

1

Least Squares Approximation

Theory attracts practice as the magnet attracts iron.

—Gauss, Karl Friedrich

The celebrated concept of least squares approximation is introduced in this chapter. Least squares can be used in a wide variety of categorical applications, including: curve fitting of data, parameter identification, and system model realization. Many examples from diverse fields fall under these categories, for instance determining the damping properties of a fluid-filled damper as a function of temperature, identification of aircraft dynamic and static aerodynamic coefficients, orbit and attitude determination, position determination using triangulation, and modal identification of vibratory systems. Even modern control strategies, for instance certain adaptive controllers, use the least squares approximation to update model parameters in the control system. The broad utility implicit in the aforementioned examples strongly confirms that the least squares approximation is worthy of study.

Before we begin analytical and mathematical discussions, let us first define some common quantities used throughout this chapter and the text. For any variable or parameter in estimation, there are three quantities of interest: the true value, the measured value, and the estimated value. The true value (or “truth”) is usually unknown in practice. This represents the actual value sought of the quantity being approximated by the estimator. Unadorned symbols are used to represent the true values. The measured value denotes the quantity which is directly determined from a sensor. For example, in orbit determination a radar is often used to obtain a measure of the range to a vehicle. In actuality, this is not a totally accurate statement since the truly measured quantity given by the radar is not the range. Radars work by “shining” a beam of energy (usually microwaves) at an object and analyzing the spectral content of the energy that gets reflected back. Signal processing of the measured return energy can yield estimates of range (or range rate). For navigation purposes, we often assume that the measured quantity is the computed range, because this is a direct function of the truly measured quantity, which is the reflected energy received by the radar. Measurements are never perfect, since they will always contain errors. Thus, measurements are usually modeled using a function of the true values plus some error. The measured values of the truth x are typically denoted by \bar{x} . Estimated values of x are determined from the estimation process itself, and are found using a combination of a static/dynamic model and the measurements. These values are denoted by \hat{x} . Other quantities used commonly in estimation are the measurement error

(measurement value minus true value) and the residual error (measurement value minus estimated value). Thus, for a measurable quantity x , the following two equations hold:

$$\begin{array}{rcccc} \text{measured value} & = & \text{true value} & + & \text{measurement error} \\ \tilde{x} & = & x & + & v \end{array}$$

and

$$\begin{array}{rcccc} \text{measured value} & = & \text{estimated value} & + & \text{residual error} \\ \tilde{x} & = & \hat{x} & + & e \end{array}$$

The actual measurement error (v), like the true value, is never known in practice. However, the errors in the mechanism that physically generate this error are usually approximated by some known process (often by a zero-mean Gaussian noise process with known variance). These assumed known statistical properties of the measurement errors are often employed to weight the relative importance of various measurements used in the estimation scheme. Unlike the measurement error, the residual error is known explicitly and is easily computed once an estimated value has been found. The residual error is often used to drive the estimator itself. It should be evident that both measurement errors and residual errors play important roles in the theoretical and computational aspects of estimation.

1.1 A Curve Fitting Example

To explore Gauss' connection between theory and practice, we introduce the concept of least squares by considering a simple example that will be used to motivate the theoretical developments of this chapter. Displayed in [Figure 1.1](#) are measurements of some process $y(t)$. At this point we do not consider the physical connotations of the particular process, but it may be useful to think of $y(t)$ as a stock quote history for a particular company. You want to determine a mathematical model for $y(t)$ in order to predict future prospects for the company. Measurements (e.g., closing stock price) of $y(t)$, denoted by $\tilde{y}(t)$, are given for a 6-month time frame. In order to insure an accurate model fit, you have been informed that the residual errors (i.e., between the measured values and estimated values) must have an absolute mean of ≤ 0.0075 and a standard deviation of ≤ 0.125 . With a large number of samples (m), the sample mean (μ) and sample standard deviation (σ) for the residual error can be computed using¹ (we will derive these later)

$$\mu = \frac{1}{m} \sum_{i=1}^m [\tilde{y}(t_i) - \hat{y}(t_i)] \quad (1.1)$$

$$\sigma^2 = \frac{1}{m-1} \sum_{i=1}^m \{[\tilde{y}(t_i) - \hat{y}(t_i)] - \mu\}^2 \quad (1.2)$$

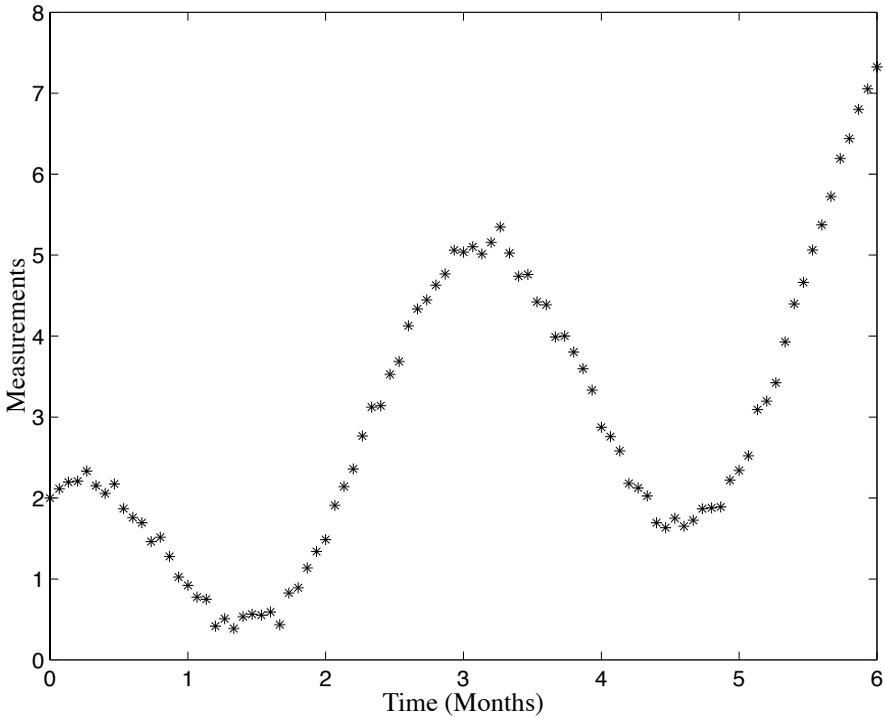


Figure 1.1: Measurements of $y(t)$

where $\hat{y}(t)$ denotes the estimate of $y(t)$.

Now in your quest to establish a model which predicts the behavior of $y(t)$, you might naturally attempt evaluation of some previously developed models. After some research you have found two models, given by

$$\text{Model 1: } y_1(t) = c_1t + c_2 \sin(t) + c_3 \cos(2t) \tag{1.3}$$

$$\text{Model 2: } y_2(t) = d_1(t + 2) + d_2t^2 + d_3t^3 \tag{1.4}$$

where t is given in months, and c_1, c_2, c_3 and d_1, d_2, d_3 are constants. The next step is to evaluate “how well” each of these models predicts the measurements with “optimum” values of c_i and d_i . The process of fitting curves, such as Models 1 and 2, to measured data is known in statistics as *regression*.

For the moment, continuing the discussion of the hypothetical problem solving situation, let us assume that you have read and digested the discussion that will come later in §1.2.1 on the method of *linear least squares*. Also, you have employed a least squares algorithm to determine the coefficients in the two models, and found that the “optimum” coefficients are

$$(\hat{c}_1, \hat{c}_2, \hat{c}_3) = (0.9967, 0.9556, 2.0030) \tag{1.5}$$

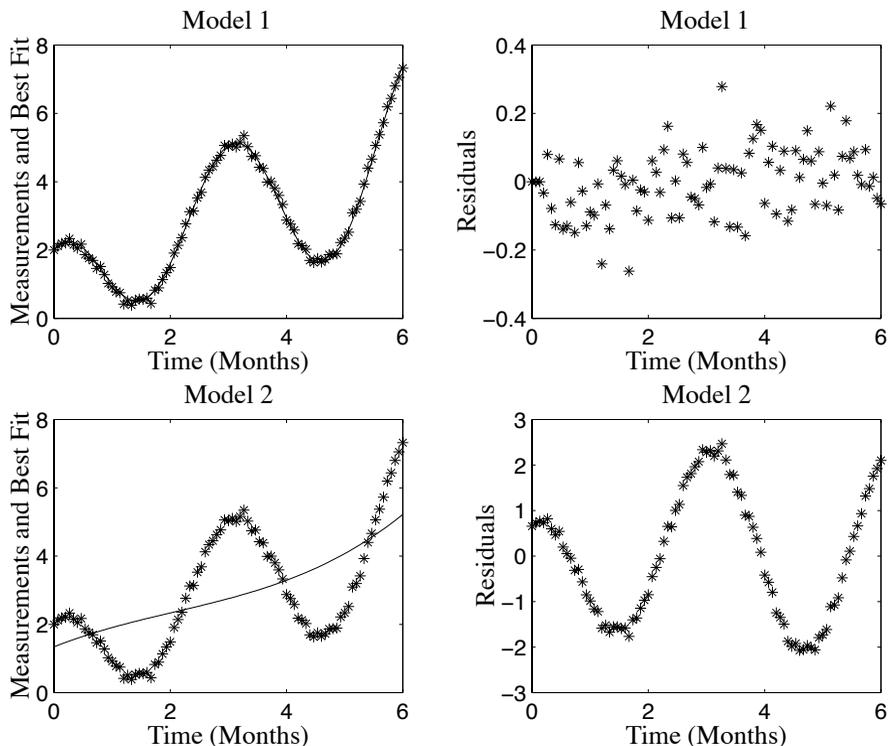


Figure 1.2: Best Fit and Residual Errors for Both Models

$$(\hat{d}_1, \hat{d}_2, \hat{d}_3) = (0.6721, -0.1303, 0.0210) \quad (1.6)$$

Plots of each model's fit superimposed on the measured data, and residual errors are shown in Figure 1.2. As is clearly evident, Model 1 is able to obtain the best fit with the determined coefficients. This can also be seen by comparing the sample mean and sample standard deviation of both fits using Equations (1.1) and (1.2). For Model 1 the sample mean is 1×10^{-5} and the sample standard deviation is 0.0921. For Model 2 the sample mean is 1×10^{-5} and the sample standard deviation is 1.3856. This shows that Model 1 meets both minimum requirements for a good fit, while Model 2 does not.

From the above analysis, you make the qualitative observation that Model 1 is a much better representation of $y(t)$'s behavior than is Model 2. From Figure 1.2, you observe that Model 1's residual errors are "random" in appearance, while Model 2's best fit failed to predict significant trends in the data. Having no reason to suspect that systematic errors are present in the measurements or in Model 1, you conclude that Model 1 can be used to provide an accurate assessment of $y(t)$'s behavior.

Since Model 1 was used to fit the measured data accurately, you might now make the logical hypothesis that this model can be used to *predict* future values for $y(t)$.

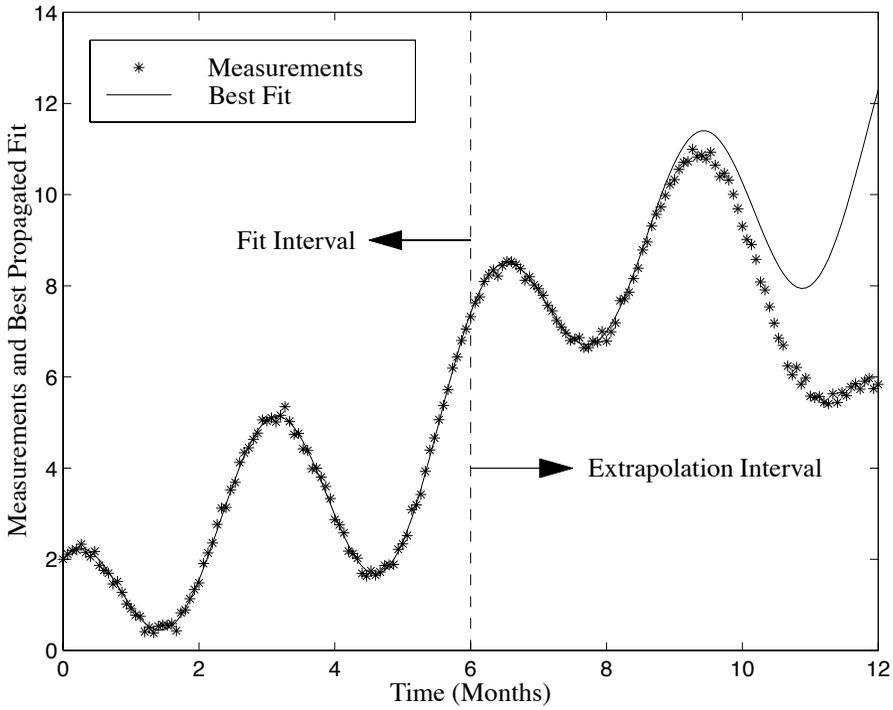


Figure 1.3: Best Fit for $y(t)$ Propagated to 12 Months

The trends in the data of the fit interval, and therefore our model, indicate that the stock prices will continue an upward trend and will more than double in 12 months. Putting your trust in this “get rich quick” scheme, suppose you invest a great amount of money in the stock. But, as is often true in many “get rich quick” schemes, this dangerous extrapolation failed. A plot of Model 1’s predictions, with coefficients given in Equation (1.5), superimposed on the measured data over a 12-month period is shown in Figure 1.3. This shows that you have actually lost money in the stock if you invest after 6 months and hold it until 12 months.

In reality, the synthetic measurements of Figure 1.1 were calculated using the following equation:

$$\tilde{y}(t) = t + \sin(t) + 2 \cos(2t) - \frac{0.4e^t}{1 \times 10^4} + v(t) \tag{1.7}$$

where the simulated measurement errors $v(t)$ were calculated by a zero-mean Gaussian noise generator with a standard deviation given by $\sigma = 0.1$. In the above example, Model 1 clearly can be used to “estimate” $y(t)$ for the first 6 months where the estimate is “supported” by many measurements, but does a poor job predicting future values. This is due to the fact that the unmodeled exponential term in Equation (1.7) begins to dominate the other terms after time $t = 10$. To further illustrate this, let us

consider the following model:

$$\text{Model 3 : } y_3(t) = x_1 t + x_2 \sin(t) + x_3 \cos(2t) + x_4 e^t \quad (1.8)$$

We observe that this model is in fact the correct model, in the absence of measurement errors. Upon applying the method of least squares using the first 6 months of measurements in [Figure 1.1](#), we find the optimal estimates of the coefficients \hat{x}_i are

$$(\hat{x}_1, \hat{x}_2, \hat{x}_3, \hat{x}_4) = (0.9958, 0.9979, 2.0117, -4.232 \times 10^{-5}) \quad (1.9)$$

It is significant to note, if we zero the measurement errors with this model, the least squares estimates give exactly the true parameter values $(1, 1, 2, -4 \times 10^{-5})$. It is also of interest to ask the question: “How well can we predict the future when we use the correct model?” This question is answered by repeating the calculation underlying [Figure 1.3](#), using the correct model (1.8) and best estimates (1.9) derived over the first 6 months of data. These results are shown in [Figure 1.4](#). Comparing Figures 1.3 and 1.4, it is evident that using the correct model (1.8) vastly improves the 6-month extrapolation accuracy. The extrapolation still diverges slowly from the subsequent measurements over months 10 to 12. This is because the coefficient estimates derived from any finite set of measurements can be expected to contain estimation errors even when the model structure is perfect. We will develop full insight into the issue: “How do measurement errors propagate into errors of the estimated parameters?”

The above contrived example demonstrates many important issues in estimation theory. First, a challenging facet of practical estimation applications is correctly specifying the system’s mathematical model. Also, the first two models contain a t term, but the corresponding numerical estimates of the t coefficient are drastically different in the two best fits. In many real-world problems, dominant terms in a mathematical model will have a correct mathematical structure, but higher-order effects may be poorly understood. Finally, unknown higher order effects and parameter estimation errors can produce erroneous results, especially outside of the measurement domain considered, as shown in [Figure 1.3](#).

Model development is the least tractable aspect of the problem setup and solution, insofar as employing universally applicable procedures. It is unlikely, indeed, that mathematically complicated physical phenomena can be correctly modeled *a priori* by anyone unfamiliar with the basic principles underlying the phenomena. In short, intelligent formulation and application of estimation algorithms require intimate knowledge of the field in which the estimation problem is embedded. In numerous cases, decisions regarding which variable should be measured, the frequency with which data should be collected, the necessary measurement accuracy, and the best mathematical model can be inferred directly from theoretical analysis of the system. *Estimation theory can be developed apart from considering a particular dynamic system, but successful applications almost invariably rely jointly upon understanding estimation theory and the principles governing the system under consideration.*

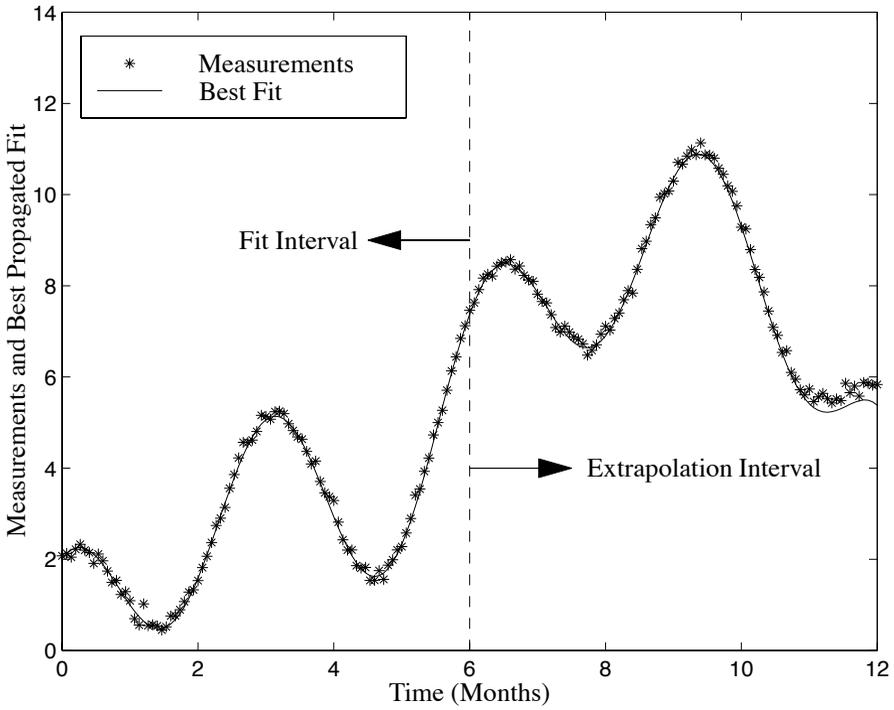


Figure 1.4: Best Fit for $y(t)$ Propagated to 12 Months

1.2 Linear Batch Estimation

In this section we formally introduce Gauss’ principle of linear least squares. This principle will be found to be central to the solution of a large family of estimation problems. Suppose that you have in hand a set (or a “batch”) of measured values, \bar{y}_j , of a process $y(t)$, taken at known discrete instants of time t_j :

$$\{\bar{y}_1, t_1; \bar{y}_2, t_2; \dots; \bar{y}_m, t_m\} \tag{1.10}$$

and a proposed mathematical model of the form

$$y(t) = \sum_{i=1}^n x_i h_i(t), \quad m \geq n \tag{1.11}$$

where

$$h_i(t) \in \{h_1(t), h_2(t), \dots, h_n(t)\} \tag{1.12}$$

are a set of independent specified *basis* functions. For example, Equations (1.3) and (1.4) each contains three basis functions in our previous work in §1.1. The x_i are a set

of constants whose numerical values are unknown. From Equation (1.11) it follows that the variables x and y are related according to a simple linear regression model. It seems altogether reasonable to select the optimum x -values based upon a measure of “how well” the proposed model (1.11) predicts the measurements (1.10). Toward this end, we seek a set of estimates, denoted by $\{\hat{x}_1, \hat{x}_2, \dots, \hat{x}_n\}$, which can be used in Equation (1.11) to predict $y(t)$. Errors, however, can arise between the “true” value $y(t)$ and the predicted (estimated) value $\hat{y}(t)$ from a number of sources, including:

- measurement errors
- incorrect choice of x -values
- modeling errors, i.e., the actual process being observed may not be accurately modeled by Equation (1.11).

In virtually every application, some combination of these error sources is present.

