

8.1. Least squares

BY E. PRINCE AND P. T. BOGGS

The process of arriving at a model for a crystal structure may usefully be considered to consist of two distinct stages. The first, which may be called determination, involves the use of chemical and physical intuition, direct methods, Fourier and Patterson methods, and other techniques to arrive at an approximate model for the structure that incorporates unit-cell dimensions, space group, chemical composition, and information with respect to the immediate environment of each atom. The second stage, which we shall call refinement, involves finding the values of adjustable parameters in the model that give the best fit between the predicted diffraction intensities and those observed in an experiment, in order to extract precise information about interatomic distances and bond angles, thermal motion, site occupancies, electron distribution, and so forth. Although there are several different criteria for the best fit to data, such as maximum likelihood and maximum entropy, one of the most commonly used is the method of least squares. This chapter discusses both numerical and statistical aspects of refinement by the method of least squares. Because both aspects make extensive use of linear algebra, we begin with a summary of definitions and fundamental operations in linear algebra (Stewart, 1973; Prince, 1994) and of basic definitions and concepts in mathematical statistics (Draper & Smith, 1981; Box, Hunter & Hunter, 1978). We then discuss the principles of linear and nonlinear least squares and conclude with an extensive discussion of numerical methods used in practical implementation of the technique.

8.1.1. Definitions

8.1.1.1. Linear algebra

A *matrix* is an ordered, rectangular array of numbers, real or complex. Matrices will be denoted by upper-case, bold italic letters, \mathbf{A} . Their individual elements will be denoted by upper-case, italic letters with subscripts. A_{ij} denotes the element in the i th row and the j th column of \mathbf{A} . A matrix with only one row is a *row vector*; a matrix with only one column is a *column vector*. Vectors will be denoted by lower-case, bold roman letters, and their elements will be denoted by lower-case, italic letters with single subscripts. Scalar constants will usually be denoted by lower-case, Greek letters.

A matrix with the same number of rows as columns is *square*. If $A_{ij} = 0$ for all $i > j$, \mathbf{A} is *upper triangular*. If $A_{ij} = 0$ for all $i < j$, \mathbf{A} is *lower triangular*. If $A_{ij} = 0$ for all $i \neq j$, \mathbf{A} is *diagonal*. If $A_{ij} = 0$ for all i and j , \mathbf{A} is *null*. A matrix, \mathbf{B} , such that $B_{ij} = A_{ji}$ for all i and j is the *transpose* of \mathbf{A} , and is denoted by \mathbf{A}^T . Matrices with the same dimensions may be added and subtracted: $(\mathbf{A} + \mathbf{B})_{ij} = A_{ij} + B_{ij}$. A matrix may be multiplied by a scalar: $(\alpha\mathbf{A})_{ij} = \alpha A_{ij}$. Multiplication of matrices is defined by $(\mathbf{AB})_{ij} = \sum_{k=1}^m A_{ik}B_{kj}$, where m is the number of columns of \mathbf{A} and the number of rows of \mathbf{B} (which must be equal). Addition and multiplication of matrices obey the associative law: $(\mathbf{A} + \mathbf{B}) + \mathbf{C} = \mathbf{A} + (\mathbf{B} + \mathbf{C})$; $(\mathbf{AB})\mathbf{C} = \mathbf{A}(\mathbf{BC})$. Multiplication of matrices obeys the distributive law: $\mathbf{A}(\mathbf{B} + \mathbf{C}) = \mathbf{AB} + \mathbf{AC}$. Addition of matrices obeys the commutative law: $\mathbf{A} + \mathbf{B} = \mathbf{B} + \mathbf{A}$, but multiplication, except in certain (important) special cases, does not: $\mathbf{AB} \neq \mathbf{BA}$. The transpose of a product is the product of the transposes of the factors in reverse order: $(\mathbf{AB})^T = \mathbf{B}^T\mathbf{A}^T$.

