\ \delta'_b &= \frac{-\delta_a + \delta_b}{\sqrt{2}} \end{aligned} \quad (107)$$

The inverse transformation is then

$$\begin{aligned} \delta_a &= \frac{\delta'_a - \delta'_b}{\sqrt{2}} \\ \delta_b &= \frac{\delta'_a + \delta'_b}{\sqrt{2}} \end{aligned} \quad (108)$$

Figure 1 shows these axes along with various $\Delta\chi^2 = 1$ ellipses.

Using the transformation to δ'_a and δ'_b in Eq. 105 for the $\Delta\chi^2 = 1$ contour gives

$$\begin{aligned} 1 &= \frac{1}{1 - \rho_{ab}^2} \left[\left(\frac{\delta'_a - \delta'_b}{\sqrt{2}} \right)^2 - 2\rho_{ab} \left(\frac{\delta'_a - \delta'_b}{\sqrt{2}} \right) \left(\frac{\delta'_a + \delta'_b}{\sqrt{2}} \right) + \left(\frac{\delta'_a + \delta'_b}{\sqrt{2}} \right)^2 \right] \\ &= \frac{1}{1 - \rho_{ab}^2} \left[\delta_a'^2 + \delta_b'^2 - 2\rho_{ab} \left(\frac{\delta_a'^2}{2} - \frac{\delta_b'^2}{2} \right) \right] \\ &= \frac{1}{(1 + \rho_{ab})(1 - \rho_{ab})} \left[\delta_a'^2 (1 - \rho_{ab}) + \delta_b'^2 (1 + \rho_{ab}) \right] \\ &= \frac{\delta_a'^2}{(1 + \rho_{ab})} + \frac{\delta_b'^2}{(1 - \rho_{ab})} \end{aligned} \quad (109)$$

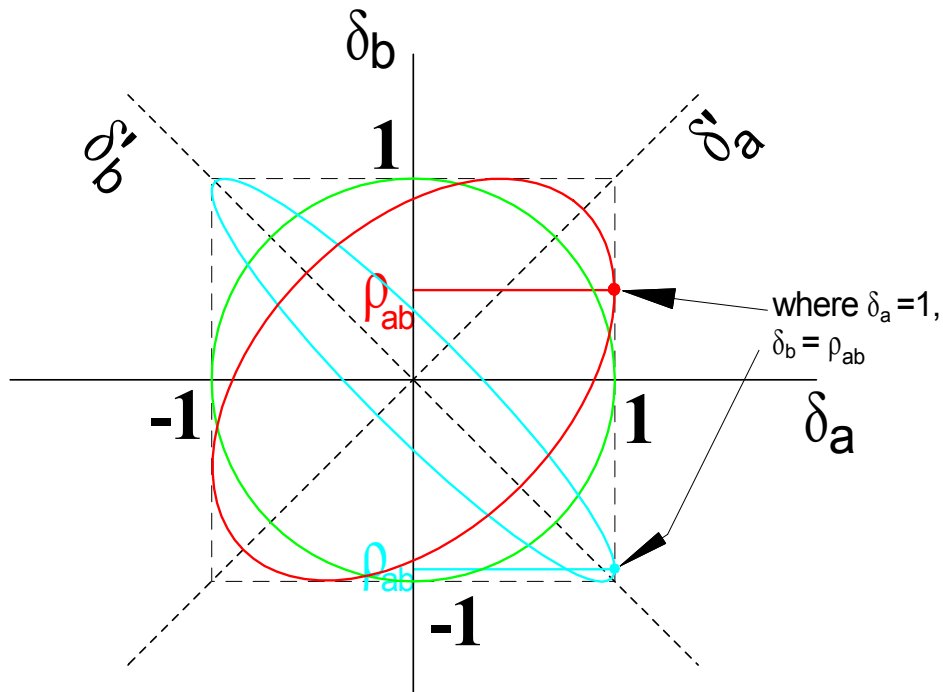


Figure 1: Three $\Delta\chi^2 = 1$ ellipses are inscribed in a unit square, each with principal axes along the square's diagonals. The green circle is for uncorrelated parameters, the red ellipse is for a medium positive correlation and the cyan ellipse is for a strong negative correlation.