We first formally relate the measurements \tilde{y}_j and the estimated output \hat{y}_j to the true and estimated x -values using the mathematical model of Equation (1.11):

$$\tilde{y}_j \equiv \tilde{y}(t_j) = \sum_{i=1}^n x_i h_i(t_j) + v_j, \quad j = 1, 2, \dots, m \quad (1.13)$$

$$\hat{y}_j \equiv \hat{y}(t_j) = \sum_{i=1}^n \hat{x}_i h_i(t_j), \quad j = 1, 2, \dots, m \quad (1.14)$$

where v_j is the measurement error. At this point of the discussion, we consider the measurement error to be some unknown process that may include random as well as deterministic characteristics (in the next chapter, we will elaborate more on v_j). It is important to remember that \tilde{y}_j is a *measured* quantity (i.e., it is the output of the measurement process). We have assumed that the measurement process is *modeled* by Equation (1.13). Next, consider the following identity:

$$\tilde{y}_j = \sum_{i=1}^n \hat{x}_i h_i(t_j) + e_j, \quad j = 1, 2, \dots, m \quad (1.15)$$

where the *residual error* e_j is defined by

$$e_j \equiv \tilde{y}_j - \hat{y}_j \quad (1.16)$$

Equation (1.15) can be rewritten in compact matrix form as

$$\tilde{\mathbf{y}} = H\hat{\mathbf{x}} + \mathbf{e} \quad (1.17)$$

where

$$\tilde{\mathbf{y}} = [\tilde{y}_1 \ \tilde{y}_2 \ \cdots \ \tilde{y}_m]^T = \text{measured } y\text{-values}$$

$$\mathbf{e} = [e_1 \ e_2 \ \cdots \ e_m]^T = \text{residual errors}$$

$$\hat{\mathbf{x}} = [\hat{x}_1 \ \hat{x}_2 \ \cdots \ \hat{x}_n]^T = \text{estimated } x\text{-values}$$

$$H = \begin{bmatrix} h_1(t_1) & h_2(t_1) & \cdots & h_n(t_1) \\ h_1(t_2) & h_2(t_2) & \cdots & h_n(t_2) \\ \vdots & \vdots & & \vdots \\ h_1(t_m) & h_2(t_m) & \cdots & h_n(t_m) \end{bmatrix}$$

and the superscript T denotes the matrix transpose operation. In a similar manner, Equations (1.13) and (1.14) can also be written in compact form as

$$\tilde{\mathbf{y}} = H\mathbf{x} + \mathbf{v} \tag{1.18}$$

$$\hat{\mathbf{y}} = H\hat{\mathbf{x}} \tag{1.19}$$

where

$$\mathbf{x} = [x_1 \ x_2 \ \cdots \ x_n]^T = \text{true } x\text{-values}$$

$$\mathbf{v} = [v_1 \ v_2 \ \cdots \ v_m]^T = \text{measurement errors}$$

$$\hat{\mathbf{y}} = [\hat{y}_1 \ \hat{y}_2 \ \cdots \ \hat{y}_m]^T = \text{estimated } y\text{-values}$$

$$\tilde{\mathbf{y}} = [\tilde{y}_1 \ \tilde{y}_2 \ \cdots \ \tilde{y}_m]^T = \text{measured } y\text{-values}$$

Equations (1.17) and (1.18) are identical, of course, if $\hat{\mathbf{x}} = \mathbf{x}$, and if the assumption of zero model errors is valid. Both of these equations, (1.17) and (1.18), are commonly referred to as the “observation equations.”

1.2.1 Linear Least Squares

Gauss’s celebrated *principle of least squares*² selects, as an optimum choice for the unknown parameters, the particular $\hat{\mathbf{x}}$ that minimizes the sum square of the residual errors, given by

$$J = \frac{1}{2} \mathbf{e}^T \mathbf{e} \tag{1.20}$$

Substituting Equation (1.17) for \mathbf{e} into Equation (1.20) and using the fact that a scalar equals its transpose yields

$$J = J(\hat{\mathbf{x}}) = \frac{1}{2} (\tilde{\mathbf{y}}^T \tilde{\mathbf{y}} - 2\tilde{\mathbf{y}}^T H\hat{\mathbf{x}} + \hat{\mathbf{x}}^T H^T H\hat{\mathbf{x}}) \tag{1.21}$$

The $1/2$ multiplier of J does have a statistical significance, as will be shown in Chapter 2. We seek to find the $\hat{\mathbf{x}}$ that minimizes J . Using the matrix calculus differentiation

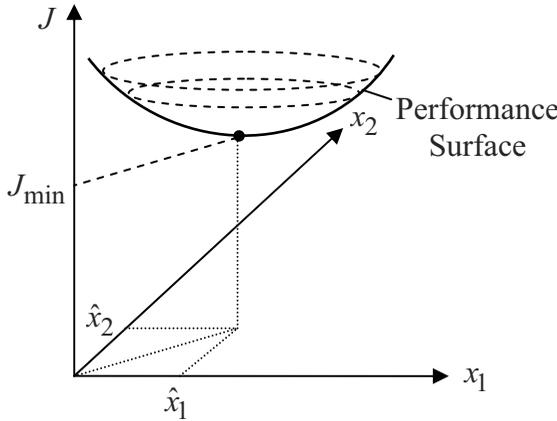


Figure 1.5: Convex Performance Surface for Order $n = 2$ Problem

rules developed in §B.5, it follows that for a global minimum of the quadratic function of Equation (1.21) we have the following requirements:

necessary condition

$$\nabla_{\hat{\mathbf{x}}} J \equiv \begin{bmatrix} \frac{\partial J}{\partial \hat{x}_1} \\ \vdots \\ \frac{\partial J}{\partial \hat{x}_n} \end{bmatrix} = H^T H \hat{\mathbf{x}} - H^T \hat{\mathbf{y}} = \mathbf{0} \tag{1.22}$$

sufficient condition

$$\nabla_{\hat{\mathbf{x}}}^2 J \equiv \frac{\partial^2 J}{\partial \hat{\mathbf{x}} \partial \hat{\mathbf{x}}^T} = H^T H \text{ must be positive definite} \tag{1.23}$$

where $\nabla_{\hat{\mathbf{x}}} J$ is the *Jacobian* and $\nabla_{\hat{\mathbf{x}}}^2 J$ is the *Hessian* (see Appendix B). Consider the sufficient condition first. Any matrix B such that

$$\mathbf{x}^T B \mathbf{x} \geq 0 \tag{1.24}$$

for all $\mathbf{x} \neq \mathbf{0}$ is called positive semi-definite. By setting $\mathbf{h} = H\mathbf{x}$ and squaring, we easily obtain the scalar $h^2 = \mathbf{h}^T \mathbf{h} \geq 0$, so $H^T H$ is always positive semi-definite. It becomes positive definite when H is of maximum rank (n).

The function J is a performance surface in $n + 1$ -dimensional space.³ This performance surface has a convex shape of an n -dimensional parabola with one *distinct* minimum. An example of this performance surface for $n = 2$ is the three-dimensional bowl-shaped surface shown in Figure 1.5.

From the necessary conditions of Equation (1.22), we now have the “normal equations”

$$(H^T H) \hat{\mathbf{x}} = H^T \hat{\mathbf{y}} \tag{1.25}$$

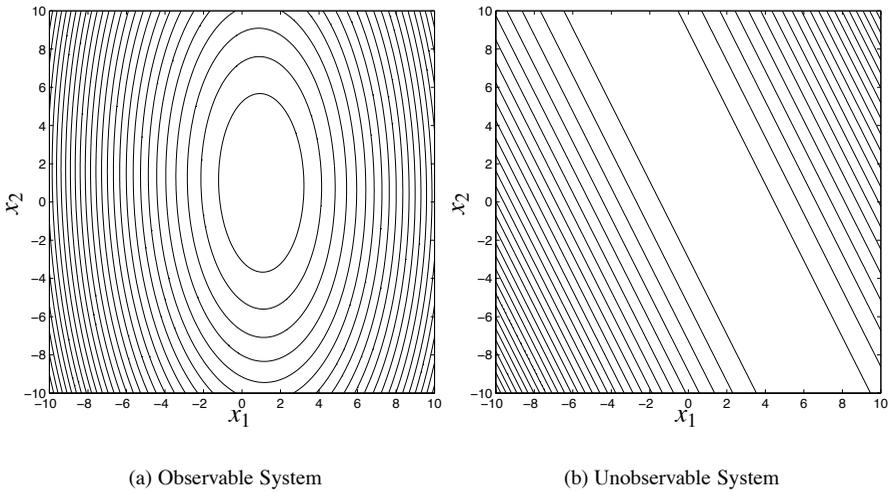


Figure 1.6: Contour Plots for an Observable and Unobservable System

If the rank of H is n (i.e., there are at least n independent observation equations), then $H^T H$ is *strictly* positive definite and can be inverted to obtain the explicit solution for the optimal estimate:

$$\hat{\mathbf{x}} = (H^T H)^{-1} H^T \tilde{\mathbf{y}} \tag{1.26}$$

Equation (1.17) is the matrix equivalent of Gauss’ original “equations of condition” which he wrote in index/summation notation.² Equation (1.26) serves as the most common basis for algorithms that solve simple least squares problems.

The inverse of $H^T H$ is required to determine $\hat{\mathbf{x}}$. This inverse exists only if the number of linearly independent observations is equal to or greater than the number of unknown x_i . To show this concept, consider a simple least squares problem with $\mathbf{x} = [1 \ 1]^T$ and two basis functions given by $H_1 = [\sin t \ 2 \cos t]$ and $H_2 = [\sin t \ 2 \sin t]$. Clearly, H_1 provides a linearly independent set of basis functions, while H_2 does not because the second column of H_2 is twice the first column. A plot of the contour lines using H_1 is shown in Figure 1.6(a), which clearly shows a minimum at the true value for $\mathbf{x} = [1 \ 1]^T$. A plot of the contour lines using H_2 is shown in Figure 1.6(b), which shows that an infinite number of solutions are possible. More details on observability for dynamic systems are discussed in §A.4.

One of the implicit advantages of least squares is that the order of the matrix inverse is equal to the number of *unknowns*, not the number of measurement observations. The explicit solution (1.26) can be seen to play a role similar to $\mathbf{x} = H^{-1} \mathbf{y}$ in solving $\mathbf{y} = H\mathbf{x}$ for the $m = n$ case. We note that Gauss introduced his method of Gaussian elimination to solve the normal equations (1.25), by reducing $(H^T H)$ to upper triangular form, then solving for $\hat{\mathbf{x}}$ by back substitution (see Appendix B).

Example 1.1: Let us illustrate the basic concept of using linear least squares for curve fitting a batch of measured data. The measurements are generated using the following model:

$$\hat{y}_i = 0.3 \sin(t_i) + 0.5 \cos(t_i) + 0.1t_i + v_i$$

with simulated measurement errors calculated using a zero-mean Gaussian noise generator with a standard deviation given by $\sigma = \sqrt{0.001}$. A total of 101 discrete measurements of the system are given sampled every 0.1 seconds.

The assumed basis function matrix is given by

$$H = \begin{bmatrix} \sin(t_0) & \cos(t_0) & t_0 & \cos(t_0) \sin(t_0) & t_0^2 \\ \sin(t_1) & \cos(t_1) & t_1 & \cos(t_1) \sin(t_1) & t_1^2 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \sin(t_{100}) & \cos(t_{100}) & t_{100} & \cos(t_{100}) \sin(t_{100}) & t_{100}^2 \end{bmatrix}$$

Note we have two “extra” basis functions as compared to the model used to generate the synthetic measurements. We thus expect that the estimated coefficients for these basis functions should be near zero in the least squares solution. Using Equation (1.26) the estimated coefficients are found to be given by

$$\hat{\mathbf{x}} = [0.3019 \ 0.5072 \ 0.1027 \ 0.0012 \ -0.0003]^T$$

Good agreement is given between the estimated coefficients and the true coefficients, and the estimated coefficients associated with the “extra” basis functions are indeed near zero as expected.

Example 1.2: In this example we employ linear least squares to estimate the parameters of a simple dynamic system. Consider the following dynamic system:

$$\dot{y} = ay + bu, \quad (\dot{}) \equiv \frac{d}{dt}()$$

where u is an exogenous (i.e., externally specified) input, and a and b are constants. The system can also be represented in discrete time with constant sampling interval Δt by (see §A.5)

$$y_{k+1} = \Phi y_k + \Gamma u_k$$

where the integer k is the sample index, and

$$\Phi = e^{a\Delta t}$$

$$\Gamma = \int_0^{\Delta t} b e^{at} dt = \frac{b}{a} (e^{a\Delta t} - 1)$$

The goal of this problem is to determine the constants Φ and Γ given a discrete set of measurements \tilde{y}_k and inputs u_k . For the particular problem in which it is known that u is given by an impulse input with magnitude 100 (i.e., $u_1 = 100$ and $u_k = 0$ for $k \geq 2$), a total of 101 discrete measurements of the system are given with $\Delta t = 0.1$, and are shown in Figure 1.7. In order to set up the least squares problem, we construct the following basis function matrix:

$$H = \begin{bmatrix} \tilde{y}_1 & u_1 \\ \tilde{y}_2 & u_2 \\ \vdots & \vdots \\ \tilde{y}_{101} & u_{101} \end{bmatrix}$$

so

$$\begin{bmatrix} \tilde{y}_2 \\ \tilde{y}_3 \\ \vdots \\ \tilde{y}_{101} \end{bmatrix} = H \begin{bmatrix} \hat{\Phi} \\ \hat{\Gamma} \end{bmatrix} + \begin{bmatrix} e_2 \\ e_3 \\ \vdots \\ e_{101} \end{bmatrix}$$

Now, estimates for Φ and Γ can be determined using Equation (1.26) directly:

$$\begin{bmatrix} \hat{\Phi} \\ \hat{\Gamma} \end{bmatrix} = (H^T H)^{-1} H^T [\tilde{y}_2 \ \tilde{y}_3 \ \dots \ \tilde{y}_{101}]^T$$

Using the measurements shown in Figure 1.7 the computed estimates are found to be

$$\begin{bmatrix} \hat{\Phi} \\ \hat{\Gamma} \end{bmatrix} = \begin{bmatrix} 0.9048 \\ 0.0950 \end{bmatrix}$$

In reality, the synthetic measurements of Figure 1.7 were generated using the following true values:

$$\begin{bmatrix} \Phi \\ \Gamma \end{bmatrix} = \begin{bmatrix} 0.9048 \\ 0.0952 \end{bmatrix}$$

with simulated measurement errors calculated using a zero-mean Gaussian noise generator with a standard deviation given by $\sigma = 0.08$.

The above example clearly involves a *dynamic* system; however, even though this system is modeled using a linear differential equation with constant coefficients, we are still able to bring the relationship (between measured quantities and constants which determine the model) to a linear algebraic equation, and therefore, we can use the principle of linear least squares. Also, the basis functions involve the measurements themselves, which is perhaps counterintuitive, but still is a valid approach, although not truly “optimal,” as discussed in §2.8.4. The measurements appear in the basis functions because one of the sought parameters, Φ , multiplies y_k in the assumed model (the other parameter multiplies the input). This example clearly shows the power of least squares for dynamic model *identification*. We note in passing that the multi-dimensional generalization and sophistication of this example lead to the

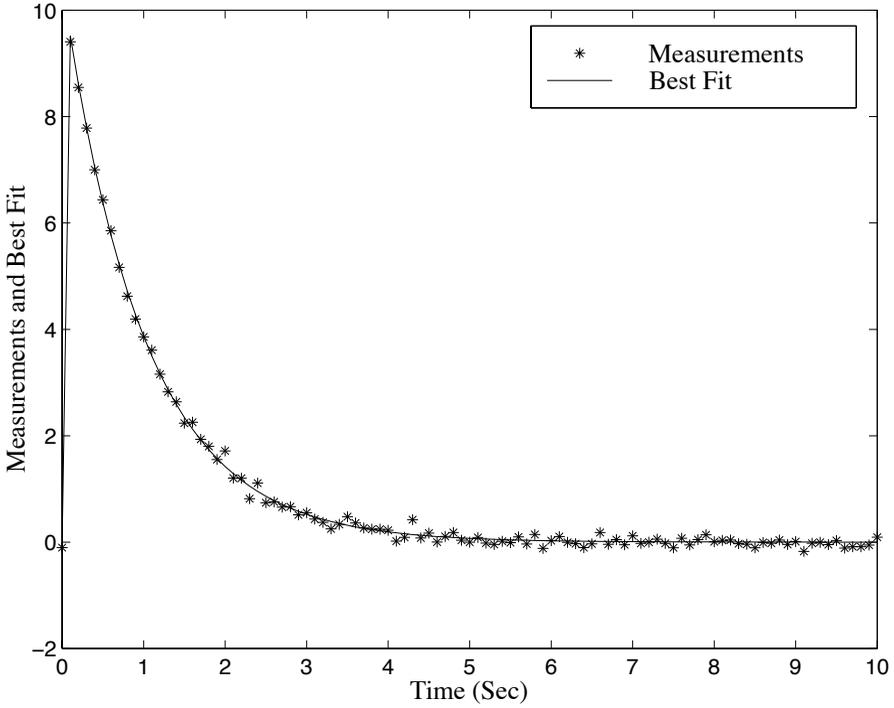


Figure 1.7: Measurements of $y(t)$ and Best Fit

Eigensystem Realization Algorithm (ERA).⁴ This algorithm is presented in Chapter 6.

1.2.2 Weighted Least Squares

The least squares criterion in Equation (1.20), minimized to determine $\hat{\mathbf{x}}$, implicitly places equal emphasis on each measurement \tilde{y}_j . For the common event that the measurements are made with unequal precision, this “equal weight” approach seems logically unsound. Thus, the question arises as to how to select proper weights. One might intuitively select weights for each measurement that are inversely proportional to the measurement’s estimated precision (i.e., a measurement with zero error should be weighted infinitely, while a measurement with infinite error should be weighted zero). Additionally, we shall see in Chapter 2 that a statistically optimal (“maximum likelihood”) choice for the weights is the reciprocal of the measurement error variance. In order to incorporate appropriate weighting, we set up a least squares

criterion of the form

$$J = \frac{1}{2} \mathbf{e}^T W \mathbf{e} \tag{1.27}$$

We now seek to determine $\hat{\mathbf{x}}$ that minimizes J , where W is an $m \times m$ symmetric matrix (it is symmetric because the terms $e_i e_j$, $i \neq j$, are always weighted equally with the corresponding $e_j e_i$ terms). In order that $\hat{\mathbf{x}}$ yield a minimum of Equation (1.27), we have the requirements:

necessary condition

$$\nabla_{\hat{\mathbf{x}}} J = H^T W H \hat{\mathbf{x}} - H^T W \tilde{\mathbf{y}} = \mathbf{0} \tag{1.28}$$

sufficient condition

$$\nabla_{\hat{\mathbf{x}}}^2 J = H^T W H \text{ must be positive definite.} \tag{1.29}$$

From the necessary condition in Equation (1.28), we obtain the solution for $\hat{\mathbf{x}}$ given by

$$\hat{\mathbf{x}} = (H^T W H)^{-1} H^T W \tilde{\mathbf{y}} \tag{1.30}$$

Also, Equation (1.29) clearly shows that W must be positive definite.

Example 1.3: To illustrate the power of weighted least squares, we will employ a subset of 31 measurements from the 91 measurements shown in [Figure 1.1](#). Also, the first three measurements are known to contain smaller measurement errors than the remaining measurements. Toward this end, the structure of the weighting matrix now becomes

$$W = \text{diag} [w \ w \ w \ 1 \ \dots \ 1]$$

where $\text{diag} [\]$ denotes a diagonal matrix. Using Model 1 in Equation (1.3) and the subset of 31 measurements with $w = 1$ (i.e., reduces to standard least squares) yields the following estimates:

$$(\hat{c}_1, \hat{c}_2, \hat{c}_3) = (1.0278, 0.8750, 1.9884)$$

Observe the unsurprising fact that the estimates are further from their true values (1, 1, 2) than the estimates (1.5) resulting from all 91 measurements. However, since we know that the first three measurements are better than the remaining measurements, we can improve the estimates using weighted least squares. A summary of the solutions for $\hat{\mathbf{x}}$ with various values of w is shown below.

w	$\hat{\mathbf{x}}$	constraint residual norm
1×10^0	(1.0278, 0.8750, 1.9884)	3.21×10^{-2}
1×10^1	(1.0388, 0.8675, 2.0018)	1.17×10^{-2}
1×10^2	(1.0258, 0.8923, 2.0049)	7.87×10^{-3}
1×10^5	(0.9047, 1.0949, 2.0000)	5.91×10^{-5}
1×10^7	(0.9060, 1.0943, 2.0000)	1.10×10^{-5}
1×10^{10}	(0.9932, 1.0068, 2.0000)	4.55×10^{-7}
1×10^{15}	(0.9970, 1.0030, 2.0000)	0.97×10^{-9}

Downloaded by [Utah State University] at 23:06 03 December 2013

One can see that the residual constraint error (i.e., the computed norm of the measurements minus the estimates for the first three observations) decreases as more weight is used. However, this does not generally guarantee that the estimates ($\hat{\mathbf{x}}$) are closer to their true values. The interaction of the basis function therefore plays an important role in weighted least squares. Still, if the weight is sufficiently large, the estimates are indeed closer to their true values, as expected. In this simulation, the first three measurements were obtained with no measurement errors. However, perfect estimates (with zero associated model error) cannot be achieved since the exponential term in Equation (1.7) is still present in the simulated measurements, which is not in the assumed model. Weighted least squares can improve the estimates if some knowledge of the relative accuracy of the measurements is known, and can obviously be used to approximately impose constraints on an estimation process.