The *trace* of a square matrix is the sum of its diagonal elements. The *determinant* of an $n \times n$ square matrix, \mathbf{A} , denoted by $|\mathbf{A}|$, is the sum of $n!$ terms, each of which is a product of the diagonal elements of a matrix derived from \mathbf{A} by permuting columns or rows (see Stewart, 1973). The *rank* of a matrix (not necessarily square) is the dimension of the largest square submatrix that can be formed from it, by selecting rows and columns, whose determinant is not equal to zero. A matrix has *full column rank* if its rank is equal to its number of columns. A square matrix whose diagonal elements are equal to one and whose off-diagonal elements are equal to zero is an *identity matrix*, denoted by \mathbf{I} . If $|\mathbf{A}| \neq 0$, \mathbf{A} is *nonsingular*, and there exists a matrix \mathbf{A}^{-1} , the *inverse* of \mathbf{A} , such that $\mathbf{AA}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{I}$. If $|\mathbf{A}| = 0$, \mathbf{A} is *singular*, and has no inverse. The *adjoint*, or *conjugate transpose*, of \mathbf{A} is a matrix, \mathbf{A}^\dagger , such that $A_{ij}^\dagger = A_{ji}^*$, where the asterisk indicates complex conjugate. If $\mathbf{A}^\dagger = \mathbf{A}^{-1}$, \mathbf{A} is *unitary*. If the elements of a unitary matrix are real, it is *orthogonal*. From this definition, if \mathbf{A} is orthogonal, it follows that

$$\sum_{i=1}^n A_{ij}^2 = 1$$

for all j , and

$$\sum_{i=1}^n A_{ij}A_{ik} = 0$$

if $j \neq k$. By analogy, two column vectors, \mathbf{x} and \mathbf{y} , are said to be orthogonal if $\mathbf{x}^T\mathbf{y} = 0$.

For any square matrix, \mathbf{A} , there exists a set of vectors, \mathbf{x}_i , such that $\mathbf{Ax}_i = \lambda_i\mathbf{x}_i$, where λ_i is a scalar. The values λ_i are the *eigenvalues* of \mathbf{A} , and the vectors \mathbf{x}_i are the corresponding *eigenvectors*. If $\mathbf{A} = \mathbf{A}^\dagger$, \mathbf{A} is *Hermitian*, and, if the elements are real, $\mathbf{A} = \mathbf{A}^T$, so that \mathbf{A} is *symmetric*. It can be shown (see, for example, Stewart, 1973) that, if \mathbf{A} is Hermitian, all eigenvalues are real, and there exists a unitary matrix, \mathbf{T} , such that $\mathbf{D} = \mathbf{T}^\dagger\mathbf{AT}$ is diagonal, with the elements of \mathbf{D} equal to the eigenvalues of \mathbf{A} , and the columns of \mathbf{T} are the eigenvectors. An $n \times n$ symmetric matrix therefore has n mutually orthogonal eigenvectors. If the product $\mathbf{x}^T\mathbf{Ax}$ is greater than (or equal to) zero for any non-null vector, \mathbf{x} , \mathbf{A} is *positive (semi)definite*. Because \mathbf{x} may be, in particular, an eigenvector, all eigenvalues of a positive (semi)definite matrix are greater than (or equal to) zero. Any matrix of the form $\mathbf{B}^T\mathbf{B}$ is positive semi-definite, and, if \mathbf{B} has full column rank, $\mathbf{A} = \mathbf{B}^T\mathbf{B}$ is positive definite. If \mathbf{A} is positive definite, there exists an upper triangular matrix, \mathbf{R} , or, equivalently, a lower triangular matrix, \mathbf{L} , with positive diagonal elements, such that $\mathbf{R}^T\mathbf{R} = \mathbf{LL}^T = \mathbf{A}$. \mathbf{R} , or \mathbf{L} , is called the *Cholesky factor* of \mathbf{A} . The *magnitude*, *length* or *Euclidean norm* of a vector, \mathbf{x} , denoted by $\|\mathbf{x}\|$, is defined by $\|\mathbf{x}\| = (\mathbf{x}^T\mathbf{x})^{1/2}$. The *induced matrix norm* of a matrix, \mathbf{B} , denoted $\|\mathbf{B}\|$, is defined as the maximum value of $\|\mathbf{Bx}\|/\|\mathbf{x}\| = (\mathbf{x}^T\mathbf{B}^T\mathbf{Bx}/\mathbf{x}^T\mathbf{x})^{1/2}$ for $\|\mathbf{x}\| > 0$. Because $\mathbf{x}^T\mathbf{B}^T\mathbf{Bx}$ will have its maximum value for a fixed value of $\mathbf{x}^T\mathbf{x}$ when \mathbf{x} is parallel to the eigenvector that corresponds to the largest eigenvalue of $\mathbf{B}^T\mathbf{B}$, this definition implies that $\|\mathbf{B}\|$ is equal to the square root of the largest eigenvalue of $\mathbf{B}^T\mathbf{B}$. The *condition number* of \mathbf{B} is the square root of the ratio of the largest and smallest eigenvalues of $\mathbf{B}^T\mathbf{B}$. (Other definitions of norms exist, with corresponding definitions of condition number. We shall not be concerned with any of these.)