With correlation, the $\Delta\chi^2 = 1$ contour is an ellipse with its principal axes along the δ'_a and δ'_b directions, i.e., oriented at 45° to the δ_a and δ_b axes and intersecting the δ'_a -axis at $\pm\sqrt{1 + \rho_{ab}}$ and intersecting the δ'_b -axis at $\pm\sqrt{1 - \rho_{ab}}$. The correlation coefficient is limited to $-1 < \rho_{ab} < 1$ and so the semimajor axis (long radius) must be between one ($\rho_{ab} = 0$) and $\sqrt{2}$ ($\rho_{ab} = \pm 1$) and the semiminor axis (short radius) must be between zero ($\rho_{ab} = \pm 1$) and one ($\rho_{ab} = 0$). Notice that the ellipse area is given by $\pi\sqrt{(1 + \rho_{ab})(1 - \rho_{ab})} = \pi\sqrt{1 - \rho_{ab}^2}$ and goes from a maximum of π for $\rho_{ab} = 0$ down to zero as ρ_{ab} goes toward ± 1 .

When $|\rho_{ab}| \ll 1$, the ellipse is only slightly squashed along one axis and slightly expanded along the other. When ρ_{ab} is positive, the major axis is along δ'_a and the ellipse has extra area in the two quadrants where δ_a and δ_b have the same sign. If one parameter turns out to be above its best estimate, the other is more likely to be above its best estimate. When ρ_{ab} is negative, the major axis is along δ'_b and the ellipse has extra area in the two quadrants where δ_a and δ_b have the opposite sign. If one parameter turns out to be above its best estimate, the other is more likely to be below its best estimate.

With nonlinear fitting functions, the linearity with fitting parameters (first order Taylor expansion) is assumed to be valid over several standard deviations of those parameters. If the data set is small or has large errors, the fitting parameter standard deviations may be large

enough that this condition is violated. It may also be violated if the fitting function derivatives depend too strongly on the parameters. Consequently, it is a good idea to test this linearity over an appropriate range. (For linear fitting functions, the formulation is exact and such a test is unnecessary.)

Finding the covariance matrix elements using Eq. 97 for the variances and Eq. 98 for the covariances should use parameter offsets that achieve values for $\Delta\chi^2$ that are neither too small nor too large. To get a $\Delta\chi^2 = 9$, Eq. 97 predicts that a_M will need to be offset by $3\sigma_M$. For values of $\Delta\chi^2$ much larger than this, the fitting function will be evaluated for parameters where the Taylor expansion of Eq. 62 need not and may not be valid.

Because it is often difficult to predict where the Taylor expansion might fail, it is good practice to try the procedure with several values of Δa_M . Choose parameter values both above and below a_M and choose various sizes for Δa_M —to get $\Delta\chi^2$ values up to at least nine or so. Results for the variances and covariances that are relatively constant over this range of approximately $\pm 3\sigma_M$ would be an indication that the Taylor expansion is probably valid over this range as well. The 68%, 95%, and 99% confidence intervals would then likely be very close to the standard ones of $a_M \pm \sigma_M$, $a_M \pm 2\sigma_M$, $a_M \pm 3\sigma_M$, respectively. On the other hand, variance and covariance calculations that differ significantly between $\Delta\chi^2$ values of 0.1 and 1 would be an indication that the Taylor expansion is already failing when a_M is only one standard deviation from the best fit. In this case, a more sophisticated statistical analysis will be needed to determine accurate confidence intervals.

Rather than inspecting covariances directly, it is preferable to divide these off-diagonal elements of the covariance matrix by the product of the standard deviations of the two parameters involved. This ratio then gives the correlation coefficients. Correlation among parameters must be taken into account in any conclusions involving more than one of them. Correlation coefficients near one indicate the two parameters involved have very similar effects on the fit and can be particularly troublesome—even for conclusions involving only one of the parameters.