1.2.3 Constrained Least Squares

Minimization of the weighted least squares criterion (1.27) allows relative emphasis to be placed upon the model agreeing with certain measurements more closely than others. Consider the limiting case of a perfect measurement where the corresponding diagonal element of the weight matrix should be ∞ . This can often be accomplished in a practical situation by replacing ∞ with a “sufficiently large” number to obtain satisfactory approximations. However, we might be motivated to seek a rigorous means for imposing equality constraints in estimation problems.⁵

Suppose the original observations in Equation (1.17) partition naturally into the sub-systems $\tilde{\mathbf{y}}_1$ and $\tilde{\mathbf{y}}_2$ as

$$\begin{bmatrix} \tilde{\mathbf{y}}_1 \\ \cdot \\ \tilde{\mathbf{y}}_2 \end{bmatrix} = \begin{bmatrix} H_1 \\ \dots \\ H_2 \end{bmatrix} \hat{\mathbf{x}} + \begin{bmatrix} \mathbf{e}_1 \\ \cdot \\ \mathbf{0} \end{bmatrix} \quad (1.31)$$

or

$$\tilde{\mathbf{y}}_1 = H_1 \hat{\mathbf{x}} + \mathbf{e}_1 \quad (1.32)$$

and

$$\tilde{\mathbf{y}}_2 = H_2 \hat{\mathbf{x}} \quad (1.33)$$

where

$\tilde{\mathbf{y}}_1$ = an $m_1 \times 1$ vector of measured y -values

H_1 = an $m_1 \times n$ basis function matrix corresponding
with the measured y -values

\mathbf{e}_1 = an $m_1 \times 1$ vector of residual errors

$\tilde{\mathbf{y}}_2$ = an $m_2 \times 1$ vector of perfectly measured y -values

H_2 = an $m_2 \times n$ basis function matrix corresponding
with the perfectly measured y -values

and further assume that the dimensions satisfy

$$\begin{aligned} n &\geq m_2 \\ n &\leq m_1 \end{aligned}$$

The absence of the residual error matrix \mathbf{e}_2 in Equations (1.31) and (1.33) reflects the fact that $H_2\hat{\mathbf{x}}$ is required to equal $\tilde{\mathbf{y}}_2$ *exactly*. Thus, we can formulate the problem as a constrained minimization problem of the type discussed in Appendix D. We seek a vector $\hat{\mathbf{x}}$ that minimizes

$$J = \frac{1}{2}\mathbf{e}_1^T W_1 \mathbf{e}_1 = \frac{1}{2}(\tilde{\mathbf{y}}_1 - H_1\hat{\mathbf{x}})^T W_1 (\tilde{\mathbf{y}}_1 - H_1\hat{\mathbf{x}}) \tag{1.34}$$

subject to the satisfaction of the equality constraint

$$\tilde{\mathbf{y}}_2 - H_2\hat{\mathbf{x}} = \mathbf{0} \tag{1.35}$$

Using the method of Lagrange multipliers (Appendix D), the necessary conditions are found by minimizing the augmented function

$$J = \frac{1}{2} [\tilde{\mathbf{y}}_1^T W_1 \tilde{\mathbf{y}}_1 - 2\tilde{\mathbf{y}}_1^T W_1 H_1 \hat{\mathbf{x}} + \hat{\mathbf{x}}^T (H_1^T W_1 H_1) \hat{\mathbf{x}}] + \boldsymbol{\lambda}^T (\tilde{\mathbf{y}}_2 - H_2\hat{\mathbf{x}}) \tag{1.36}$$

where

$$\boldsymbol{\lambda} = [\lambda_1 \ \lambda_2 \ \dots \ \lambda_{m_2}]^T \tag{1.37}$$

is a vector of Lagrange multipliers. As necessary conditions for constrained minimization of J , we have the requirements:

$$\nabla_{\hat{\mathbf{x}}} J = -H_1^T W_1 \tilde{\mathbf{y}}_1 + (H_1^T W_1 H_1) \hat{\mathbf{x}} - H_2^T \boldsymbol{\lambda} = \mathbf{0} \tag{1.38}$$

and

$$\nabla_{\boldsymbol{\lambda}} J = \tilde{\mathbf{y}}_2 - H_2\hat{\mathbf{x}} = \mathbf{0}, \quad \rightarrow \tilde{\mathbf{y}}_2 = H_2\hat{\mathbf{x}} \tag{1.39}$$

Solving Equation (1.38) for $\hat{\mathbf{x}}$ yields

$$\hat{\mathbf{x}} = (H_1^T W_1 H_1)^{-1} H_1^T W_1 \tilde{\mathbf{y}}_1 + (H_1^T W_1 H_1)^{-1} H_2^T \boldsymbol{\lambda} \tag{1.40}$$

Substituting Equation (1.40) into Equation (1.39) allows for solution of the Lagrange multipliers as

$$\boldsymbol{\lambda} = [H_2(H_1^T W_1 H_1)^{-1} H_2^T]^{-1} [\tilde{\mathbf{y}}_2 - H_2(H_1^T W_1 H_1)^{-1} H_1^T W_1 \tilde{\mathbf{y}}_1] \tag{1.41}$$

Finally, substituting Equation (1.41) into Equation (1.40) allows for elimination of $\boldsymbol{\lambda}$, yielding an explicit solution for the equality constrained least squares coefficient estimates as

$$\hat{\mathbf{x}} = \bar{\mathbf{x}} + K(\tilde{\mathbf{y}}_2 - H_2\bar{\mathbf{x}}) \tag{1.42}$$

where

$$K = (H_1^T W_1 H_1)^{-1} H_2^T [H_2(H_1^T W_1 H_1)^{-1} H_2^T]^{-1} \tag{1.43}$$

and

$$\bar{\mathbf{x}} = (H_1^T W_1 H_1)^{-1} H_1^T W_1 \tilde{\mathbf{y}}_1 \quad (1.44)$$

Observe that $\bar{\mathbf{x}}$, the first term of Equation (1.42), is the least squares estimate of \mathbf{x} in the absence of the constraint equations (1.33). The second term is an additive correction in which an optimal “gain matrix” K multiplies the constraint residual $(\tilde{\mathbf{y}}_2 - H_2 \bar{\mathbf{x}})$ prior to the correction. This general “update form” (1.42) is seen often in estimation theory and is therefore an important result.

Due to the more complicated structure of Equations (1.42), (1.43), and (1.44), in comparison to algorithms for solution of the weighted least squares problem, it often proves more expedient to simply use a least squares solution with a large weight on the constraint equation. However, if the number m_2 of constraint equations is small, the number of arithmetic operations in Equations (1.42) and (1.43) can be much less than Equation (1.30). In the limit, of $m_2 = 1$ constraint, then the matrix inverse in Equation (1.43) simplifies to a scalar division.

As another important special case, consider $m_2 = n$. In this case H_2 is a square matrix, so Equation (1.43) reduces to

$$K = H_2^{-1} \quad (1.45)$$

Thus, the constrained least squares estimate becomes

$$\hat{\mathbf{x}} = H_2^{-1} \tilde{\mathbf{y}}_2 \quad (1.46)$$

This shows that the solution is dependent on the perfectly measured values and H_2 only, which is the same result obtained using a square H matrix in the standard least squares solution. Thus, if $m_2 = n$ perfect measurements are available, the solution is unaffected by an arbitrary number m of erroneous measurements.

Example 1.4: In example 1.3, weighted least squares was used to improve the estimates by incorporating knowledge of the perfectly known measurements. This result can also be obtained using constrained least squares. Again, a subset of 31 measurements is used. Three cases have been examined for the equality constraint, summarized by

$$\begin{aligned} \text{case 1: } \tilde{\mathbf{y}}_1 &= [\tilde{y}_2 \ \tilde{y}_3 \ \cdots \ \tilde{y}_{31}]^T, & \tilde{\mathbf{y}}_2 &= y_1 \\ \text{case 2: } \tilde{\mathbf{y}}_1 &= [\tilde{y}_3 \ \tilde{y}_4 \ \cdots \ \tilde{y}_{31}]^T, & \tilde{\mathbf{y}}_2 &= [y_1 \ y_2]^T \\ \text{case 3: } \tilde{\mathbf{y}}_1 &= [\tilde{y}_4 \ \tilde{y}_5 \ \cdots \ \tilde{y}_{31}]^T, & \tilde{\mathbf{y}}_2 &= [y_1 \ y_2 \ y_3]^T \end{aligned}$$

Results using constrained least squares for $\bar{\mathbf{x}}$ and $\hat{\mathbf{x}}$ are summarized for each case below.

case	$\bar{\mathbf{x}}$	$\hat{\mathbf{x}}$
1	(1.0261, 0.8766, 1.9869)	(1.0406, 0.8629, 2.0000)
2	(1.0233, 0.8789, 1.9840)	(0.9039, 1.0901, 2.0000)
3	(1.0192, 0.8820, 1.9793)	(0.9970, 1.0030, 2.0000)

We see that when one perfect measurement is used (case 1), the solution is not substantially improved over conventional least squares since $\bar{\mathbf{x}} \approx \hat{\mathbf{x}}$. However, when two perfect measurements are used (case 2), the estimates are closer to their true values. When three perfect measurements are used (case 3), which implies that $n = m_2$, the estimates are even closer to their true values. In fact, the estimates are identical within several significant digits to the case of $w = 1 \times 10^{15}$ in example 1.3. Were it not for the unaccounted error term $-0.4e^t / 1 \times 10^4$ in the simulated measurements, these would be found to agree exactly with the true coefficients (1, 1, 2).

The theoretical equivalence of an infinitely weighted measurement to an equality constraint, from the viewpoint that Equations (1.30) and (1.42) are equivalent for this limiting case, is algebraically difficult to establish. It is possible, however, and is an intuitively pleasing truth. In practical applications, one can often obtain satisfactory solutions of constrained least squares problems in a fashion analogous to this example.

1.3 Linear Sequential Estimation

In the developments of the previous section, an implicit assumption is present, namely, that all measurements are available for simultaneous (“batch”) processing. In numerous real-world applications, the measurements become available sequentially in subsets and, immediately upon receipt of a new data subset, it may be desirable to determine new estimates based upon all previous measurements (including the current subset). To simplify the initial discussion, consider only two subsets:

$$\tilde{\mathbf{y}}_1 = [\tilde{y}_{11} \ \tilde{y}_{12} \ \cdots \ \tilde{y}_{1m_1}]^T = \text{an } m_1 \times 1 \text{ vector of measurements} \tag{1.47a}$$

$$\tilde{\mathbf{y}}_2 = [\tilde{y}_{21} \ \tilde{y}_{22} \ \cdots \ \tilde{y}_{2m_2}]^T = \text{an } m_2 \times 1 \text{ vector of measurements} \tag{1.47b}$$

and the associated observation equations

$$\tilde{\mathbf{y}}_1 = H_1 \mathbf{x} + \mathbf{v}_1 \tag{1.48a}$$

$$\tilde{\mathbf{y}}_2 = H_2 \mathbf{x} + \mathbf{v}_2 \tag{1.48b}$$

where

$H_1 = \text{an } m_1 \times n \text{ known coefficient matrix of maximum rank } n \leq m_1$

$H_2 = \text{an } m_2 \times n \text{ known coefficient matrix}$

$\mathbf{v}_1, \mathbf{v}_2 = \text{vectors of measurement errors}$

$\mathbf{x} = \text{the } n \times 1 \text{ vector of unknown parameters}$

The least squares estimate, $\hat{\mathbf{x}}$, of \mathbf{x} based upon the *first* measurement subset (1.47a) follows from Equation (1.30) as

$$\hat{\mathbf{x}}_1 = (H_1^T W_1 H_1)^{-1} H_1^T W_1 \tilde{\mathbf{y}}_1 \quad (1.49)$$

where W_1 is an $m_1 \times m_1$ symmetric, positive definite matrix associated with measurements $\tilde{\mathbf{y}}_1$. It is possible to consider $\tilde{\mathbf{y}}_1$ and $\tilde{\mathbf{y}}_2$ *simultaneously* and determine an estimate $\hat{\mathbf{x}}_2$ of \mathbf{x} based upon *both* measurement subsets (1.47a) and (1.47b). Toward this end, we form the *merged* observation equations

$$\tilde{\mathbf{y}} = H\mathbf{x} + \mathbf{v} \quad (1.50)$$

where

$$\tilde{\mathbf{y}} = \begin{bmatrix} \tilde{\mathbf{y}}_1 \\ \dots \\ \tilde{\mathbf{y}}_2 \end{bmatrix}, \quad H = \begin{bmatrix} H_1 \\ \dots \\ H_2 \end{bmatrix}, \quad \mathbf{v} = \begin{bmatrix} \mathbf{v}_1 \\ \dots \\ \mathbf{v}_2 \end{bmatrix} \quad (1.51)$$

Next, we assume that the merged weight matrix is in block diagonal structure, so that*

$$W = \begin{bmatrix} W_1 & \vdots & 0 \\ \dots & \dots & \\ 0 & \vdots & W_2 \end{bmatrix} \quad (1.52)$$

Then, the optimal least squares estimate based upon the first two measurement subsets follows from Equation (1.30) as

$$\hat{\mathbf{x}}_2 = (H^T W H)^{-1} H^T W \tilde{\mathbf{y}} \quad (1.53)$$

Now, since W is block diagonal, Equation (1.53) can be expanded as

$$\hat{\mathbf{x}}_2 = [H_1^T W_1 H_1 + H_2^T W_2 H_2]^{-1} (H_1^T W_1 \tilde{\mathbf{y}}_1 + H_2^T W_2 \tilde{\mathbf{y}}_2) \quad (1.54)$$

It is clearly possible, in principle, to continue forming merged normal equations using the above procedure (upon receipt of each data subset) and solving for new optimal estimates as in Equation (1.54). However, the above route does not take efficient advantage of the calculations done in processing the previous subsets of data. The essence of the *sequential* approach to the least squares problem is to simply arrange calculations for the new estimate (e.g., $\hat{\mathbf{x}}_2$) to make efficient use of previous estimates and the associated side calculations. We begin the derivation of this approach by defining the following variables:

$$P_1 \equiv [H_1^T W_1 H_1]^{-1} \quad (1.55)$$

$$P_2 \equiv [H_1^T W_1 H_1 + H_2^T W_2 H_2]^{-1} \quad (1.56)$$

*In Chapter 2 and Appendix C, we will see that an implicit assumption here is that measurement errors can be *correlated* only to other measurements belonging to the same subset.

From these definitions it immediately follows that (assuming that both P_1^{-1} and P_2^{-1} exist)

$$P_2^{-1} = P_1^{-1} + H_2^T W_2 H_2 \quad (1.57)$$

We now rewrite Equations (1.49) and (1.54) using the definitions in Equations (1.55) and (1.56) as

$$\hat{\mathbf{x}}_1 = P_1 H_1^T W_1 \tilde{\mathbf{y}}_1 \quad (1.58)$$

$$\hat{\mathbf{x}}_2 = P_2 (H_1^T W_1 \tilde{\mathbf{y}}_1 + H_2^T W_2 \tilde{\mathbf{y}}_2) \quad (1.59)$$

Pre-multiplying Equation (1.58) by P_1^{-1} yields

$$P_1^{-1} \hat{\mathbf{x}}_1 = H_1^T W_1 \tilde{\mathbf{y}}_1 \quad (1.60)$$

Next, from Equation (1.57) we have

$$P_1^{-1} = P_2^{-1} - H_2^T W_2 H_2 \quad (1.61)$$

Substituting Equation (1.61) into Equation (1.60) leads to

$$H_1^T W_1 \tilde{\mathbf{y}}_1 = P_2^{-1} \hat{\mathbf{x}}_1 - H_2^T W_2 H_2 \hat{\mathbf{x}}_1 \quad (1.62)$$

Finally, substituting Equation (1.62) into Equation (1.59) and collecting terms gives

$$\hat{\mathbf{x}}_2 = \hat{\mathbf{x}}_1 + K_2 (\tilde{\mathbf{y}}_2 - H_2 \hat{\mathbf{x}}_1) \quad (1.63)$$

where

$$K_2 \equiv P_2 H_2^T W_2 \quad (1.64)$$

We now have a mechanism to *sequentially* provide an updated estimate, $\hat{\mathbf{x}}_2$, based upon the previous estimate, $\hat{\mathbf{x}}_1$, and associated side calculations. We can easily generalize Equations (1.63) and (1.64) to use the k^{th} estimate to determine the estimate at $k + 1$ from the $k + 1$ subset of measurements, which leads to a most important result in sequential estimation theory:

$$\hat{\mathbf{x}}_{k+1} = \hat{\mathbf{x}}_k + K_{k+1} (\tilde{\mathbf{y}}_{k+1} - H_{k+1} \hat{\mathbf{x}}_k) \quad (1.65)$$

where

$$K_{k+1} = P_{k+1} H_{k+1}^T W_{k+1} \quad (1.66)$$

$$P_{k+1}^{-1} = P_k^{-1} + H_{k+1}^T W_{k+1} H_{k+1} \quad (1.67)$$

Equation (1.65) modifies the previous best correction $\hat{\mathbf{x}}_k$ by an additional correction to account for the information contained in the $k + 1$ measurement subset. This equation is a *Kalman update equation*⁶ for computing the improved estimate $\hat{\mathbf{x}}_{k+1}$. Also, notice the similarity between Equation (1.65) and Equation (1.42). Equation (1.66) is the correction term, known as the *Kalman gain matrix*. The sequential least squares

algorithm plays an important role for linear (and nonlinear) dynamic *state* estimation, as will be seen in the Kalman filter in §3.3. Equation (1.65) is in fact a linear difference equation, commonly found in digital control analysis. This equation may be rearranged as

$$\hat{\mathbf{x}}_{k+1} = [I - K_{k+1}H_{k+1}]\hat{\mathbf{x}}_k + K_{k+1}\tilde{\mathbf{y}}_{k+1} \quad (1.68)$$

which clearly is in the form of a time-varying dynamic system. Therefore, linear tools can be used to check stability, dynamic response times, etc.

The specific form for P^{-1} in Equation (1.67) is known as the *information matrix recursion*.[†] The current approach for computing P_{k+1} involves computing the inverse of Equation (1.67), which offers no advantage over inverting the normal equations in their original *batch* processing in Equation (1.53). This is due to the fact that an $n \times n$ inverse must still be performed. We might wonder if there is an easier way to compute P_{k+1} given that we have computed P_k previously. As it turns out, when the number of measurements m in the new data subset is small compared to n (as is usually the case), a *small rank adjustment* to the already computed P_k can be calculated efficiently using the Sherman-Morrison-Woodbury *matrix inversion lemma*.⁷ Let

$$F = [A + BCD]^{-1} \quad (1.69)$$

where

F = an arbitrary $n \times n$ matrix

A = an arbitrary $n \times n$ matrix

B = an arbitrary $n \times m$ matrix

C = an arbitrary $m \times m$ matrix

D = an arbitrary $m \times n$ matrix

Then, assuming all inverses exist

$$F = A^{-1} - A^{-1}B(DA^{-1}B + C^{-1})^{-1}DA^{-1} \quad (1.70)$$

The matrix inversion lemma can be proved by showing that $F^{-1}F = I$. Brute force calculation of $F^{-1}F$ gives

$$\begin{aligned} F^{-1}F = I - B \left[(DA^{-1}B + C^{-1})^{-1} - C \right. \\ \left. + CDA^{-1}B(DA^{-1}B + C^{-1})^{-1} \right] DA^{-1} \end{aligned} \quad (1.71)$$

To prove the matrix inversion lemma, it is enough to show that the quantity inside the square brackets of Equation (1.71) is identically zero. Therefore, we need to prove that

$$(DA^{-1}B + C^{-1})^{-1} = C - CDA^{-1}B(DA^{-1}B + C^{-1})^{-1} \quad (1.72)$$

[†]As is evident in Chapter 2, the interpretation of P^{-1} as the *information matrix* (and P as the *covariance matrix*) hinges upon several assumptions, most notably that W_k is the inverse of the measurement error covariance.

Right multiplying both sides of Equation (1.72) by $(DA^{-1}B + C^{-1})$ reduces Equation (1.72) to

$$I = C(DA^{-1}B + C^{-1}) - CDA^{-1}B = I \tag{1.73}$$

This completes the proof.

Our next step is to apply the matrix inversion lemma to Equation (1.67). The “judicious choices” for F , A , B , C , and D are

$$F = P_{k+1} \tag{1.74a}$$

$$A = P_k^{-1} \tag{1.74b}$$

$$B = H_{k+1}^T \tag{1.74c}$$

$$C = W_{k+1} \tag{1.74d}$$

$$D = H_{k+1} \tag{1.74e}$$

The matrix information recursion now becomes

$$P_{k+1} = P_k - P_k H_{k+1}^T (H_{k+1} P_k H_{k+1}^T + W_{k+1}^{-1})^{-1} H_{k+1} P_k \tag{1.75}$$

Thus, P_{k+1} , which is used in Equation (1.66), can be obtained by “updating” P_k , and the update process usually requires inverting a matrix with rank less than n . A large number of successive applications of the recursion (1.75) occasionally introduces arithmetic errors which can invalidate the estimates (1.65). In connection with the applications of Chapter 6, alternatives to (1.75) which are numerically superior are presented.