8.1. LEAST SQUARES

We shall make extensive use of the so-called QR decomposition, which is defined as follows: For any $n \times p$ ($n \geq p$) real matrix, \mathbf{Z} , there exists an $n \times n$ orthogonal matrix, \mathbf{Q} , such that

$$\mathbf{Q}^T \mathbf{Z} = \begin{pmatrix} \mathbf{R} \\ \mathbf{O} \end{pmatrix}, \quad (8.1.1.1)$$

where \mathbf{R} is a $p \times p$ upper triangular matrix, and \mathbf{O} denotes an $(n - p) \times p$ null matrix. Thus, we have

$$\mathbf{Z} = \mathbf{Q} \begin{pmatrix} \mathbf{R} \\ \mathbf{O} \end{pmatrix}, \quad (8.1.1.2)$$

which is known as the QR decomposition of \mathbf{Z} . If we partition \mathbf{Q} as $(\mathbf{Q}_Z, \mathbf{Q}_\perp)$, where \mathbf{Q}_Z has dimensions $n \times p$, and \mathbf{Q}_\perp has dimensions $n \times (n - p)$, (8.1.1.2) becomes

$$\mathbf{Z} = \mathbf{Q}_Z \mathbf{R}, \quad (8.1.1.3)$$

which is known as the QR factorization. We shall make use of the following facts. First, \mathbf{R} is nonsingular if and only if the columns of \mathbf{Z} are linearly independent; second, the columns of \mathbf{Q}_Z form an orthonormal basis for the range space of \mathbf{Z} , that is, they span the same space as \mathbf{Z} ; and, third, the columns of \mathbf{Q}_\perp form an orthonormal basis for the null space of \mathbf{Z}^T , that is, $\mathbf{Z}^T \mathbf{Q}_\perp = \mathbf{O}$.

There are two common procedures for computing the QR factorization. The first makes use of *Householder transformations*, which are defined by

$$\mathbf{H} = \mathbf{I} - 2\mathbf{x}\mathbf{x}^T, \quad (8.1.1.4)$$

where $\mathbf{x}^T \mathbf{x} = 1$. \mathbf{H} is symmetric, and $\mathbf{H}^2 = \mathbf{I}$, so that \mathbf{H} is orthogonal. In three dimensions, \mathbf{H} corresponds to a reflection in a mirror plane perpendicular to \mathbf{x} , because of which Stewart (1973) has suggested the alternative term *elementary reflector*. A vector \mathbf{v} is transformed by $\mathbf{H}\mathbf{v}$ into the vector $\|\mathbf{v}\|\mathbf{e}$, where \mathbf{e} represents a vector with $e_1 = 1$, and $e_i = 0$ for $i \neq 1$, if

$$\mathbf{x} = [\mathbf{v} - \|\mathbf{v}\|\mathbf{e}] / \|\mathbf{v} - \|\mathbf{v}\|\mathbf{e}\|. \quad (8.1.1.5)$$