The “update equation” (1.65) can also be rearranged in several alternate forms. One of the more common is obtained by substituting Equation (1.75) into Equation (1.66) to obtain

$$K_{k+1} = \left[P_k - P_k H_{k+1}^T (H_{k+1} P_k H_{k+1}^T + W_{k+1}^{-1})^{-1} H_{k+1} P_k \right] \times H_{k+1}^T W_{k+1} \tag{1.76a}$$

$$= P_k H_{k+1}^T \left[I - (H_{k+1} P_k H_{k+1}^T + W_{k+1}^{-1})^{-1} H_{k+1} P_k H_{k+1}^T \right] W_{k+1} \tag{1.76b}$$

Now, factoring $(H_{k+1} P_k H_{k+1}^T + W_{k+1}^{-1})^{-1}$ outside of the square brackets leads directly to

$$K_{k+1} = P_k H_{k+1}^T (H_{k+1} P_k H_{k+1}^T + W_{k+1}^{-1})^{-1} \times [W_{k+1}^{-1} + H_{k+1} P_k H_{k+1}^T - H_{k+1} P_k H_{k+1}^T] W_{k+1} \tag{1.77}$$

This leads to the *covariance recursion form*, given by

$$\hat{\mathbf{x}}_{k+1} = \hat{\mathbf{x}}_k + K_{k+1}(\tilde{\mathbf{y}}_{k+1} - H_{k+1} \hat{\mathbf{x}}_k) \tag{1.78}$$

where

$$K_{k+1} = P_k H_{k+1}^T [H_{k+1} P_k H_{k+1}^T + W_{k+1}^{-1}]^{-1} \tag{1.79}$$

$$P_{k+1} = [I - K_{k+1} H_{k+1}] P_k \tag{1.80}$$

Downloaded by [Utah State University] at 23:06 03 December 2013

The covariance form of sequential least squares is most commonly used in practice, because it is more computationally efficient. However, the information form may be numerically superior in the initialization stage. The process may be initiated at any step by an *a priori* estimate, $\hat{\mathbf{x}}_1$, and covariance estimate P_1 . If *a priori* estimates are not available, then the first data subset can be used for initialization by using a batch least squares to determine $\hat{\mathbf{x}}_q$ and P_q , where $q \geq n$. Then, the sequential least squares algorithm can be invoked for $k \geq q$. However, sequential least squares can still be used for $k = 1, 2, \dots, q - 1$ if one uses

$$P_1 = \left[\frac{1}{\alpha^2} I + H_1^T W_1 H_1 \right]^{-1} \quad (1.81)$$

$$\hat{\mathbf{x}}_1 = P_1 \left[\frac{1}{\alpha} \boldsymbol{\beta} + H_1^T W_1 \tilde{\mathbf{y}}_1 \right] \quad (1.82)$$

where α is a very “large” number and $\boldsymbol{\beta}$ is a vector of very “small” numbers. It can be shown that the resulting recursive least squares values of P_n and $\hat{\mathbf{x}}_n$ are very close to the corresponding batch values at time t_n .

If the model is in fact linear and if there is no correlation between measurement errors of different measurement subsets (so that the assumed block structure of W is strictly valid), then the sequential solution for $\hat{\mathbf{x}}$ in Equation (1.65) will agree exactly with the batch solution in Equation (1.30), to within arithmetic errors. This is because Equation (1.65) is simply an algebraic rearrangement of the normal equations (1.30).

Example 1.5: In example 1.2, we used a batch least squares process to estimate the parameters of a simple dynamic system. We now will use this same system to determine the parameters sequentially using recursive least squares with one measurement \tilde{y}_k at a time. In order to initialize the routine, we will use Equations (1.81) and (1.82) with $\alpha = 1 \times 10^3$ and $\boldsymbol{\beta} = [1 \times 10^{-2} \ 1 \times 10^{-2}]^T$. As mentioned in example 1.2, the measurement errors were simulated using a zero-mean Gaussian noise generator with a standard deviation given by $\sigma = 0.08$. We will see in Chapter 2 that an “optimal” choice for W_k is given by $W_k = \sigma^{-2}$. The calculated initial values for P_1 and $\hat{\mathbf{x}}_1$ are given by

$$P_1 = \begin{bmatrix} 1.000 \times 10^6 & 1.038 \times 10^3 \\ 1.038 \times 10^3 & 1.077 \times 10^0 \end{bmatrix}$$

$$\hat{\mathbf{x}}_1 = \begin{bmatrix} 10.010 \\ 0.014 \end{bmatrix}$$

Plots of the estimates $\hat{\mathbf{x}}_k$ and diagonal elements of P_k are shown in Figure 1.8. As can be seen from these plots, convergence is reached very quickly for this example. This is not the case in all systems, but is typical for well-conditioned linear systems. The sequential estimates at the final time agree exactly with the batch estimates in example 1.2. The diagonal elements of P_k actually have a physical meaning, as shown in Chapter 2, which can be used to develop a suitable stopping criterion. This example

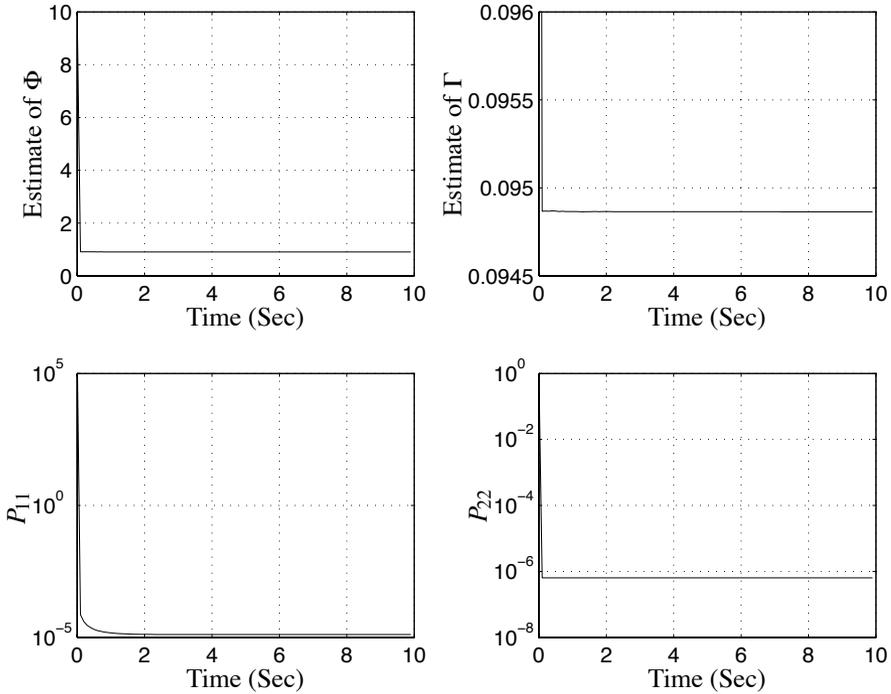


Figure 1.8: Estimates and Diagonal Elements of P_k

clearly shows the power of sequential least squares to identify the parameters of a dynamic system in *real time*.

1.4 Nonlinear Least Squares Estimation

It is a fact of life that most real-world estimation problems are nonlinear. The preceding developments of this chapter apply rigorously to only a small subset of problems encountered in practice. Fortunately, most nonlinear estimation problems can be accurately solved by a judiciously chosen successive approximation procedure. In this section we develop the most widely used successive approximation procedure, *nonlinear least squares*, otherwise known as *Gaussian least squares differential correction*. This method was originally developed by Gauss and employed to determine planetary orbits (during the early 1800s) from telescope measurements of the “line

Downloaded by [Utah State University] at 23:06 03 December 2013

of sight angles” to the planets.²

The method to be developed here is an $m \times n$ generalization of Newton’s root solving method⁸ for finding x -values satisfying $y - f(x) = 0$. As with Newton’s method, convergence of the multi-dimensional generalization is guaranteed only under rather strict requirements on the functions and their first two partial derivatives as well as on the closeness of the starting estimates. Let us not be concerned with convergence at this stage (although be informed, convergence difficulties do occasionally occur!). Rather, let us proceed with formulating the method and look at typical applications.

Assume m observable quantities modeled as

$$y_j = f_j(x_1, x_2, \dots, x_n); \quad j = 1, 2, \dots, m; \quad m \geq n \quad (1.83)$$

where the $f_j(x_1, x_2, \dots, x_n)$ are m arbitrary independent functions of the unknown parameters x_i . These should be interpreted as “functions” in the general sense, as specifying “whatever process one must go through” to compute the y_j given the x_i (including, for example, numerical solution of differential equations). We do require that $f_j(x_1, x_2, \dots, x_n)$ and at least its first partial derivatives be single-valued, continuous, and at least once differentiable. Additionally, suppose that a set of observed values of the variables y_j is available:

$$y_j \in \{y_1, y_2, \dots, y_m\} \quad (1.84)$$

As done in §1.2, we can rewrite the measurement model with Equation (1.84) in compact form as

$$\tilde{\mathbf{y}} = \mathbf{f}(\mathbf{x}) + \mathbf{v} \quad (1.85)$$

where

$$\begin{aligned} \tilde{\mathbf{y}} &= [\tilde{y}_1 \ \tilde{y}_2 \ \cdots \ \tilde{y}_m]^T = \text{measured } y\text{-values} \\ \mathbf{f}(\mathbf{x}) &= [f_1 \ f_2 \ \cdots \ f_m]^T = \text{independent functions} \\ \mathbf{x} &= [x_1 \ x_2 \ \cdots \ x_n]^T = \text{true } x\text{-values} \\ \mathbf{v} &= [v_1 \ v_2 \ \cdots \ v_m]^T = \text{measurement errors} \end{aligned}$$

Likewise, the estimated y -values, denoted by \hat{y}_j and residual errors $e_j = \tilde{y}_j - \hat{y}_j$, can also be written in compact form as

$$\hat{\mathbf{y}} = \mathbf{f}(\hat{\mathbf{x}}) \quad (1.86)$$

$$\mathbf{e} = \tilde{\mathbf{y}} - \hat{\mathbf{y}} \equiv \Delta \mathbf{y} \quad (1.87)$$

where

$$\begin{aligned} \hat{\mathbf{y}} &= [\hat{y}_1 \ \hat{y}_2 \ \cdots \ \hat{y}_m]^T = \text{estimated } y\text{-values} \\ \mathbf{e} &= [e_1 \ e_2 \ \cdots \ e_m]^T = \text{residual errors} \\ \hat{\mathbf{x}} &= [\hat{x}_1 \ \hat{x}_2 \ \cdots \ \hat{x}_n]^T = \text{estimated } x\text{-values} \end{aligned}$$

The measurement model in Equation (1.86) can again be written using the residual errors \mathbf{e} as

$$\tilde{\mathbf{y}} = \mathbf{f}(\hat{\mathbf{x}}) + \mathbf{e} \tag{1.88}$$

As done in §1.2, we seek an estimate $(\hat{\mathbf{x}})$ for \mathbf{x} that minimizes

$$J = \frac{1}{2} \mathbf{e}^T W \mathbf{e} = \frac{1}{2} [\tilde{\mathbf{y}} - \mathbf{f}(\hat{\mathbf{x}})]^T W [\tilde{\mathbf{y}} - \mathbf{f}(\hat{\mathbf{x}})] \tag{1.89}$$

where W is an $m \times m$ weighting matrix again used to weight the relative importance of each measurement.

In most practical problems, J cannot be directly minimized by application of ordinary calculus to Equation (1.89), in the sense that explicit closed form solutions for $\hat{\mathbf{x}}$ result. The case where $\mathbf{f}(\hat{\mathbf{x}}) = H\hat{\mathbf{x}}$ reduces to the standard linear least squares solution; however, general nonlinear functions for $\mathbf{f}(\hat{\mathbf{x}})$ typically make the solution difficult to find explicitly. For this reason, attention is directed to construction of a successive approximation procedure due to Gauss, that is designed to converge to accurate least squares estimates, given approximate starting values (the determination of sufficiently close starting estimates is a problem that cannot be dealt with in general, but can usually be overcome, as seen in applications of Chapter 6 and in §1.6.3).

Assume that the *current* estimates of the unknown \mathbf{x} -values are available, denoted by

$$\mathbf{x}_c = [x_{1c} \ x_{2c} \ \cdots \ x_{nc}]^T \tag{1.90}$$

Whatever the unknown objective \mathbf{x} -values $\hat{\mathbf{x}}$ are, we assume that they are related to their respective current estimates, \mathbf{x}_c , by an also unknown set of corrections, $\Delta\mathbf{x}$, as

$$\hat{\mathbf{x}} = \mathbf{x}_c + \Delta\mathbf{x} \tag{1.91}$$

If the components of $\Delta\mathbf{x}$ are sufficiently small, it may be possible to solve for approximations to them and thereby update \mathbf{x}_c with an improved estimate of \mathbf{x} from Equation (1.91). With this assumption, we may *linearize* $\mathbf{f}(\hat{\mathbf{x}})$ in Equation (1.86) about \mathbf{x}_c using a first-order Taylor series expansion as

$$\mathbf{f}(\hat{\mathbf{x}}) \approx \mathbf{f}(\mathbf{x}_c) + H\Delta\mathbf{x} \tag{1.92}$$

where

$$H \equiv \left. \frac{\partial \mathbf{f}}{\partial \mathbf{x}} \right|_{\mathbf{x}_c} \tag{1.93}$$

The gradient matrix H is known as a *Jacobian* matrix (see Appendix B). The measurement residual “after the correction” can now be linearly approximated as

$$\Delta\mathbf{y} \equiv \tilde{\mathbf{y}} - \mathbf{f}(\hat{\mathbf{x}}) \approx \tilde{\mathbf{y}} - \mathbf{f}(\mathbf{x}_c) - H\Delta\mathbf{x} = \Delta\mathbf{y}_c - H\Delta\mathbf{x} \tag{1.94}$$

where the residual “before the correction” is

$$\Delta\mathbf{y}_c \equiv \tilde{\mathbf{y}} - \mathbf{f}(\mathbf{x}_c) \tag{1.95}$$

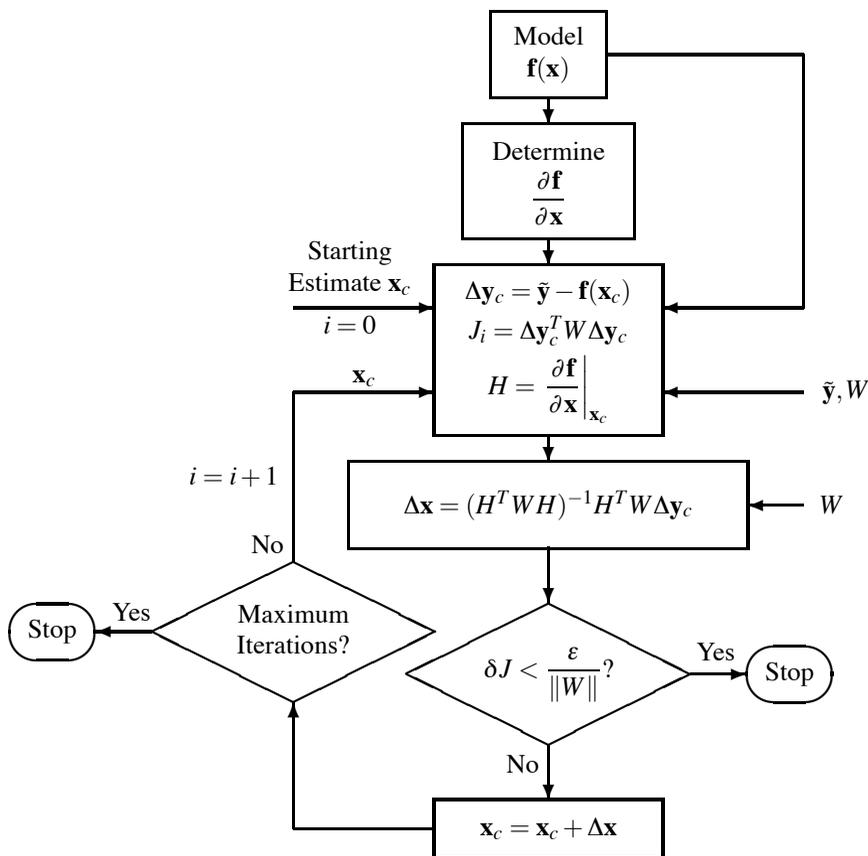


Figure 1.9: Nonlinear Least Squares Algorithm

Recall that the objective is to minimize the weighted sum squares, J , given by Equation (1.89). The local strategy for determining the approximate corrections (“differential corrections”) in $\Delta \mathbf{x}$ is to select the particular corrections that lead to the *minimum sum of squares of the linearly predicted residuals* J_p :

$$J = \frac{1}{2} \Delta \mathbf{y}^T W \Delta \mathbf{y} \approx J_p \equiv \frac{1}{2} (\Delta \mathbf{y}_c - H \Delta \mathbf{x})^T W (\Delta \mathbf{y}_c - H \Delta \mathbf{x}) \quad (1.96)$$

Before carrying out the minimization, we note (to the approximation that the linearization implicit in the prediction (1.92) is valid) that the minimization of J_p in Equation (1.96) is equivalent to the minimization of J in Equation (1.89). If the process is convergent, then $\Delta \mathbf{x}$ determined by minimizing Equation (1.96) would be expected to decrease on successive iterations until (on the final iteration) the linearization is an extremely good approximation.

Observe that the minimization of Equation (1.96) is completely analogous to the previously minimized quadratic form (1.27). Thus, any algorithm for solving the

weighted least squares problem directly applies to solving for $\Delta \mathbf{x}$ in Equation (1.96). Therefore, the appropriate version of the normal equations follows as in the development of Equations (1.28)-(1.30), as

$$\Delta \mathbf{x} = (H^T W H)^{-1} H^T W \Delta \mathbf{y}_c \tag{1.97}$$

The complete nonlinear least squares algorithm is summarized in Figure 1.9. An initial guess \mathbf{x}_c is required to begin the algorithm. Equation (1.97) is then calculated using the residual measurements ($\Delta \mathbf{y}_c$), Jacobian matrix (H), and weighting matrix (W), so that the current estimate can be updated. A stopping condition with an accuracy dependent tolerance for the minimization of J is given by

$$\delta J \equiv \frac{|J_i - J_{i-1}|}{J_i} < \frac{\epsilon}{\|W\|} \tag{1.98}$$

where ϵ is a prescribed small value. If Equation (1.98) is not satisfied, then the update procedure is iterated with the new estimate as the current estimate until the process converges, or unsatisfactory convergence progress is evident (e.g., a maximum allowed number of iterations is exceeded, or J increases on successive iterations).

The above least squares differential correction process, while far from fail-safe, has been successfully applied to an extremely wide variety of nonlinear estimation problems. Convergence difficulties usually stem from one of the following sources: (1) the initial \mathbf{x} -estimate is too far from the minimizing $\hat{\mathbf{x}}$ (for the nonlinearity of the particular application), resulting in the implicit local linearity assumption being invalid; (2) numerical difficulties are encountered in solving for the corrections, $\Delta \mathbf{x}$, due to (2a) arithmetic errors corrupting the particular algorithm used to calculate the $\Delta \mathbf{x}$, or (2b) the H matrix having fewer than n linearly independent rows or columns (i.e., rank deficient). The difficulties (1) and (2a) can usually be overcome by a resourceful analyst; however, the least squares criterion does not uniquely define $\Delta \mathbf{x}$ in the (2b) case, and therefore some other criterion must be employed to select $\Delta \mathbf{x}$. The initial estimate convergence difficulty can also be overcome by using the Levenberg-Marquardt algorithm shown in §1.6.3, which combines the least squares differential correction process with a gradient search.

Example 1.6: In this simple example, we consider the 1×1 special case of nonlinear least squares with $m = n = 1$. Suppose we have the following model:

$$y = x^3 + 6x^2 + 11x + 6 = 0$$

For this model, we can assume that

$$\begin{aligned} \mathbf{y} &= y = 0 \\ \mathbf{f}(\mathbf{x}) &= f(x) = x^3 + 6x^2 + 11x + 6 \end{aligned}$$

For this case, Equation (1.97) becomes simply

$$x = x_c - \left[\left. \frac{\partial f}{\partial x} \right|_{x_c} \right]^{-1} f(x_c)$$

where

$$\frac{\partial f}{\partial x} = 3x^2 + 12x + 11$$

As seen in the above equations, this special scalar case reduces to the classical Newton root solving method. Therefore, Equation (1.97) actually represents an $m \times n$ generalization of Newton's root solver. Seven iterations for three different starting values of x are given below.

iteration	x	x	x
0	0.0000	-1.6000	-5.0000
1	-0.5455	-2.2462	-4.0769
2	-0.8490	-1.9635	-3.5006
3	-0.9747	-2.0001	-3.1742
4	-0.9991	-2.0000	-3.0324
5	-1.0000	-2.0000	-3.0015
6	-1.0000	-2.0000	-3.0000
7	-1.0000	-2.0000	-3.0000

This clearly shows that different solutions are possible for various starting conditions. In this case, we know this to be true since we are solving a cubic equation, which has three possible solutions, and obviously, we have converged to all three roots. More generally, complex algebra would have to be used to find complex roots.