The factorization procedure for an $n \times p$ matrix, \mathbf{A} (Stewart, 1973; Anderson *et al.*, 1992), takes as \mathbf{v} in the first step the first column of \mathbf{A} , and forms $\mathbf{A}_1 = \mathbf{H}_1 \mathbf{A}$, which has zeros in all elements of the first column below the diagonal. In the second step, \mathbf{v} has a zero as the first element and is filled out by those elements of the second column of \mathbf{A}_1 on or below the diagonal. $\mathbf{A}_2 = \mathbf{H}_2 \mathbf{A}_1$ then has zeros in all elements below the diagonal in the first two columns. This process is repeated $(p - 2)$ more times, after which $\mathbf{Q} = \mathbf{H}_p \dots \mathbf{H}_2 \mathbf{H}_1$, and $\mathbf{R} = \mathbf{A}_p$ is upper triangular.

QR factorization by Householder transformations requires for efficiency that the entire $n \times p$ matrix be stored in memory, and requires of order np^2 operations. A procedure that requires storage of only the upper triangle makes use of *Givens rotations*, which are 2×2 matrices of the form

$$\mathbf{G} = \begin{pmatrix} \cos \alpha & \sin \alpha \\ -\sin \alpha & \cos \alpha \end{pmatrix}. \quad (8.1.1.6)$$

Multiplication of a $2 \times m$ matrix, \mathbf{B} , by \mathbf{G} will put a zero in the B_{21} element if $\alpha = \arctan B_{21}/B_{11}$. The factorization of \mathbf{A} involves reading, or computing, the rows of \mathbf{A} one at a time. In the first step, matrix \mathbf{B}_1 consists of the first row of \mathbf{R} and the current row of \mathbf{A} , from which the first element is eliminated. In the second step, \mathbf{B}_{21} is the second row of \mathbf{R} and the $(p - 1)$ non-zero elements of the second row of the transformed \mathbf{B}_1 . After the first p rows have been treated, each additional row of \mathbf{A} requires $2p(p + 1)$ multiplications to fill it with zeros. However, because the operation is easily vectorized, the time required may be a

small proportion of the total computing time on a vector oriented computer.

8.1.1.2. Statistics

A *probability density function*, which will be abbreviated p.d.f., is a function, $\Phi(x)$, such that the probability of finding the random variable x in the interval $a \leq x \leq b$ is given by

$$p(a \leq x \leq b) = \int_a^b \Phi(x) dx.$$

A p.d.f. has the properties

$$\Phi(x) \geq 0, \quad -\infty < x < +\infty,$$

and

$$\int_{-\infty}^{+\infty} \Phi(x) dx = 1.$$

A *cumulative distribution function*, which will be abbreviated c.d.f., is defined by

$$\Psi(x) = \int_{-\infty}^x \Phi(t) dt.$$

The properties of $\Phi(x)$ imply that $0 \leq \Psi(x) \leq 1$, and $\Phi(x) = d\Psi(x)/dx$. The *expected value* of a function, $f(x)$, of random variable x is defined by

$$\langle f(x) \rangle = \int_{-\infty}^{+\infty} f(x)\Phi(x) dx.$$

If $f(x) = x^n$, $\langle f(x) \rangle = \langle x^n \rangle$ is the *n*th moment of $\Phi(x)$. The first moment, often denoted by μ , is the *mean* of $\Phi(x)$. The second moment about the mean, $\langle (x - \langle x \rangle)^2 \rangle$, usually denoted by σ^2 , is the *variance* of $\Phi(x)$. The positive square root of the variance is the *standard deviation*.