Example 1.7: In example 1.2, we used linear least squares to estimate the parameters of a simple dynamic system. Recall that the system is given by

$$y_{k+1} = \left[e^{a\Delta t} \right] y_k + \left[\frac{b}{a} (e^{a\Delta t} - 1) \right] u_k$$

Suppose that we now wish to determine a and b directly from the above equation. To accomplish this task, we must now use nonlinear least squares, with

$$\begin{aligned} \mathbf{x} &= [a \ b]^T \\ \tilde{\mathbf{y}} &= [\tilde{y}_2 \ \tilde{y}_3 \ \cdots \ \tilde{y}_{101}]^T \\ f_k &= \left[e^{a\Delta t} \right] y_k + \left[\frac{b}{a} (e^{a\Delta t} - 1) \right] u_k \end{aligned}$$

The appropriate partials are given by

$$\frac{\partial f_k}{\partial a} = \Delta t \left[e^{a\Delta t} \right] y_k + \left[\frac{b}{a^2} (1 - e^{a\Delta t}) + \frac{b}{a} \Delta t e^{a\Delta t} \right] u_k$$

$$\frac{\partial f_k}{\partial b} = \frac{1}{a} (e^{a\Delta t} - 1) u_k$$

Then, the H matrix is given by

$$H = \begin{bmatrix} \Delta t [e^{a\Delta t}] \tilde{y}_1 + \left[\frac{b}{a^2}(1 - e^{a\Delta t}) + \frac{b}{a}\Delta t e^{a\Delta t} \right] u_1 & \frac{1}{a}(e^{a\Delta t} - 1)u_1 \\ \Delta t [e^{a\Delta t}] \tilde{y}_2 + \left[\frac{b}{a^2}(1 - e^{a\Delta t}) + \frac{b}{a}\Delta t e^{a\Delta t} \right] u_2 & \frac{1}{a}(e^{a\Delta t} - 1)u_2 \\ \vdots & \vdots \\ \Delta t [e^{a\Delta t}] \tilde{y}_{100} + \left[\frac{b}{a^2}(1 - e^{a\Delta t}) + \frac{b}{a}\Delta t e^{a\Delta t} \right] u_{100} & \frac{1}{a}(e^{a\Delta t} - 1)u_{100} \end{bmatrix}$$

The nonlinear least squares algorithm in Figure 1.9 can now be used to determine a and b . The starting guess for the iteration is given by

$$\mathbf{x}_c = [5 \ 5]^T$$

Also, the stopping criterion is given by $\varepsilon = 1 \times 10^{-8}$. Results are tabulated below.

iteration	\hat{a}	\hat{b}
0	5.0000	5.0000
1	0.4876	1.9540
2	-0.8954	1.0634
3	-1.0003	0.9988
4	-1.0009	0.9985
5	-1.0009	0.9985
6	-1.0009	0.9985

If we convert the final values for \hat{a} and \hat{b} into their discrete time equivalents, we see that $\hat{\Phi} = 0.9048$ and $\hat{\Gamma} = 0.0950$, which agree with the results obtained in example 1.2. This example clearly shows that the *form* of the model chosen can have a highly significant impact on the complexity of the required estimator. If we choose to determine Φ and Γ directly, then *linear* least squares may be employed. However, if we choose to determine a and b , then nonlinear least squares must be used. Clearly, by using creative system model choices, one can greatly simplify the overall solution process. This point is further explored in §1.5 and in Chapter 6.

Example 1.8: Under certain approximations, the pitch (θ) and yaw (ψ) attitude dynamics of an inertially and aerodynamically symmetric projectile can be modeled via a pair of equations

$$\begin{aligned} \theta(t) &= k_1 e^{\lambda_1 t} \cos(\omega_1 t + \delta_1) + k_2 e^{\lambda_2 t} \cos(\omega_2 t + \delta_2) \\ &\quad + k_3 e^{\lambda_3 t} \cos(\omega_3 t + \delta_3) + k_4 \\ \psi(t) &= k_1 e^{\lambda_1 t} \sin(\omega_1 t + \delta_1) + k_2 e^{\lambda_2 t} \sin(\omega_2 t + \delta_2) \\ &\quad + k_3 e^{\lambda_3 t} \sin(\omega_3 t + \delta_3) + k_5 \end{aligned}$$

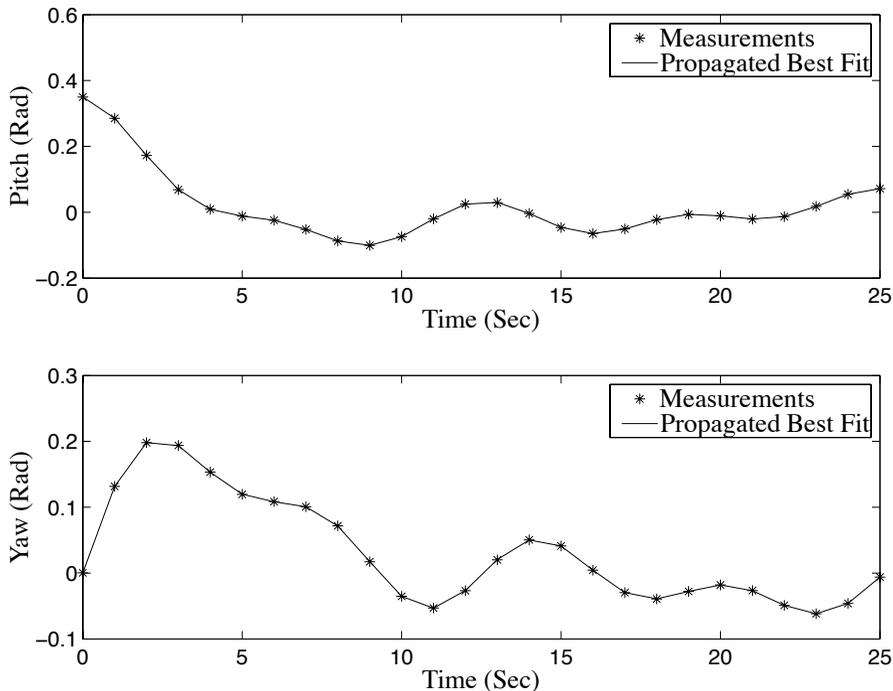


Figure 1.10: Simulated Pitch and Yaw Measurements and Best Fits

where $k_1, k_2, k_3, k_4, k_5, \lambda_1, \lambda_2, \lambda_3, \omega_1, \omega_2, \omega_3, \delta_1, \delta_2, \delta_3$ are 14 constants which can be related to the aerodynamic and mass characteristics of the projectile and to the initial motion conditions. These constants are often estimated by nonlinear least squares to “best fit” measured pitch and yaw histories modeled by the above equations.

As an example of such a data reduction process, consider the simulated measurements of $\theta(t)$ and $\psi(t)$ with the measurement error generated by using a zero-mean Gaussian noise process with a standard deviation given by $\sigma = 0.0002$. The measurements are sampled at 1 sec intervals, shown in Figure 1.10. The *a priori* constant estimates and true values are given by

Constant Parameter	Start Value	True Value
k_1	0.5000	0.2000
k_2	0.2500	0.1000
k_3	0.1250	0.0500
k_4	0.0000	0.0001
k_5	0.0000	0.0001
λ_1	-0.1500	-0.1000
λ_2	-0.0600	-0.0500
λ_3	-0.0300	-0.0250
ω_1	0.2600	0.2500
ω_2	0.5500	0.5000
ω_3	0.9500	1.0000
δ_1	0.0100	0.0000
δ_2	0.0100	0.0000
δ_3	0.0100	0.0000

For the problem at hand the necessary conditions in Equation (1.97) are defined as

$$\mathbf{x}^{(14 \times 1)} = [k_1 \ k_2 \ k_3 \ k_4 \ k_5 \ \lambda_1 \ \lambda_2 \ \lambda_3 \ \omega_1 \ \omega_2 \ \omega_3 \ \delta_1 \ \delta_2 \ \delta_3]^T$$

$$\mathbf{y}^{(52 \times 1)} = [\tilde{\theta}(0) \ \tilde{\psi}(0) \ \tilde{\theta}(1) \ \tilde{\psi}(1) \ \dots \ \tilde{\theta}(25) \ \tilde{\psi}(25)]^T$$

$$H^{(52 \times 14)} = \begin{bmatrix} \left. \frac{\partial \theta(0)}{\partial x_1} \right|_{\mathbf{x}_c} & \dots & \left. \frac{\partial \theta(0)}{\partial x_{14}} \right|_{\mathbf{x}_c} \\ \left. \frac{\partial \psi(0)}{\partial x_1} \right|_{\mathbf{x}_c} & \dots & \left. \frac{\partial \psi(0)}{\partial x_{14}} \right|_{\mathbf{x}_c} \\ \vdots & & \vdots \\ \left. \frac{\partial \theta(25)}{\partial x_1} \right|_{\mathbf{x}_c} & \dots & \left. \frac{\partial \theta(25)}{\partial x_{14}} \right|_{\mathbf{x}_c} \\ \left. \frac{\partial \psi(25)}{\partial x_1} \right|_{\mathbf{x}_c} & \dots & \left. \frac{\partial \psi(25)}{\partial x_{14}} \right|_{\mathbf{x}_c} \end{bmatrix}$$

$$W^{(52 \times 52)} = 10^8 \begin{bmatrix} 0.25 & & 0 \\ & 0.25 & \\ & & \ddots \\ 0 & & & 0.25 \end{bmatrix}$$

and the 28 partial derivative expressions (needed to fill the H -matrix) are given by

$$\frac{\partial \theta(t_j)}{\partial k_i} = e^{\lambda_i t_j} \cos(\omega_i t_j + \delta_i), \quad i = 1, 2, 3$$

$$\frac{\partial \psi(t_j)}{\partial k_i} = e^{\lambda_i t_j} \sin(\omega_i t_j + \delta_i), \quad i = 1, 2, 3$$

$$\frac{\partial \theta(t_j)}{\partial k_4} = 1, \quad \frac{\partial \psi(t_j)}{\partial k_4} = 0, \quad \frac{\partial \theta(t_j)}{\partial k_5} = 0, \quad \frac{\partial \psi(t_j)}{\partial k_5} = 1$$

$$\frac{\partial \theta(t_j)}{\partial \lambda_i} = t_j k_i e^{\lambda_i t_j} \cos(\omega_i t_j + \delta_i), \quad i = 1, 2, 3$$

$$\frac{\partial \psi(t_j)}{\partial \lambda_i} = t_j k_i e^{\lambda_i t_j} \sin(\omega_i t_j + \delta_i), \quad i = 1, 2, 3$$

$$\frac{\partial \theta(t_j)}{\partial \omega_i} = -t_j k_i e^{\lambda_i t_j} \sin(\omega_i t_j + \delta_i), \quad i = 1, 2, 3$$

$$\frac{\partial \psi(t_j)}{\partial \omega_i} = t_j k_i e^{\lambda_i t_j} \cos(\omega_i t_j + \delta_i), \quad i = 1, 2, 3$$

$$\frac{\partial \theta(t_j)}{\partial \delta_i} = -k_i e^{\lambda_i t_j} \sin(\omega_i t_j + \delta_i), \quad i = 1, 2, 3$$

$$\frac{\partial \psi(t_j)}{\partial \delta_i} = k_i e^{\lambda_i t_j} \cos(\omega_i t_j + \delta_i), \quad i = 1, 2, 3$$

Results in the convergence history are summarized below.

Parameter	Iteration Number					σ
	0	1	2	...	5	
k_1	0.5000	0.1852	0.1975		0.1999	0.0006
k_2	0.2500	0.1075	0.1012		0.0997	0.0005
k_3	0.1250	0.0567	0.0505		0.0500	0.0001
k_4	0.0000	-0.0006	0.0001		0.0002	0.0001
k_5	0.0000	-0.0018	-0.0005		0.0001	0.0001
λ_1	-0.1500	-0.1234	-0.0954		-0.0998	0.0004
λ_2	-0.0600	-0.0661	-0.0585		-0.0497	0.0004
λ_3	-0.0300	-0.0398	-0.0338		-0.0250	0.0002
ω_1	0.2600	0.2490	0.2471		0.2500	0.0004
ω_2	0.5500	0.5300	0.4955		0.4999	0.0004
ω_3	0.9500	0.9697	1.0068		0.9998	0.0002
δ_1	0.0100	0.0344	0.0143		0.0010	0.0031
δ_2	0.0100	-0.0447	0.0051		0.0001	0.0048
δ_3	0.0100	0.0024	-0.0570		-0.0001	0.0024

Observe the rather dramatic convergence progress shown in the results. The right-most column is obtained by taking the square root of the 14 diagonal elements of $(H^TWH)^{-1}$ on the final iteration. We prove this interpretation of $(H^TWH)^{-1}$ in Chapter 2. Thus, a by-product of the least squares algorithm is an uncertainty measure of the answer! Note that the convergence errors are comparable in size to the corresponding σ . Also, for this example the weighted sum square of residuals (i.e., the value of J) at each iteration is given by

Cost	Iteration Number				
	0	1	2	...	5
J	1.08×10^7	2.51×10^5	1.17×10^4		1.93×10^1

Clearly, the dramatic convergence is evidenced by the decrease of the weighted sum square of the residuals by six orders of magnitude in five iterations. Also, observe that the final converged values of the fifth iteration are in reasonable agreement with their respective true values.

1.5 Basis Functions

This section gives an overview of some common basis functions used in least squares. Although the discussion here is not exhaustive, it will serve to introduce the subject matter. As seen in previous examples from this chapter, various basis functions have been used to identify system parameters. How to choose these basis functions usually comes from experience and knowledge of the particular dynamic system under investigation. Still, some commonly used basis functions can be used for a wide variety of systems. A very common choice for the linearly independent basis functions (1.12) are the powers of t :

$$\{1, t, t^2, t^3, \dots\} \tag{1.99}$$

in which case the model (1.11) is a power series polynomial

$$y(t) = x_1 + x_2t + x_3t^2 + \dots = \sum_{i=1}^n x_i t^{i-1} \tag{1.100}$$

The least squares coefficients estimates then follow from Equation (1.26) with the coefficient matrix

$$H = \begin{bmatrix} 1 & t_1 & t_1^2 & \dots & t_1^{n-1} \\ 1 & t_2 & t_2^2 & \dots & t_2^{n-1} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & t_m & t_m^2 & \dots & t_m^{n-1} \end{bmatrix} \tag{1.101}$$

known as the *Vandermonde matrix*.^{7,9} Often, one encounters a nonlinear system where the basis functions are not polynomials. However, through a change of variables, one may be able to transform the original basis functions into powers of t .¹⁰ Examples of such a change are given in [Table 1.1](#).

Downloaded by [Utah State University] at 23:06 03 December 2013

Table 1.1: Change of Variables into Powers of t

Basis Function	New Form	Change of Variables
$y = x_1 + \frac{x_2}{a} + \frac{x_3}{a^2} + \dots$	$y = x_1 + x_2t + x_3t^2 + \dots$	$t = \frac{1}{a}, a \neq 0$
$y = Be^{at}$	$z = x_1 + x_2t$	$z = \ln y, y > 0$ $x_1 = \ln B, B > 0$ $x_2 = a$
$y = x_1w^{-m} + x_2w^n$	$z = x_1 + x_2t$	$z = yw^m$ $t = w^{m+n}$
$y = B \exp \left[-\frac{(1-at)^2}{2\sigma^2} \right]$	$z = x_1 + x_2t + x_3t^2$	$z = \ln y, y > 0$ $x_1 = \ln B - \frac{\ln e}{2\sigma^2}, B > 0$ $x_2 = \frac{a \ln e}{\sigma^2}$ $x_3 = -\frac{\ln e}{2\sigma^2} a^2$

Therefore, linear least squares may often be used to determine the parameters that appear to be nonlinear in nature. Through judicious change of variables, a linear solution is now possible. But one must take care because singular conditions may arise by the change of variables. For example, using the change of variables approach for $y = Be^{at}$ shown in Table 1.1 creates a singular condition when B is negative. Note that the Vandermonde matrix may have numerical problems due to ill-conditioning for $n > 10$, but this headache may be partially overcome by using least squares matrix decompositions, which are discussed in §1.6.1.

Another common choice for the linearly independent basis functions (1.12) are harmonic series, which can be used to approximate y :

$$\begin{aligned}
 y_j &= a_0 + a_1 \cos(\omega t_j) + b_1 \sin(\omega t_j) + \dots \\
 &\quad + a_n \cos(n\omega t_j) + b_n \sin(n\omega t_j), \tag{1.102} \\
 &j = 1, \dots, m; \quad m \geq 2n + 1
 \end{aligned}$$

where the amplitudes (a_i, b_i) are the sought parameters. Suppose we are given $\bar{y}_j, t_j, W = (W_{ij})$, and $\omega = 2\pi/T$, where T is the period under consideration. Then, the

desired least squares estimate (\hat{a}_i, \hat{b}_i) is computable as

$$\hat{\mathbf{x}} = \begin{bmatrix} \hat{a}_0 \\ \hat{a}_1 \\ \hat{b}_1 \\ \vdots \\ \hat{a}_n \\ \hat{b}_n \end{bmatrix} = (H^T W H)^{-1} H^T W \tilde{\mathbf{y}} \tag{1.103}$$

where

$$H = \begin{bmatrix} 1 & \cos(\omega t_1) & \sin(\omega t_1) & \cdots & \cos(n\omega t_1) & \sin(n\omega t_1) \\ 1 & \cos(\omega t_2) & \sin(\omega t_2) & \cdots & \cos(n\omega t_2) & \sin(n\omega t_2) \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ 1 & \cos(\omega t_m) & \sin(\omega t_m) & \cdots & \cos(n\omega t_m) & \sin(n\omega t_m) \end{bmatrix} \tag{1.104}$$

In the case above, if W is chosen as an identity matrix and the sample points $\{t_1, t_2, \dots\}$ are chosen such that the off-diagonal elements of $(H^T W H)$ vanish, then the least squares solution is reduced to its most elegant form. This leads to a simple solution, given by

$$\hat{x}_i = \left[\sum_{j=1}^m h_i^2(t_j) \right]^{-1} \sum_{j=1}^m h_i(t_j) \tilde{y}_j, \quad i = 1, 2, \dots, n \tag{1.105}$$

where

$$\begin{aligned} \mathbf{h}(t) &\equiv [h_1(t) \ h_2(t) \ h_3(t) \ \cdots]^T \\ &= [1 \ \cos(\omega t) \ \sin(\omega t) \ \cdots \ \cos(n\omega t) \ \sin(n\omega t)]^T \end{aligned} \tag{1.106}$$

A significant advantage of the uncoupled solution for the coefficients in Equation (1.105) is that adding another $(n + 1)$ basis function (which has the same form as any of the first n) does not affect the first n solutions for \hat{x}_i .

The least squares estimate for the coefficients has a strong connection to the continuous approximation for $\tilde{y}(t)$. Before we formally prove this, let us review the concept of an *orthogonal* set of functions.^{11, 12} An infinite system of real functions

$$\{\varphi_1(t), \varphi_2(t), \varphi_3(t), \dots, \varphi_n(t), \dots\} \tag{1.107}$$

is said to be orthogonal on the interval $[\alpha, \beta]$ if

$$\int_{\alpha}^{\beta} \varphi_p(t) \varphi_q(t) dt = 0 \quad (p \neq q, p, q = 1, 2, 3, \dots) \tag{1.108}$$

and

$$\int_{\alpha}^{\beta} \varphi_p^2(t) dt \equiv c_p \neq 0 \quad (p = 1, 2, 3, \dots) \tag{1.109}$$

Downloaded by [Utah State University] at 23:06 03 December 2013

The series given in Equation (1.106) can be shown to be orthogonal over any interval centered on $t = T/2$. We further note the distinction between the continuous orthogonality conditions of Equations (1.108) and the corresponding discrete orthogonality conditions

$$\sum_{j=1}^m \varphi_p(t_j) \varphi_q(t_j) = c_p \delta_{pq} \quad (1.110)$$

where the Kronecker delta δ_{pq} is defined as

$$\begin{aligned} \delta_{pq} &= 0 & \text{if } p \neq q \\ &= 1 & \text{if } p = q \end{aligned} \quad (1.111)$$

For the discrete orthogonality case, a specific pattern of sample points underlies this condition. We also mention that the most general forms of the continuous and discrete orthogonality conditions are

$$\int_{\alpha}^{\beta} w(t) \varphi_p(t) \varphi_q(t) dt = c_p \delta_{pq} \quad (1.112)$$

and

$$\sum_{j=1}^m w(t_j) \varphi_p(t_j) \varphi_q(t_j) = c_p \delta_{pq} \quad (1.113)$$

where $w(t)$ is an associated weight function.

The orthogonality condition on the individual integrals of the terms $\sin(2\pi pt/T)$ and $\cos(2\pi pt/T)$ are trivial to prove on the interval $[0, T]$. A slightly more complex case involves the integral of $\sin(ct) \sin(dt)$ for any $c \neq d$ on the interval $[0, T]$:

$$\begin{aligned} \int_0^T \sin(ct) \sin(dt) dt &= \frac{1}{2} \int_0^T [\cos(ct - dt) - \cos(ct + dt)] dt \\ &= \left[\frac{\sin(ct - dt)}{2(c - d)} - \frac{\sin(ct + dt)}{2(c + d)} \right] \Big|_0^T \end{aligned} \quad (1.114)$$

If we let $c = 2\pi p/T$ and $d = 2\pi q/T$, then it is easy to see that Equation (1.114) is identically zero for any $p \neq q$. Therefore, this system is orthogonal with the associated weight function $w(t) = 1$. It can also be shown that all integrals of any combinations of the functions in Equation (1.106) are orthogonal on the interval $[0, T]$. Of course, we may also replace the integral with a summation; for symmetrically located samples, we have discrete orthogonality and this leads directly to the solution in Equation (1.105).

The *Fourier series* of a function is a harmonic expansion of sines and cosines, given by

$$y(t) = a_0 + \sum_{n=1}^{\infty} a_n \cos(n\omega t) + \sum_{n=1}^{\infty} b_n \sin(n\omega t) \quad (1.115)$$

To compute a coefficient such as a_1 , multiply both sides of Equation (1.115) by $\cos(\omega t)$ and integrate from 0 to T (the function y is given on this interval). This

leads to

$$\int_0^T y(t) \cos(\omega t) dt = a_0 \int_0^T \cos(\omega t) dt + a_1 \int_0^T [\cos(\omega t)]^2 dt + \dots + b_1 \int_0^T \cos(\omega t) \sin(\omega t) dt + \dots \tag{1.116}$$

Every integral on the right-hand side of Equation (1.116) is zero (since the sines and cosines are mutually orthogonal) except the one in which $\cos(\omega t)$ multiplies itself. Therefore, a_1 is given by

$$a_1 = \frac{\int_0^T y(t) \cos(\omega t) dt}{\int_0^T [\cos(\omega t)]^2 dt} \tag{1.117}$$

The coefficient b_1 would have $\sin(\omega t)$ in place of $\cos(\omega t)$, and b_2 would use $\sin(2\omega t)$, and so on. Evaluating the integral in the denominator of Equation (1.117) and likewise for the other coefficients leads to the *Fourier coefficients*,^{13, 14} given by

$$a_0 = \frac{1}{T} \int_0^T y(t) dt \tag{1.118a}$$

$$a_n = \frac{2}{T} \int_0^T y(t) \cos(n\omega t) dt \tag{1.118b}$$

$$b_n = \frac{2}{T} \int_0^T y(t) \sin(n\omega t) dt \tag{1.118c}$$

The Fourier coefficients can also be determined using linear least squares, and in the process, we establish that the determined coefficients are simply a special case of least squares approximation. For this development we will assume that our measurement model, $\tilde{y}(t)$, is given by Equation (1.115), so that $\tilde{y}(t) = y(t)$. Consider minimizing the following function:

$$J = \frac{1}{2} \int_0^T [y(t) - \hat{\mathbf{x}}^T \mathbf{h}(t)]^T [y(t) - \hat{\mathbf{x}}^T \mathbf{h}(t)] dt \tag{1.119}$$

or

$$J = \frac{1}{2} \int_0^T [y(t)]^2 dt - \left[\int_0^T y(t) \mathbf{h}^T(t) dt \right] \hat{\mathbf{x}} + \frac{1}{2} \hat{\mathbf{x}}^T \left[\int_0^T \mathbf{h}(t) \mathbf{h}^T(t) dt \right] \hat{\mathbf{x}} \tag{1.120}$$

The necessary condition $\nabla_{\hat{\mathbf{x}}} J = \mathbf{0}$ leads to

$$\hat{\mathbf{x}} = \left[\int_0^T \mathbf{h}(t) \mathbf{h}^T(t) dt \right]^{-1} \left[\int_0^T y(t) \mathbf{h}(t) dt \right] \tag{1.121}$$

Since $\mathbf{h}(t)$ represents a set of orthogonal functions on the interval $[0, T]$, i.e., the functions satisfy Equations (1.108) and (1.109), so that $\int_0^T \mathbf{h}(t) \mathbf{h}^T(t) dt$ is a diagonal

matrix with elements given by $\int_0^T [h_i(t)]^2 dt$, then the individual components of $\hat{\mathbf{x}}$ are simply given by the uncoupled equations

$$\hat{x}_i = \frac{\int_0^T y(t)h_i(t) dt}{\int_0^T [h_i(t)]^2 dt}, \quad i = 1, 2, \dots, n \quad (1.122)$$

This is identical to the solution shown in Equation (1.118). Therefore, the Fourier coefficients are just “least square” estimates using the particular orthogonal basis function in Equation (1.106). On several occasions herein, we will make use of orthogonal basis functions; however, this subject is not treated comprehensively within the scope of this text. Most standard mathematical handbooks, such as Abramowitz and Stegun,¹⁵ and Ledermann,¹⁶ summarize a large family of orthogonal polynomials and discuss their use in approximation.

1.6 Advanced Topics

In this section we will show some advanced topics used in least squares. Although an exhaustive treatment is beyond the scope of this text, we hope that the subjects presented herein will motivate the interested reader to pursue them in the referenced literature.

1.6.1 Matrix Decompositions in Least Squares

The core component of any least squares algorithm is $(H^T H)^{-1}$. As an alternative to direct computation of this inverse, it is common to decompose H in some way which simplifies the calculations and/or is more robust with respect to near singularity conditions. A more detailed mathematical development of some of the topics presented here is provided in §B.4.

A particularly useful decomposition of the matrix H is the QR decomposition. Before we discuss this decomposition, let us first review the definition and properties of orthogonal vectors and matrices. Two vectors, \mathbf{u} and \mathbf{v} , are *orthogonal* if the angle between them is $\pi/2$. This can be true if and only if $\mathbf{u}^T \mathbf{v} = 0$. An *orthogonal matrix*^{7, 17} Q is a square matrix with *orthonormal* column vectors. Orthonormal vectors are orthogonal vectors each with unit lengths. Since the columns of an orthogonal matrix Q are orthonormal, then $Q^T Q = I$ (where $Q^T Q$ is a matrix of vector-space inner-products) and $Q^T = Q^{-1}$. This clearly shows that the inverse of an orthogonal matrix is given by its transpose!

An example of an orthogonal matrix in dynamic systems is the *rotation* matrix. For example, let

$$Q = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos \phi & \sin \phi \\ 0 & -\sin \phi & \cos \phi \end{bmatrix} \quad (1.123)$$

This matrix is clearly orthogonal, since the column vectors are orthonormal.

The QR decomposition factors a full rank matrix H as the product of an orthogonal matrix Q and an upper-triangular matrix R , given by

$$H = QR \tag{1.124}$$

where Q is an $m \times n$ matrix with $Q^T Q = I$, and R is an upper triangular $n \times n$ matrix with all elements $R_{ij} = 0$ for $i > j$. The QR decomposition can be accomplished using the *modified Gram-Schmidt* algorithm (see §B.4). The advantage of the QR decomposition is that it greatly simplifies the least squares problem. The term $H^T H$ in the normal equations is easier to invert since

$$H^T H = R^T Q^T QR = R^T R \tag{1.125}$$

Therefore, the normal equations (1.26) simplify to

$$R^T R \hat{\mathbf{x}} = R^T Q^T \tilde{\mathbf{y}} \tag{1.126}$$

or

$$\boxed{R \hat{\mathbf{x}} = Q^T \tilde{\mathbf{y}}} \tag{1.127}$$

The solution to Equation (1.127) can easily be accomplished since R is upper triangular (see Appendix B). The real cost is in the $2mn^2$ operations in the modified Gram-Schmidt algorithm, which are required to compute Q and R . The QR decomposition can also be used in linear least squares to improve an approximate solution using iterative refinement.¹⁸ Notice it is not necessary to square H (i.e., form $H^T H$); the QR algorithm operates directly on H . If H is poorly conditioned, it is easy to verify that $H^T H$ is much more poorly conditioned than H itself.

Another decomposition of the matrix H is the *singular-value decomposition*,^{7,17} which decomposes a matrix into a diagonal matrix and two orthogonal matrices:

$$H = USV^T \tag{1.128}$$

where U is the $m \times n$ matrix with orthonormal columns, S is an $n \times n$ diagonal matrix such that $S_{ij} = 0$ for $i \neq j$, and V is an $n \times n$ orthogonal matrix. Note that $U^T U = I$, but it is no longer possible to make the same statement for $U U^T$. Now, substitute Equation (1.128) into Equation (1.25):

$$(H^T H) \hat{\mathbf{x}} = H^T \tilde{\mathbf{y}} \tag{1.129a}$$

$$(VSU^T USV^T) \hat{\mathbf{x}} = VSU^T \tilde{\mathbf{y}} \tag{1.129b}$$

$$(VSSV^T) \hat{\mathbf{x}} = VSU^T \tilde{\mathbf{y}} \tag{1.129c}$$

$$(SV^T) \hat{\mathbf{x}} = U^T \tilde{\mathbf{y}} \tag{1.129d}$$

Therefore, the solution for $\hat{\mathbf{x}}$ is simply given by

$$\boxed{\hat{\mathbf{x}} = VS^{-1}U^T \tilde{\mathbf{y}}} \tag{1.130}$$

Note that the inverse of S is easy to compute since it is a diagonal matrix (i.e., $S = \text{diag}[s_1 \cdots s_n]$). The elements of S are known as the *singular values* of H .

The singular value decomposition can also be used to perform a least squares minimization subject to a spherical (ball) constraint on $\hat{\mathbf{x}}$.⁷ Consider the minimization of

$$J = \frac{1}{2}(\tilde{\mathbf{y}} - H\hat{\mathbf{x}})^T(\tilde{\mathbf{y}} - H\hat{\mathbf{x}}) \quad (1.131)$$

subject to the following constraint:

$$\sqrt{\hat{\mathbf{x}}^T \hat{\mathbf{x}}} \leq \gamma \quad (1.132)$$

where γ is some known constant. Equation (1.132) constrains $\hat{\mathbf{x}}$ to lie within or on a sphere. The solution to this problem can be given using a singular value decomposition as follows⁷

$$H = USV^T \quad (1.133a)$$

$$[\mathbf{v}_1, \dots, \mathbf{v}_n] = V \quad (1.133b)$$

$$\mathbf{z} = U^T \tilde{\mathbf{y}} \quad (1.133c)$$

$$r = \text{rank}(H) \quad (1.133d)$$

If the following inequality is true:

$$\sum_{i=1}^r \left(\frac{z_i}{s_i} \right)^2 > \gamma^2 \quad (1.134)$$

then find λ^* such that

$$\sum_{i=1}^r \left(\frac{s_i z_i}{s_i^2 + \lambda^*} \right)^2 = \gamma^2 \quad (1.135)$$

and the optimal estimate is given by

$$\hat{\mathbf{x}} = \sum_{i=1}^r \left(\frac{s_i z_i}{s_i^2 + \lambda^*} \right) \mathbf{v}_i \quad (1.136)$$

If the inequality in Equation (1.134) is not satisfied, then the optimal estimate is given by

$$\hat{\mathbf{x}} = \sum_{i=1}^r \left(\frac{z_i}{s_i} \right) \mathbf{v}_i \quad (1.137)$$

It can be shown that there exists a unique positive solution for λ^* which can be found using Newton's root solving method. A more general case of the quadratic inequality constraint can be found in Golub and Van Loan.⁷

Example 1.9: Consider the following model:

$$y = x_1 + x_2 t + x_3 t^2$$

Given a set of 101 measurements, shown in [Figure 1.11](#), we are asked to determine $\hat{\mathbf{x}}$ such that $\hat{\mathbf{x}}^T \hat{\mathbf{x}} \leq 14$. After forming the H matrix, we determine that the rank of H is $r = 3$, and the singular values are given by

$$S = \text{diag} [456.3604 \ 15.5895 \ 3.1619]$$

The singular values clearly show that this least squares problem is well posed since the condition number is given by $456.36/3.16 = 144.33$. Forming the \mathbf{z} vector, and with $\gamma^2 = 14$, we see that the inequality in Equation (1.134) is satisfied with the given measurements. The optimal value for λ^* in Equation (1.135) was determined using Newton's root solving with a starting value of 0, and converged to a value of $\lambda^* = 0.245$. The optimal estimate in Equation (1.136) is given by

$$\hat{\mathbf{x}} = \begin{bmatrix} 3.0209 \\ 1.9655 \\ 1.0054 \end{bmatrix}$$

The inequality constraint in Equation (1.132) is clearly satisfied since $\hat{\mathbf{x}}^T \hat{\mathbf{x}} = 14$ (in this case the equality condition is actually satisfied). It is interesting to note that the solution using standard least squares in Equation (1.26) is given by

$$\hat{\mathbf{x}}_{ls} = \begin{bmatrix} 3.0686 \\ 1.9445 \\ 1.0067 \end{bmatrix}$$

We can see that the solutions are nearly identical; however, the standard least squares solution violates the inequality constraint since $\hat{\mathbf{x}}_{ls}^T \hat{\mathbf{x}}_{ls} = 14.2109 \geq 14$. Also, since the standard least squares solution gives a condition that violates the constraint, we expect that the optimal solution should give estimates that lie on the surface of the sphere (i.e., on the equality constraint).

This section has introduced some popular matrix decompositions used in linear least squares. Choosing which decomposition to use is primarily dependent upon the particular application, numerical concerns, and desired level of accuracy. For example, the singular value decomposition is one of the most robust algorithms to compute the least squares estimates. However, it is also one of the most computationally expensive algorithms. The decompositions presented in this section do not represent an exhaustive treatise of the subject. For the interested reader, the many references cited throughout this section give more thorough treatments of the subject matter. In particular, both the QR and singular-value decomposition algorithms can be generalized to include the case that H is either row or column rank deficient.¹⁸

1.6.2 Kronecker Factorization and Least Squares

The Singular Value Decomposition (SVD) approach of [§1.6.1](#) can be used to improve the numerical accuracy of the solution over the equivalent standard least

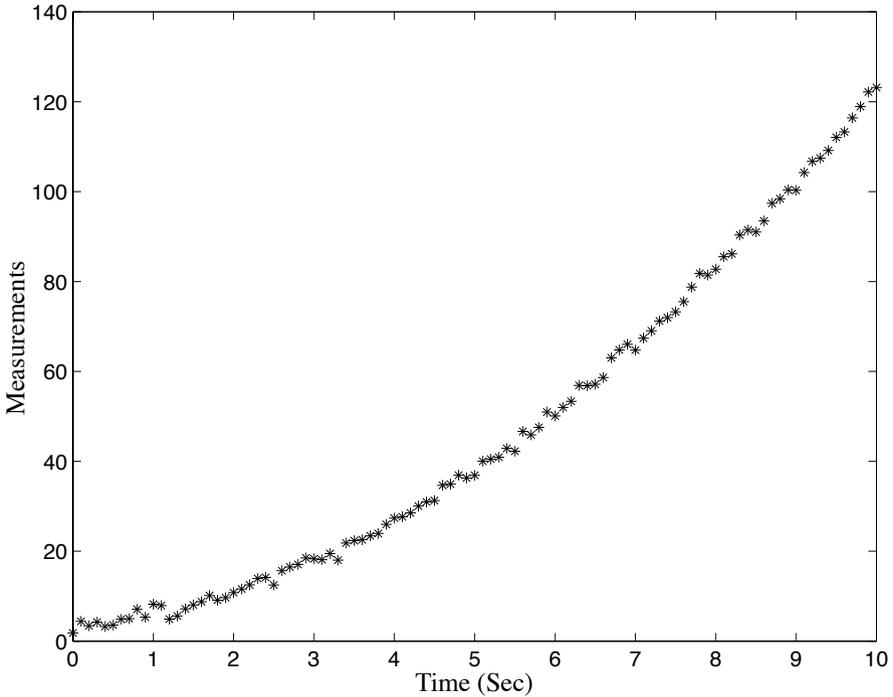


Figure 1.11: Measurements of $y(t)$

squares solution. However, this comes at a significant computational cost. In this section another approach based on the Kronecker factorization¹⁹ is shown that can be used to improve the accuracy and reduce the computational costs for a certain class of problems. The Kronecker product is defined as

$$H = A \otimes B \equiv \begin{bmatrix} a_{11}B & a_{12}B & \cdots & a_{1\beta}B \\ a_{21}B & a_{22}B & \cdots & a_{2\beta}B \\ \vdots & \vdots & \ddots & \vdots \\ a_{\alpha 1}B & a_{\alpha 2}B & \cdots & a_{\alpha\beta}B \end{bmatrix} \quad (1.138)$$

where H is an $M \times N$ dimension matrix, A is an $\alpha \times \beta$ matrix, and B is a $\gamma \times \delta$ matrix. The Kronecker product is only valid when $M = \alpha \gamma$ and $N = \beta \delta$. The key result for least squares problems is that if $H = A \otimes B$, then Equation (1.26) reduces down to

$$\hat{\mathbf{x}} = \{[(A^T A)^{-1} A^T] \otimes [(B^T B)^{-1} B^T]\} \tilde{\mathbf{y}} \quad (1.139)$$

In essence the Kronecker product takes the square root of the matrix dimensions in regard to the computational difficulty.

A key question now arises: “Under what conditions can a matrix be factored as a Kronecker product of smaller matrices?” This is a difficult question to answer, but

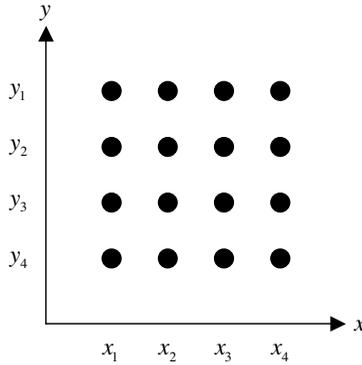


Figure 1.12: Gridded Data

fortunately it is easy to show that some important curve fitting problems lead to a Kronecker factorization, such as the case of gridded data depicted in Figure 1.12. We first consider the case of fitting a two-variable polynomial to data on an x - y grid:

$$z = f(x, y) = \sum_{p=0}^M \sum_{q=0}^N c_{pq} x^p y^q \tag{1.140}$$

where the measurements are now defined by

$$z_{ij} = f(x_i, y_j) + v_{ij} \tag{1.141}$$

for $i = 1, 2, \dots, n_x$ and $j = 1, 2, \dots, n_y$. Now consider the special case of $M = 2$, $N = 1$, $n_x = 4$, and $n_y = 3$. The quantity z in Equation (1.140) is given by

$$z = c_{00} + c_{01}y + c_{10}x + c_{11}xy + c_{20}x^2 + c_{21}x^2y \tag{1.142}$$

The least squares measurement model is now given by

$$\begin{bmatrix} z_{11} \\ z_{12} \\ z_{13} \\ \vdots \\ z_{41} \\ z_{42} \\ z_{43} \end{bmatrix} = \begin{bmatrix} 1 & y_1 & x_1 & x_1 y_1 & x_1^2 & x_1^2 y_1 \\ 1 & y_2 & x_1 & x_1 y_2 & x_1^2 & x_1^2 y_2 \\ 1 & y_3 & x_1 & x_1 y_3 & x_1^2 & x_1^2 y_3 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & y_1 & x_4 & x_4 y_1 & x_4^2 & x_4^2 y_1 \\ 1 & y_2 & x_4 & x_4 y_2 & x_4^2 & x_4^2 y_2 \\ 1 & y_3 & x_4 & x_4 y_3 & x_4^2 & x_4^2 y_3 \end{bmatrix} \begin{bmatrix} c_{00} \\ c_{01} \\ c_{10} \\ c_{11} \\ c_{20} \\ c_{21} \end{bmatrix} + \begin{bmatrix} v_{11} \\ v_{12} \\ v_{13} \\ \vdots \\ v_{41} \\ v_{42} \\ v_{43} \end{bmatrix} \equiv H\mathbf{c} + \mathbf{v} \tag{1.143}$$

where H , \mathbf{c} , and \mathbf{v} have dimensions of 12×6 , 6×1 , and 12×1 , respectively. We can now easily verify that the matrix H has a Kronecker factorization given by

$$H = \begin{bmatrix} 1 & x_1 & x_1^2 \\ 1 & x_2 & x_2^2 \\ 1 & x_3 & x_3^2 \\ 1 & x_4 & x_4^2 \end{bmatrix} \otimes \begin{bmatrix} 1 & y_1 \\ 1 & y_2 \\ 1 & y_3 \end{bmatrix} \equiv H_x \otimes H_y \tag{1.144}$$

Downloaded by [Utah State University] at 23:06 03 December 2013

where H_x and H_y have dimensions of 4×3 and 3×2 , respectively. Thus, perhaps, it is not surprising that the two-variable Vandermonde matrix can be produced by the Kronecker product of the corresponding one-variable Vandermonde matrices. The consequences in the least squares solution are enormous, since the estimate for the coefficient vector, \mathbf{c} , can be computed by

$$\hat{\mathbf{c}} = (H^T H)^{-1} H^T \tilde{\mathbf{z}} = \{[(H_x^T H_x)^{-1} H_x^T] \otimes [(H_y^T H_y)^{-1} H_y^T]\} \tilde{\mathbf{z}} \tag{1.145}$$

Hence, only inverses of 3×3 and 2×2 matrices need to be computed instead of an inverse of a 6×6 matrix. In general, for H of dimension $M \times N$, and H_x and H_y of dimensions about $\sqrt{M}/2$ and $\sqrt{N}/2$, respectively, the least squares computational burden is reduced from an order of n^3 operations to an order of $(\sqrt{n})^3$ operations! Furthermore, as will be shown in example 1.10, the accuracy of the solution is also vastly improved.