For a vector, \mathbf{x} , of random variables, x_1, x_2, \dots, x_n , the *joint probability density function*, or *joint p.d.f.*, is a function, $\Phi_J(\mathbf{x})$, such that

$$p(a_1 \leq x_1 \leq b_1; a_2 \leq x_2 \leq b_2; \dots; a_n \leq x_n \leq b_n) = \int_{a_1}^{b_1} \int_{a_2}^{b_2} \dots \int_{a_n}^{b_n} \Phi_J(\mathbf{x}) dx_1 dx_2 \dots dx_n. \quad (8.1.1.7)$$

The *marginal p.d.f.* of an element (or a subset of elements), x_i , is a function, $\Phi_M(x_i)$, such that

$$p(a_i \leq x_i \leq b_i) = \int_{a_i}^{b_i} \Phi_M(x_i) dx_i = \int_{-\infty}^{+\infty} \dots \int_{a_i}^{b_i} \dots \int_{-\infty}^{\infty} \Phi_J(\mathbf{x}) dx_1 \dots dx_i \dots dx_n. \quad (8.1.1.8)$$

This is a p.d.f. for x_i alone, irrespective of the values that may be found for any other element of \mathbf{x} . For two random variables, x and y (either or both of which may be vectors), the *conditional p.d.f. of x given $y = y_0$* is defined by

$$\Phi_C(x|y_0) = c\Phi_J(x, y)_{y=y_0},$$

where $c = 1/\Phi_M(y_0)$ is a *renormalizing factor*. This is a p.d.f. for x when it is known that $y = y_0$. If $\Phi_C(x|y) = \Phi_M(x)$ for all y , or, equivalently, if $\Phi_J(x, y) = \Phi_M(x)\Phi_M(y)$, the random variables x and y are said to be *statistically independent*.

8. REFINEMENT OF STRUCTURAL PARAMETERS

Moments may be defined for multivariate p.d.f.s in a manner analogous to the one-dimensional case. The mean is a vector defined by

$$\mu_i = \langle x_i \rangle = \int x_i \Phi(\mathbf{x}) \, d\mathbf{x},$$

where the volume of integration is the entire domain of \mathbf{x} . The *variance-covariance matrix* is defined by

$$\begin{aligned} V_{ij} &= \langle (x_i - \langle x_i \rangle)(x_j - \langle x_j \rangle) \rangle \\ &= \int (x_i - \langle x_i \rangle)(x_j - \langle x_j \rangle) \Phi_j(\mathbf{x}) \, d\mathbf{x}. \end{aligned} \quad (8.1.1.9)$$

The diagonal elements of \mathbf{V} are the variances of the marginal p.d.f.s of the elements of \mathbf{x} , that is, $V_{ii} = \sigma_i^2$. It can be shown that, if x_i and x_j are statistically independent, $V_{ij} = 0$ when $i \neq j$. If two vectors of random variables, \mathbf{x} and \mathbf{y} , are related by a linear transformation, $\mathbf{x} = \mathbf{B}\mathbf{y}$, the means of their joint p.d.f.s are related by $\mu_x = \mathbf{B}\mu_y$, and their variance-covariance matrices are related by $\mathbf{V}_x = \mathbf{B}\mathbf{V}_y\mathbf{B}^T$.

8.1.2. Principles of least squares

The method of least squares may be formulated as follows: Given a set of n observations, y_i ($i = 1, 2, \dots, n$), that are measurements of quantities that can be described by differentiable model functions, $M_i(\mathbf{x})$, where \mathbf{x} is a vector of parameters, x_j ($j = 1, 2, \dots, p$), find the values of the parameters for which the sum

$$S = \sum_{i=1}^n w_i [y_i - M_i(\mathbf{x})]^2 \quad (8.1.2.1)$$

is minimum. Here, w_i represents a weight assigned to the i th observation. The values of the parameters that give the minimum value of S are called *estimates* of the parameters, and a function of the data that locates the minimum is an *estimator*. A necessary condition for S to be a minimum is for the gradient to vanish, which gives a set of simultaneous equations, the *normal equations*, of the form

$$\frac{\partial S}{\partial x_j} = -2 \sum_{i=1}^n w_i [y_i - M_i(\mathbf{x})] \frac{\partial M_i(\mathbf{x})}{\partial x_j} = 0. \quad (8.1.2.2)$$

The model functions, $M_i(\mathbf{x})$, are, in general, nonlinear, and there are no direct ways to solve these systems of equations. Iterative methods for solving them are discussed in Section 8.1.4. Much of the analysis of results, however, is based on the assumption that linear approximations to the model functions are good approximations in the vicinity of the minimum, and we shall therefore begin with a discussion of linear least squares.

To express linear relationships, it is convenient to use matrix notation. Let $\mathbf{M}(\mathbf{x})$ and \mathbf{y} be column vectors whose i th elements are $M_i(\mathbf{x})$ and y_i . Similarly, let \mathbf{b} be a vector and \mathbf{A} be a matrix such that a linear approximation to the i th model function can be written

$$M_i(\mathbf{x}) = b_i + \sum_{j=1}^p A_{ij} x_j. \quad (8.1.2.3)$$

Equations (8.1.2.3) can be written, in matrix form,

$$\mathbf{M}(\mathbf{x}) = \mathbf{b} + \mathbf{A}\mathbf{x}, \quad (8.1.2.4)$$

and, for this linear model, (8.1.2.1) becomes

$$S = [(\mathbf{y} - \mathbf{b}) - \mathbf{A}\mathbf{x}]^T \mathbf{W} [(\mathbf{y} - \mathbf{b}) - \mathbf{A}\mathbf{x}], \quad (8.1.2.5)$$

where \mathbf{W} is a diagonal matrix whose diagonal elements are $W_{ii} = w_i$. In this notation, the normal equations (8.1.2.2) can be written

$$\mathbf{A}^T \mathbf{W} \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{W} (\mathbf{y} - \mathbf{b}), \quad (8.1.2.6)$$

and their solution is

$$\hat{\mathbf{x}} = (\mathbf{A}^T \mathbf{W} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{W} (\mathbf{y} - \mathbf{b}). \quad (8.1.2.7)$$

If $W_{ii} > 0$ for all i , and \mathbf{A} has full column rank, then $\mathbf{A}^T \mathbf{W} \mathbf{A}$ will be positive definite, and S will have a unique minimum at $\mathbf{x} = \hat{\mathbf{x}}$. The matrix $\mathbf{H} = (\mathbf{A}^T \mathbf{W} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{W}$ is a $p \times n$ matrix that relates the n -dimensional observation space to the p -dimensional parameter space and is known as the *least-squares estimator*; because each element of $\hat{\mathbf{x}}$ is a linear function of the observations, it is a *linear estimator*. [Note that, in actual practice, the matrix \mathbf{H} is not actually evaluated, except, possibly, in very small problems. Rather, the linear system $\mathbf{A}^T \mathbf{W} \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{W} (\mathbf{y} - \mathbf{b})$ is solved using the methods of Section 8.1.3.]