The previous Kronecker factorization solution in the least squares problem can be expanded to the n -dimensional case, where data are at the vertices of an n -dimensional grid:

$$z = f(x_1, x_2, \dots, x_n) = \sum_{i_1=1}^{N_1} \sum_{i_2=1}^{N_2} \dots \sum_{i_n=1}^{N_n} c_{i_1 i_2 \dots i_n} \phi_{i_1}(x_1) \phi_{i_2}(x_2) \dots \phi_{i_n}(x_n) \tag{1.146}$$

where $\phi_{i_j}(x_j)$ are basis functions. The measurements now follow

$$\tilde{z}_{j_1 j_2 \dots j_n} \quad \text{at} \quad (x_{1j_1}, x_{2j_2}, \dots, x_{nj_n}) \tag{1.147}$$

for $j_1 = 1, 2, \dots, M_1$ through $j_n = 1, 2, \dots, M_n$. The vectors $\tilde{\mathbf{z}}$ and \mathbf{c} are now denoted by

$$\tilde{\mathbf{z}} = [\tilde{z}_{11\dots 11} \dots \tilde{z}_{11\dots 1M_n} \dots \tilde{z}_{M_1 M_2 \dots M_{n-1} 1} \dots \tilde{z}_{M_1 M_2 \dots M_{n-1} M_n}]^T \tag{1.148a}$$

$$\mathbf{c} = [c_{11\dots 11} \dots c_{11\dots 1N_1} \dots c_{N_1 N_2 \dots N_{n-1} 1} \dots c_{N_1 N_2 \dots N_{n-1} N_n}]^T \tag{1.148b}$$

The matrix H is given by

$$H = H_1 \otimes H_2 \otimes \dots \otimes H_n \tag{1.149}$$

with

$$H_i = \begin{bmatrix} \Phi_1(x_{i_1}) & \Phi_2(x_{i_1}) & \dots & \Phi_{N_i}(x_{i_1}) \\ \vdots & \vdots & \ddots & \vdots \\ \Phi_1(x_{i_{M_i}}) & \Phi_2(x_{i_{M_i}}) & \dots & \Phi_{N_i}(x_{i_{M_i}}) \end{bmatrix}, \quad i = 1, 2, \dots, n \tag{1.150}$$

where the Φ 's are sub-matrices composed of the basis functions $\phi_{i_1}(x_1)$ through $\phi_{i_n}(x_n)$. The estimate for the coefficient vector, \mathbf{c} , can be computed by

$$\hat{\mathbf{c}} = \{[(H_1^T H_1)^{-1} H_1^T] \otimes \dots \otimes [(H_n^T H_n)^{-1} H_n^T]\} \tilde{\mathbf{z}} \tag{1.151}$$

Downloaded by [Utah State University] at 23:06 03 December 2013

Therefore, the least squares solution is given by a Kronecker product of sub-matrices with much smaller dimension than the original problem.

Example 1.10: In this simple example, the power of the Kronecker product in least squares problems is illustrated. We consider a 21×21 grid over the intervals $-2 \leq x \leq 2$ and $-2 \leq y \leq 2$ with functions given by

$$\begin{bmatrix} 1 & x & x^2 & x^3 & x^4 & x^5 \\ 1 & y & y^2 & y^3 & y^4 & y^5 \end{bmatrix}$$

The 21×6 matrices H_x and H_y are given by

$$H_x = \begin{bmatrix} 1 & x_1 & x_1^2 & x_1^3 & x_1^4 & x_1^5 \\ 1 & x_2 & x_2^2 & x_2^3 & x_2^4 & x_2^5 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{21} & x_{21}^2 & x_{21}^3 & x_{21}^4 & x_{21}^5 \end{bmatrix}, \quad H_y = \begin{bmatrix} 1 & y_1 & y_1^2 & y_1^3 & y_1^4 & y_1^5 \\ 1 & y_2 & y_2^2 & y_2^3 & y_2^4 & y_2^5 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & y_{21} & y_{21}^2 & y_{21}^3 & y_{21}^4 & y_{21}^5 \end{bmatrix}$$

The 441×36 matrix H is just the Kronecker product of H_x and H_y , so that $H = H_x \otimes H_y$. The true coefficient vector, \mathbf{c} , has elements simply given by 1 in this formulation. As shown previously, the Kronecker factorization gives a substantial savings in numerical computations. We also wish to investigate the accuracy of this approach. To accomplish this task, no noise is added to form the 441×1 vector of measurements, which is simply given by $\bar{\mathbf{z}} = H \mathbf{c}$.

The numerical accuracy is shown by computing $\epsilon \equiv \|\hat{\mathbf{c}} - \mathbf{c}\|$, which is ideally zero. Using the standard least squares solution of §1.2.1, which takes the inverse of a 36×36 matrix, gives $\epsilon = 7.15 \times 10^{-10}$. Using the SVD solution of §1.6.1 gives $\epsilon = 1.15 \times 10^{-12}$, which provides more accuracy but at a price of a substantial computational cost over the standard least squares solution. Using the Kronecker factorization gives $\epsilon = 1.66 \times 10^{-13}$, which provides even better accuracy than the SVD solution, but is more computationally efficient than the standard least squares solution. An SVD solution for each inverse in the Kronecker factorization can also be used instead of the standard inverse. This approach gives $\epsilon = 1.20 \times 10^{-13}$, which provides the most accurate solution with only a modest increase in computational cost over the standard Kronecker factorization solution. This example clearly shows the power of the Kronecker factorization for curve fitting problems with gridded data.

This section summarized a powerful solution to the curve fitting problem involving gridded data. The Kronecker factorization leads to substantial computational savings, while improving the numerical accuracy of the solution, over the standard least squares solution. This is especially significant for systems involving polynomial models, which have a tendency to be ill conditioned. This approach has substantial advantages for applications in many systems, such as satellite imagery, terrain modeling, and photogrammetry. More details on the usefulness of the Kronecker factorization in least squares applications can be found in Ref. [19].

1.6.3 Levenberg-Marquardt Method

The differential correction algorithm in §1.4 may not be suitable for some nonlinear problems since convergence cannot be guaranteed, unless the *a priori* estimate is close to a minimum in the loss function. This difficulty may be overcome by using the *method of steepest descent* (see Appendix D). This method adjusts the current estimate so that the most favorable direction is given (i.e., the direction of steepest descent), which is along the negative gradient of J . The method of steepest descent often converges rapidly for the first few iterations, but has difficulty converging to a solution because the slope becomes more and more shallow as the number of iterations increases.

The Levenberg-Marquardt algorithm²⁰ overcomes both the difficulties of the standard differential correction approach when an accurate initial estimate is not given, and the slow convergence problems of the method of steepest descent when the solution is close to minimizing the nonlinear least squares loss function (1.89). The paper by Marquardt develops the entire algorithm; however, a significant acknowledgment is given to Levenberg.²¹ Hence, the algorithm is usually referred to by both authors. This algorithm performs an optimum interpolation between the differential correction, which approximates a second-order Taylor series expansion of J , and the method of steepest descent, which uses a first-order approximation of local J behavior.

We first derive an expression for the gradient correction. Consider the loss function given by Equation (1.96):

$$J = \frac{1}{2} \Delta \mathbf{y}^T W \Delta \mathbf{y} \quad (1.152)$$

The gradient of Equation (1.152) is given by

$$\nabla_{\hat{\mathbf{x}}} J = -H^T W \Delta \mathbf{y}_c \quad (1.153)$$

where

$$H \equiv \left. \frac{\partial \mathbf{f}}{\partial \mathbf{x}} \right|_{\hat{\mathbf{x}}} \quad (1.154)$$

The method of gradients seeks corrections down the gradient:

$$\Delta \mathbf{x} = -\frac{1}{\eta} \nabla_{\hat{\mathbf{x}}} J = \frac{1}{\eta} H^T W \Delta \mathbf{y}_c \quad (1.155)$$

where $1/\eta$ is a scalar which controls the step size. The poor terminal convergence of the first-order gradient and the less reliable early convergence of the second-order differential correction algorithm can be compromised, as in the Levenberg-Marquardt algorithm, with the modified normal equations:

$$\Delta \mathbf{x} = (H^T W H + \eta \mathcal{H})^{-1} H^T W \Delta \mathbf{y}_c \quad (1.156)$$

where \mathcal{H} is a diagonal matrix with entries given by the diagonal elements of $H^T W H$ or in some cases simply the identity matrix. By using the algorithm in Equation (1.156) the search direction is an intermediate between the steepest descent and

the differential correction direction. As $\eta \rightarrow 0$, Equation (1.156) is equivalent to the differential correction method; however, as $\eta \rightarrow \infty$, if $\mathcal{H} = I$, Equation (1.156) reduces to a steepest descent search along the negative gradient of J .

Controlling η (and therefore both the magnitude and direction of $\Delta \mathbf{x}$) is a heuristic art form that can be tuned by the user. Generally η is large in early iterations and should definitely be reduced toward zero in the region near the minimum. To capture the spirit of the approach, here is a typical recipe for implementing the Levenberg-Marquardt algorithm:

1. Compute Equation (1.89) using an initial estimate for $\hat{\mathbf{x}}$, denoted by \mathbf{x}_c .
2. Use Equations (1.156) and (1.91) to update the current estimate with a large value for η (usually much larger than the norm of H^TWH , typically 10 to 100 times the norm).
3. Recompute Equation (1.89) with the new estimate. If the new value for Equation (1.89) is \geq the value computed in step 1, then the new estimate is disregarded and η is replaced by $f\eta$, where f is a fixed positive constant, usually between 1 and 10 (we suggest a default of 5). Otherwise, retain the estimate, and replace η with η/f .
4. After each subsequent iteration, compare the new value of Equation (1.89) with its value using the previous estimate and replace η with $f\eta$ or η/f as in step 3. The estimate $\hat{\mathbf{x}}$ is retained if J in Equation (1.89) continues to decrease and discarded if (1.89) increases.

This procedure continues until the difference in Equation (1.89) between two consecutive iterations is small. The Levenberg-Marquardt method is heuristic, seeking to find the middle ground between the method of steepest descent and the Gaussian differential correction, tending toward the Gaussian differential correction in the terminal corrections. However, a little effort in tuning this algorithm often leads to a significantly enhanced domain of convergence.

Example 1.11: In example 1.8, we used nonlinear least squares to determine the parameters of an inertially and aerodynamically symmetric projectile. In this example we begin with the same start values, except that the start value for λ_1 is equal to -0.8500 instead of -0.1500 . For this initial value, the standard least squares solution diverges rapidly with each iteration. Therefore, we must use a different starting set or, in this case, we choose to use the Levenberg-Marquardt algorithm. For this algorithm, we set the initial value for η to 1×10^6 . Results in the convergence history are summarized below.

Parameter	Iteration Number				
	0	10	15	...	20
k_1	0.5000	0.3601	0.0844		0.1999
k_2	0.2500	0.1946	0.2099		0.0997
k_3	0.1250	0.0905	0.0620		0.0500
k_4	0.0000	-0.0062	0.0111		0.0002
k_5	0.0000	-0.0047	-0.0004		0.0001
λ_1	-0.8500	-0.7977	-0.0436		-0.0998
λ_2	-0.0600	-0.0760	-0.1270		-0.0497
λ_3	-0.0300	-0.0418	-0.0436		-0.0250
ω_1	0.2600	0.1094	0.1621		0.2500
ω_2	0.5500	0.5505	0.4950		0.4999
ω_3	0.9500	0.9582	0.9874		0.9998
δ_1	0.0100	0.0060	0.5068		0.0010
δ_2	0.0100	-0.1234	-0.3482		0.0001
δ_3	0.0100	0.1225	0.1918		-0.0001
η	10^6	0.5120	0.0041		10^{-6}

Clearly, the Levenberg-Marquardt algorithm converges to the correct estimates for this case, where the classical Gaussian differential correction fails.

1.6.4 Projections in Least Squares

In this section we give a geometrical interpretation of least squares. The term “normal” in Normal Equations implies that there is a geometrical interpretation to least squares. In fact, we will show that the least squares solution for $\hat{\mathbf{x}}$ provides the *orthogonal projection*, hence normal, of $\tilde{\mathbf{y}}$ onto a *subspace* which is spanned by columns of the matrix H . Let us illustrate this concept using the simple scalar case of least squares. Say we wish to determine \hat{x} which minimizes

$$J = \frac{1}{2}(\tilde{\mathbf{y}} - \hat{x}\mathbf{h})^T(\tilde{\mathbf{y}} - \hat{x}\mathbf{h}) \quad (1.157)$$

where \mathbf{h} is the basis function vector. The necessary conditions yield the following simple solution:

$$\hat{x} = \frac{\mathbf{h}^T \tilde{\mathbf{y}}}{\mathbf{h}^T \mathbf{h}} \quad (1.158)$$

The residual error is given by

$$\mathbf{e} = (\tilde{\mathbf{y}} - \hat{x}\mathbf{h}) \quad (1.159)$$

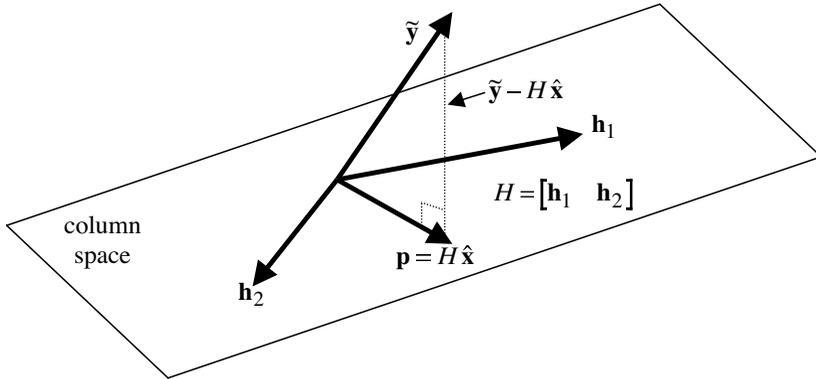


Figure 1.13: Projection onto the Column Space of a 3×2 Matrix

Now, left multiply the residual error by \mathbf{h}^T in Equation (1.159) and substitute Equation (1.158) into Equation (1.159). This yields

$$\begin{aligned}
 \mathbf{h}^T \mathbf{e} &= \mathbf{h}^T (\tilde{\mathbf{y}} - \hat{\mathbf{x}}\mathbf{h}) \\
 &= \mathbf{h}^T \left(\tilde{\mathbf{y}} - \frac{\mathbf{h}^T \tilde{\mathbf{y}}}{\mathbf{h}^T \mathbf{h}} \mathbf{h} \right) \\
 &= \mathbf{h}^T \tilde{\mathbf{y}} - \frac{\mathbf{h}^T \tilde{\mathbf{y}}}{\mathbf{h}^T \mathbf{h}} \mathbf{h}^T \mathbf{h} \\
 &= 0
 \end{aligned}
 \tag{1.160}$$

This shows that the angle between \mathbf{h} and \mathbf{e} is 90 degrees, so that the line connecting $\tilde{\mathbf{y}}$ to $\hat{\mathbf{x}}\mathbf{h}$ must be *perpendicular* to \mathbf{h} .

The aforementioned scalar case is easily expanded to the multi-dimensional case where $\tilde{\mathbf{y}}$ is *projected* onto a subspace rather than just onto a line. In this case, the vector $\mathbf{p} \equiv H\hat{\mathbf{x}}$ must be the projection of $\tilde{\mathbf{y}}$ onto the column space of H , and the residual error \mathbf{e} must be perpendicular to that space.²² This is illustrated for a simple 3×2 case in Figure 1.13. In other words, the residual error must be perpendicular to every column (\mathbf{h}_i) of H , so that

$$\begin{aligned}
 \mathbf{h}_1^T (\tilde{\mathbf{y}} - H\hat{\mathbf{x}}) &= 0 \\
 \mathbf{h}_2^T (\tilde{\mathbf{y}} - H\hat{\mathbf{x}}) &= 0 \\
 &\vdots \\
 \mathbf{h}_n^T (\tilde{\mathbf{y}} - H\hat{\mathbf{x}}) &= 0
 \end{aligned}
 \tag{1.161}$$

or

$$H^T (\tilde{\mathbf{y}} - H\hat{\mathbf{x}}) = 0
 \tag{1.162}$$

which gives the normal equations again. The projection of $\tilde{\mathbf{y}}$ onto the column space is therefore given by

$$\mathbf{p} = H(H^T H)^{-1} H^T \tilde{\mathbf{y}}
 \tag{1.163}$$

Geometrically, this means that the closest point to $\tilde{\mathbf{y}}$ on the column space of H is \mathbf{p} . Equation (1.163) expresses in matrix terms the construction of a perpendicular line from $\tilde{\mathbf{y}}$ to the column space of H .²² The *projection matrix* is given by

$$\mathcal{P} = H(H^T H)^{-1} H^T \quad (1.164)$$

The projection matrix \mathcal{P} can readily be seen to be symmetric. More importantly, the projection matrix has another property, known as *idempotence*, which states

$$\mathcal{P}\tilde{\mathbf{y}} = [\mathcal{P} \mathcal{P} \dots \mathcal{P}]\tilde{\mathbf{y}} \quad (1.165)$$

The idempotence property shows that once a vector has been obtained as the projection onto a subspace using \mathcal{P} , it can never be modified by any further application of \mathcal{P} .³ The corresponding prediction error, \mathbf{e}_{min} , once the solution for $\hat{\mathbf{x}}$ has been found, is given by

$$\mathbf{e}_{min} = (I - \mathcal{P})\tilde{\mathbf{y}} \quad (1.166)$$

where the matrix $(I - \mathcal{P})$ is the *orthogonal complement* of \mathcal{P} . It is easy to show that $(I - \mathcal{P})$ must also be a projection matrix, since it projects $\tilde{\mathbf{y}}$ onto the orthogonal complement.

1.7 Summary

With some reluctance, the curve fitting example of §1.1 was presented prior to discussion of the methods of §1.2 necessary to carry out the calculations. On several subsequent occasions herein, theoretical development of *methods* follows typical *results*, to provide motivation and to allow some *a priori* evaluation by the reader of the role played by the methodology under development.

The results developed in §1.2 are among the most important in estimation theory. Indeed, the bulk of estimation theory could be viewed as extensions, modifications, or generalizations of these basic results that address a wider variety of mathematical models and measurement strategies. We shall see, however, that the results of §1.2 can be placed upon a more rigorous foundation and several important new insights gained through study of the developments of Chapter 2 and Appendices B and C.

The sequential estimation results in §1.3 are the simplest version of a class of procedures known as *Kalman Filter* algorithms. Indeed, with the advancement of computer technology in today's age, sequential algorithms have found their way into mainstream applications in a wide variety of areas. Numerous investigators have extended/applied these algorithms since the most fundamental results were published by Kalman and Bucy.⁶ The constrained least squares solution⁵ in Equation (1.42) is closely related to the sequential estimation solution in Equation (1.78), and can in fact be obtained from it by limiting arguments (allowing the weight of the constraint "observation" equations to approach infinity). A substantial portion of the present text deals with sequential estimation methodology and applications thereof.

The differential correction procedures documented in §1.4 are most fundamental whenever estimation methods must be applied to a nonlinear problem. It is interesting to note that the original estimation problem motivating Gauss (i.e., determination of the planetary orbits from telescope/sextant observations) was nonlinear, and his methods (essentially §1.4) have survived as a standard operating procedure to this day. Other *mathematical programming* methods (Appendix D), such as the gradient method, can also be employed in minimizing the sum square residuals.

A summary of the key formulas presented in this chapter is given below.

- Linear Least Squares

$$\begin{aligned} \tilde{\mathbf{y}} &= H\mathbf{x} + \mathbf{v} \\ \hat{\mathbf{x}} &= (H^T H)^{-1} H^T \tilde{\mathbf{y}} \end{aligned}$$

- Weighted Least Squares

$$\begin{aligned} \tilde{\mathbf{y}} &= H\mathbf{x} + \mathbf{v} \\ \hat{\mathbf{x}} &= (H^T W H)^{-1} H^T W \tilde{\mathbf{y}} \end{aligned}$$

- Constrained Least Squares

$$\begin{aligned} \tilde{\mathbf{y}}_1 &= H_1 \mathbf{x} + \mathbf{v} \\ \tilde{\mathbf{y}}_2 &= H_2 \hat{\mathbf{x}} \\ \hat{\mathbf{x}} &= \bar{\mathbf{x}} + K(\tilde{\mathbf{y}}_2 - H_2 \bar{\mathbf{x}}) \\ K &= (H_1^T W_1 H_1)^{-1} H_1^T [H_2 (H_1^T W_1 H_1)^{-1} H_2^T]^{-1} \\ \bar{\mathbf{x}} &= (H_1^T W_1 H_1)^{-1} H_1^T W_1 \tilde{\mathbf{y}}_1 \end{aligned}$$

- Sequential Least Squares

$$\begin{aligned} \hat{\mathbf{x}}_{k+1} &= \hat{\mathbf{x}}_k + K_{k+1}(\tilde{\mathbf{y}}_{k+1} - H_{k+1} \hat{\mathbf{x}}_k) \\ K_{k+1} &= P_k H_{k+1}^T [H_{k+1} P_k H_{k+1}^T + W_{k+1}^{-1}]^{-1} \\ P_{k+1} &= [I - K_{k+1} H_{k+1}] P_k \end{aligned}$$

- Nonlinear Least Squares (see Figure 1.9)

$$\begin{aligned} \tilde{\mathbf{y}} &= \mathbf{f}(\mathbf{x}) + \mathbf{v} \\ H &\equiv \left. \frac{\partial \mathbf{f}}{\partial \mathbf{x}} \right|_{\mathbf{x}_c} \\ \Delta \mathbf{y} &\equiv \tilde{\mathbf{y}} - \mathbf{f}(\mathbf{x}_c) \\ \Delta \mathbf{x} &= (H^T W H)^{-1} H^T W \Delta \mathbf{y} \\ \hat{\mathbf{x}} &= \mathbf{x}_c + \Delta \mathbf{x} \end{aligned}$$

- QR Decomposition

$$H = QR$$

$$R\hat{\mathbf{x}} = Q^T \tilde{\mathbf{y}}$$

- Singular Value Decomposition

$$H = USV^T$$

$$\hat{\mathbf{x}} = VS^{-1}U^T \tilde{\mathbf{y}}$$

- Kronecker Factorization

$$\hat{\mathbf{c}} = \{[(H_1^T H_1)^{-1} H_1^T] \otimes \dots \otimes [(H_N^T H_N)^{-1} H_N^T]\} \tilde{\mathbf{z}}$$

- The Levenberg-Marquardt Algorithm

$$\Delta \mathbf{x} = (H^T W H + \eta \mathcal{H})^{-1} H^T W \Delta \mathbf{y}_c$$

$$\mathcal{H} = \text{diag}[H^T W H]$$

- Projection Matrix and Idempotence

$$\mathcal{P} = H(H^T H)^{-1} H^T$$

$$\mathcal{P} \tilde{\mathbf{y}} = [\mathcal{P} \ \mathcal{P} \ \dots \ \mathcal{P}] \tilde{\mathbf{y}}$$

Exercises

- 1.1 Prove that $H^T H$ is a symmetric matrix.
- 1.2 Prove that if W is a symmetric positive definite matrix, then $H^T W H$ will always be positive semi-definite (hint: any positive definite matrix W can be factored into $W = R^T R$, where R is an upper triangular matrix, known as the Cholesky Decomposition).
- 1.3 Following the notation of §1.2 consider the m dimensional observation equation

$$\tilde{\mathbf{y}} = Hx + \mathbf{v}$$

$$\tilde{\mathbf{y}} = H\hat{x} + \mathbf{e}$$

with

$$H = [1 \ 1 \ \dots \ 1]^T$$

These observation equations hold for the simplest situation in which an unknown scalar parameter x is *directly* measured m times (assume that the

measurements errors have zero mean and known, equal variances). From the normal equations (1.26), establish the well-known truth that the optimum least squares estimate \hat{x} of x is the sample mean

$$\hat{x} = \frac{1}{m} \sum_{i=1}^m \tilde{y}_i$$

- 1.4 Suppose that \mathbf{v} in exercise 1.3 is a constant vector (i.e., a *bias error*). Evaluate the loss function (1.21) in terms of v_i only and discuss how the value of the loss function changes with a bias error in the measurements instead of a zero mean assumption.
- 1.5 Show that the mean of the linear least squares residuals, given by Equation (1.1), vanishes identically if *one* of the linearly independent basis functions is a constant.
- 1.6 In this problem we will consider a simple linear *regression* model. The vertical deviation of a point (z_j, y_j) from the line $y = a + bz$ is $e_j = y_j - (a + bz_j)$. Determine closed-form least squares estimates of a and b given measurement sets for z_j and y_j .
- 1.7 Using the simple model

$$y = x_1 + x_2 \sin 10t + x_3 e^{-2t^2}$$

with $x_1 = x_2 = x_3 = 1.0$, generate four sets of “synthetic data” at the instants $t = 0, 0.1, 0.2, 0.3, \dots, 1.0$ by truncating each y value after 6, 4, 2, and 1 significant figures, respectively, to simulate (crudely) measurement errors. Use the normal equations (1.26) to process the measurements and derive \hat{x}_i estimates for each of the four cases. Compare the estimates with the true values (1, 1, 1) in each case.

- 1.8 Use the sequential estimation algorithm (1.78) to (1.80) to process the first three measurements of exercise 1.7 as a single measurement subset and then consider the remaining measurements to become available one at a time, for each of the four synthetic data sets of exercise 1.7.
- 1.9 Consider the following partitioned matrix (assume that $|A_{11}| \neq 0$ and $|A_{22}| \neq 0$):

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$$

Prove that the following matrices are all valid inverses:

$$A^{-1} = \begin{bmatrix} A_{11}^{-1} + A_{11}^{-1}A_{12}B_{22}^{-1}A_{21}A_{11}^{-1} & -A_{11}^{-1}A_{12}B_{22}^{-1} \\ -B_{22}^{-1}A_{21}A_{11}^{-1} & B_{22}^{-1} \end{bmatrix}$$

$$A^{-1} = \begin{bmatrix} B_{11}^{-1} & -B_{11}^{-1}A_{12}A_{22}^{-1} \\ -A_{22}^{-1}A_{21}B_{11}^{-1} & A_{22}^{-1} + A_{22}^{-1}A_{21}B_{11}^{-1}A_{12}A_{22}^{-1} \end{bmatrix}$$

Downloaded by [Utah State University] at 23:06 03 December 2013

$$A^{-1} = \begin{bmatrix} B_{11}^{-1} & -A_{11}^{-1}A_{12}B_{22}^{-1} \\ -A_{22}^{-1}A_{21}B_{11}^{-1} & B_{22}^{-1} \end{bmatrix}$$

where B_{ii} is the Schur complement of A_{ii} , given by

$$B_{11} = A_{11} - A_{12}A_{22}^{-1}A_{21}$$

$$B_{22} = A_{22} - A_{21}A_{11}^{-1}A_{12}$$

Also, prove the matrix inversion lemma from these matrix inverses.

- 1.10** Create 101 synthetic measurements \tilde{y} at 0.1 second intervals of the following:

$$\tilde{y}_j = a \sin t_j - b \cos t_j + v_j$$

where $a = b = 1$, and v is a zero-mean Gaussian noise process with standard deviation given by 0.01. Determine the unweighted least squares estimates for a and b . Using the same measurements, find a value of \tilde{y} that is near zero (near time $\pi/4$), and set that “measurement” value to 1. Compute the unweighted least squares solution, and compare it to the original solution. Then, use weighted least squares to “deweight” the measurement.

- 1.11** In the derivation of the weighted least squares estimator of §1.2.2, the weight matrix W is assumed to be symmetric. How does the solution change if W is no longer symmetric (but still positive definite)?

- 1.12** Using the method of Lagrange multipliers, find all solutions \mathbf{x} of the first necessary conditions for extremals of the function

$$J(\mathbf{x}) = (\mathbf{x} - \mathbf{a})^T W (\mathbf{x} - \mathbf{a})$$

subject to $\mathbf{b}^T \mathbf{x} = c$

where \mathbf{a} and \mathbf{b} are constant vectors, c is a scalar, and W is a symmetric, positive definite matrix.

- 1.13** Consider the following dynamic model:

$$y_k = \sum_{i=1}^n \phi_i y_{k-i} + \sum_{i=1}^p \gamma_i u_{k-i}$$

where u_i is a known input. This ARX (AutoRegressive model with eXogenous input) model extends the simple scalar model given in example 1.2. Given measurements of y_i and the known inputs u_i recast the above model into least squares form and determine estimates for ϕ_i and γ_i .

- 1.14** Program a sequential estimation algorithm to determine in real time the parameters of the ARX model shown in exercise 1.13. Develop some synthetic data with various system models, and verify your algorithm.

- 1.15** One of the most important mathematical equations in history is given by Kepler's equation, which provides powerful geometrical insights into orbiting bodies. This equation is given by

$$M = E - e \sin E$$

where M and E are known as the mean anomaly and eccentric anomaly, respectively, both given in radians, and e is the eccentricity of the orbit. For elliptical orbits $0 < e < 1$. To date, no one has found a closed-form solution for E in terms of M and e . Pick various values for M and e and use nonlinear least squares, which reduces to Newton's method for this equation, to determine E .

1.16 Consider the following dynamic model:

$$\begin{bmatrix} z_1 \\ z_2 \end{bmatrix}_{k+1} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}_k$$

and measurement model

$$\tilde{y}_k = [\sin(\omega_0 \Delta t k) \quad \cos(\omega_0 \Delta t k)] \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}_k + v_k$$

where ω_0 is the harmonic frequency, and Δt is the sampling interval. Create synthetic measurements of the above process with $\omega_0 = 0.4\pi$ rad/sec and $\Delta t = 0.1$ seconds. Also, create different synthetic measurement sets using various values for the standard deviation of v in the measurement errors. Use nonlinear least squares to find an estimate for ω_0 for each synthetic measurement set.

1.17 A measurement process used in three-axis magnetometers for low-Earth attitude determination involves the following measurement model:

$$\mathbf{b}_j = A_j \mathbf{r}_j + \mathbf{c} + \epsilon_j$$

where \mathbf{b}_j is the measurement of the magnetic field (more exactly, magnetic induction) by the magnetometer at time t_j , \mathbf{r}_j is the corresponding value of the geomagnetic field with respect to some reference coordinate system, A_j is the orthogonal attitude matrix (see §A.7.1), \mathbf{c} is the magnetometer bias, and ϵ_j is the measurement error. We can eliminate the dependence on the attitude by transposing terms and computing the square, and can define an effective measurement by

$$\tilde{y}_j = \mathbf{b}_j^T \mathbf{b}_j - \mathbf{r}_j^T \mathbf{r}_j$$

which can be rewritten to form the following measurement model:

$$\tilde{y}_j = 2\mathbf{b}_j^T \mathbf{c} - \mathbf{c}^T \mathbf{c} + v_j$$

where v_j is the effective measurement error, whose closed-form expression is not required for this problem. For this exercise assume that

$$\mathbf{A} \mathbf{r} = \begin{bmatrix} 10 \sin(0.001t) \\ 5 \sin(0.002t) \\ 10 \cos(0.001t) \end{bmatrix}, \quad \mathbf{c} = \begin{bmatrix} 0.5 \\ 0.3 \\ 0.6 \end{bmatrix}$$

Also, assume that ϵ is given by a zero-mean Gaussian noise process with standard deviation given by 0.05 in each component. Using the above values

create 1001 synthetic measurements of \mathbf{b} and \bar{y} at 5-second intervals. The estimated output is computed from

$$\hat{y}_j = 2\mathbf{b}_j^T \hat{\mathbf{c}} - \hat{\mathbf{c}}^T \hat{\mathbf{c}}$$

where $\hat{\mathbf{c}}$ is the estimated solution from the nonlinear least square iterations. Use nonlinear least squares to determine $\hat{\mathbf{c}}$ for a starting value of $\mathbf{x}_c = [0 \ 0 \ 0]^T$. Also, try various starting values to check convergence. Note: $\mathbf{r}^T \mathbf{r} = \mathbf{r}^T \mathbf{A}^T \mathbf{A} \mathbf{r}$, since $\mathbf{A}^T \mathbf{A} = \mathbf{I}$.

- 1.18** An approximate linear solution to exercise 1.17 is possible. The original loss function is quartic in $\hat{\mathbf{c}}$. But this can be approximated by a quadratic loss function using a process known as *centering*.²³ The linearized solution proceeds as follows. First, compute the following averaged values:

$$\bar{y} = \frac{1}{m} \sum_{j=1}^m \tilde{y}_j$$

$$\bar{\mathbf{b}} = \frac{1}{m} \sum_{j=1}^m \mathbf{b}_j$$

where m is the total number of measurements, which is equal to 1001 from exercise 1.17. Next, define the following variables:

$$\check{y}_j = \tilde{y}_j - \bar{y}$$

$$\check{\mathbf{b}}_j = \mathbf{b}_j - \bar{\mathbf{b}}$$

The centered estimate now minimizes the following loss function:

$$\bar{J}(\hat{\mathbf{c}}) = \frac{1}{2} \sum_{j=1}^m (\check{y}_j - 2\check{\mathbf{b}}_j^T \hat{\mathbf{c}})^2$$

Minimizing this function yields

$$\hat{\mathbf{c}} = P \sum_{j=1}^m 2\check{y}_j \check{\mathbf{b}}_j$$

where

$$P \equiv \left[\sum_{j=1}^m 4\check{\mathbf{b}}_j \check{\mathbf{b}}_j^T \right]^{-1}$$

Using the parameters described in exercise 1.17, compare the linear solution described here to the solution obtained by nonlinear least squares. Furthermore, find solutions for $\hat{\mathbf{c}}$ using both approaches with the following trajectory for $\mathbf{A}\mathbf{r}$:

$$\mathbf{A}\mathbf{r} = \begin{bmatrix} 10 \sin(0.001t) \\ 5 \\ 10 \cos(0.001t) \end{bmatrix}$$

Discuss the performance of the linear solution using this assumed trajectory for $\mathbf{A}\mathbf{r}$.

1.19 ♣ Convert the linear batch solution shown in exercise 1.18 to a sequential form (hint: use the matrix inversion lemma in Equation (1.69) to find a sequential form for P). Perform a simulation using the parameters in exercise 1.17 to test your algorithm.

1.20 Consider the following measurement model:

$$\tilde{y}_j = B \exp \left[-\frac{(1-at)^2}{2\sigma^2} \right] + v_j$$

with $a = 1$, $B = 2$, $\sigma = 3$, and let v be represented by a zero-mean Gaussian noise process with standard deviation given by 0.001. Create 101 synthetic measurements at 0.1-second intervals. Use the change of variables in Table 1.1 to determine *linear* least squares estimates for a , B , and σ .

1.21 Analytically expand $y = |\sin t|$ in a Fourier series. Compute the Fourier coefficients using least squares with the basis functions in Equation (1.104) for $n = 10$ and compare the numerical solutions to the analytically derived solutions.

1.22 Consider the following matrix commonly used to describe attitude motion:

$$A = \begin{bmatrix} \cos \theta & \sin \theta & 0 \\ -\sin \theta & \cos \theta & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Prove that the columns of the A matrix are orthonormal.

1.23 Show that the vector $(\mathbf{x} - \mathbf{y})$ is orthogonal to the vector $(\mathbf{x} + \mathbf{y})$ if and only if $\|\mathbf{x}\| = \|\mathbf{y}\|$.

1.24 Prove that the Kronecker product in Equation (1.144) is indeed equivalent to the matrix H given in Equation (1.143).

1.25 Reproduce the results of example 1.10. Try some higher-order polynomials to further show the importance of the solution using the Kronecker factorization.

1.26 Find starting values in exercise 1.17 that cause the standard nonlinear least squares problem to diverge using the following trajectory for \mathbf{Ar} :

$$\mathbf{Ar} = \begin{bmatrix} 10 \sin(0.001t) \\ 5 \\ 10 \cos(0.001t) \end{bmatrix}$$

For example, try starting values of $\mathbf{x}_c = [10 \ 10 \ 10]^T$. Program the Levenberg-Marquardt method, and check convergence for this starting condition as well as various other starting conditions. Also, check the performance of the Levenberg-Marquardt method for various values of η and f (start with $\eta = 10\|H^T H\|$ and $f = 5$).

- 1.27** Consider the projection onto the θ -direction in the $x - y$ plane. Find the projection matrix for the line through $\mathbf{h} = [\cos \theta \ \sin \theta]^T$. Is this matrix invertible? Explain.
- 1.28** Prove that $(I - \mathcal{P})$, with \mathcal{P} given by Equation (1.164), has the idempotence property.

References

- [1] Devore, J.L., *Probability and Statistics for Engineering and Sciences*, Duxbury Press, Pacific Grove, CA, 1995.
- [2] Gauss, K.F., *Theory of the Motion of the Heavenly Bodies Moving about the Sun in Conic Sections, A Translation of Theoria Motus*, Dover Publications, New York, NY, 1963.
- [3] Strobach, P., *Linear Prediction Theory*, Springer-Verlag, Berlin, 1990.
- [4] Juang, J.N. and Pappa, R.S., "An Eigensystem Realization Algorithm for Modal Parameter Identification and Model Reduction," *Journal of Guidance, Control, and Dynamics*, Vol. 8, No. 5, Sept.-Oct. 1985, pp. 620–627.
- [5] Junkins, J.L., "On the Optimization and Estimation of Powered Rocket Trajectories Using Parametric Differential Correction Processes," Tech. Rep. SM G1793, McDonnell Douglas Astronautics Co., 1969.
- [6] Kalman, R.E. and Bucy, R.S., "New Results in Linear Filtering and Prediction Theory," *Journal of Basic Engineering*, March 1961, pp. 95–108.
- [7] Golub, G.H. and Van Loan, C.F., *Matrix Computations*, The Johns Hopkins University Press, Baltimore, MD, 3rd ed., 1996.
- [8] Saaty, T.L., *Modern Nonlinear Equations*, Dover Publications, New York, NY, 1981.
- [9] Mirsky, L., *An Introduction to Linear Algebra*, Dover Publications, New York, NY, 1990.
- [10] Sveshnikov, A.A., *Problems in Probability Theory, Mathematical Statistics and Theory of Random Functions*, Dover Publications, New York, NY, 1978.
- [11] Chihara, T.S., *An Introduction to Orthogonal Polynomials*, Gordan and Breach Science Publishers, New York, NY, 1978.
- [12] Datta, K.B. and Mohan, B.M., *Orthogonal Functions in Systems and Control*, World Scientific, Singapore, 1995.
- [13] Tolstov, G.P., *Fourier Series*, Dover Publications, New York, NY, 1972.

- [14] Gasquet, C. and Witomski, P., *Fourier Analysis and Applications: Filtering, Numerical Computations, Wavelets*, Springer-Verlag, New York, NY, 1978.
- [15] Abramowitz, M. and Stegun, I.A., *Handbook of Mathematical Functions with Formulas, Graphs and Mathematical Tables*, Applied Mathematics Series - 55, National Bureau of Standards, Washington, DC, 1964.
- [16] Ledermann, W., *Handbook of Applicable Mathematics: Analysis*, Vol. 4, John Wiley & Sons, New York, NY, 1982.
- [17] Horn, R.A. and Johnson, C.R., *Matrix Analysis*, Cambridge University Press, Cambridge, MA, 1985.
- [18] Stewart, G.W., *Introduction to Matrix Computations*, Academic Press, New York, NY, 1973.
- [19] Snay, R.A., "Applicability of Array Algebra," *Reviews of Geophysics and Space Physics*, Vol. 16, No. 3, Aug. 1978, pp. 459–464.
- [20] Marquardt, D.W., "An Algorithm for Least-Squares Estimation of Nonlinear Parameters," *Journal of the Society for Industrial and Applied Mathematics*, Vol. 11, No. 2, June 1963, pp. 431–441.
- [21] Levenberg, K., "A Method for the Solution of Certain Nonlinear Problems in Least Squares," *Quarterly of Applied Mathematics*, Vol. 2, 1944, pp. 164–168.
- [22] Strang, G., *Linear Algebra and its Applications*, Saunders College Publishing, Fort Worth, TX, 1988.
- [23] Alonso, R. and Shuster, M.D., "A New Algorithm for Attitude-Independent Magnetometer Calibration," *Proceedings of the Flight Mechanics/Estimation Theory Symposium*, NASA-Goddard Space Flight Center, Greenbelt, MD, May 1994, pp. 513–527.