The least-squares estimator has some special properties in statistical analysis. Suppose that the elements of \mathbf{y} are experimental observations drawn at random from populations whose means are given by the model, $\mathbf{M}(\mathbf{x})$, for some unknown \mathbf{x} , which we wish to estimate. This may be written

$$\langle \mathbf{y} - \mathbf{b} \rangle = \mathbf{A}\mathbf{x}. \quad (8.1.2.8)$$

The expected value of the least-squares estimate is

$$\begin{aligned} \langle \hat{\mathbf{x}} \rangle &= \langle (\mathbf{A}^T \mathbf{W} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{W} (\mathbf{y} - \mathbf{b}) \rangle \\ &= (\mathbf{A}^T \mathbf{W} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{W} \langle \mathbf{y} - \mathbf{b} \rangle \\ &= (\mathbf{A}^T \mathbf{W} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{W} \mathbf{A} \mathbf{x} \\ &= \mathbf{x}. \end{aligned} \quad (8.1.2.9)$$

If the expected value of an estimate is equal to the variable to be estimated, the estimator is said to be *unbiased*. Equation (8.1.2.9) shows that the least-squares estimator is an unbiased estimator for \mathbf{x} , independent of \mathbf{W} , provided only that \mathbf{y} is an unbiased estimate of $\mathbf{M}(\mathbf{x})$, the matrix $\mathbf{A}^T \mathbf{W} \mathbf{A}$ is nonsingular, and the elements of \mathbf{W} are constants independent of \mathbf{y} and $\mathbf{M}(\mathbf{x})$. Let \mathbf{V}_x and \mathbf{V}_y be the variance-covariance matrices for the joint p.d.f.s of the elements of \mathbf{x} and \mathbf{y} , respectively. Then, $\mathbf{V}_x = \mathbf{H}\mathbf{V}_y\mathbf{H}^T$. Let \mathbf{G} be the matrix $(\mathbf{A}^T \mathbf{V}_y^{-1} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{V}_y^{-1}$, so that $\hat{\mathbf{x}} = \mathbf{G}(\mathbf{y} - \mathbf{b})$ is the particular least-squares estimate for which $\mathbf{W} = \mathbf{V}_y^{-1}$. Then, $\mathbf{V}_x = \mathbf{G}\mathbf{V}_y\mathbf{G}^T$. If \mathbf{V}_y is positive definite, its lower triangular Cholesky factor, \mathbf{L} , exists, so that $\mathbf{L}\mathbf{L}^T = \mathbf{V}_y$. [If \mathbf{V} is diagonal, \mathbf{L} is also diagonal, with $L_{ii} = (\mathbf{V}_y)_{ii}^{1/2}$.] It is readily verified that the matrix product $[(\mathbf{H} - \mathbf{G})\mathbf{L}][(\mathbf{H} - \mathbf{G})\mathbf{L}]^T = \mathbf{H}\mathbf{V}_y\mathbf{H}^T - \mathbf{G}\mathbf{V}_y\mathbf{G}^T$, but the diagonal elements of this product are the sums of squares of the elements of rows of $(\mathbf{H} - \mathbf{G})\mathbf{L}$, and are therefore greater than or equal to zero. Therefore, the diagonal elements of \mathbf{V}_x , which are the variances of the marginal p.d.f.s of the elements of $\hat{\mathbf{x}}$, are minimum when $\mathbf{W} = \mathbf{V}_y^{-1}$.

Thus, the least-squares estimator is unbiased for any positive-definite weight matrix, \mathbf{W} , but the variances of the elements of the vector of estimated parameters are minimized if $\mathbf{W} = \mathbf{V}_y^{-1}$. [Note also that $\mathbf{V}_x = (\mathbf{A}^T \mathbf{W} \mathbf{A})^{-1}$ if, and only if, $\mathbf{W} = \mathbf{V}_y^{-1}$.] For this reason, the least-squares estimator with weights proportional to the reciprocals of the variances of the observations is referred to as the *best linear unbiased estimator* for the parameters of a model describing those observations. (These specific results are included in a more general result known as the *Gauss-Markov theorem*.)

The analysis up to this point has assumed that the model is linear, that is that the expected values of the observations can be expressed by $\langle \mathbf{y} \rangle = \mathbf{b} + \mathbf{A}\mathbf{x}$, where \mathbf{A} is some matrix. In crystallography, of course, the model is highly nonlinear, and this assumption is not valid. The principles of linear least squares