

# DATA FUSION III: Estimation Theory

**Date:** February 23, 2006      **Time:** 5:00 – 7:30 PM

**Instructor:** Dr. James K Beard

**Credits:** 1      **Course Code:** 0901-501-05      **Registration Number:** 13460

## Today's topics:

1	Definition of Estimation theory .....	1
2	Course Overview .....	2
3	Accuracy and Confidence Limits.....	2
4	Properties of Estimators.....	2
4.1	The Unbiased and Consistency Properties.....	2
4.2	Sufficiency .....	2
4.3	Invariance.....	3
4.4	Efficiency .....	3
5	Types of Estimators .....	3
5.1	Ad Hoc .....	3
5.2	Averages .....	3
5.3	Least Squares .....	3
5.4	Maximum Likelihood .....	6
5.5	Maximum a Posteriori.....	6
5.6	Bayesian Mean.....	7
6	The Method of Maximum Likelihood .....	7
6.1	Definition and properties .....	7
6.2	The Likelihood Equations and the MLE.....	8
6.3	The MLE for Mean and Covariance of Correlated Gaussian Variables.....	8
6.4	Covariance of Estimates .....	10
7	The Kalman Update as an MLE.....	10
8	Estimation Theory Library.....	11
8.1	Foundations of Modern Estimation Theory .....	11
8.2	Major References .....	12

## 1 Definition of Estimation theory

Estimation theory is the science of solutions to the problem of estimating the value of unknown parameters from noisy measurements of related quantities. This includes the definition and analysis of estimation methods, their accuracy, and their use in practical problems.

A function of data is called a *statistic*. A function of the data which approximates an unknown parameter is an *estimator*. Thus, estimation theory is the design, analysis, evaluation, and use of estimators, which are statistics computed from available data.

## 2 Course Overview

We have only five sessions, and we must define the scope of our material carefully to provide a clear treatment of selected topics in that time. We will

- Define estimation theory and its basic types.
- Define and explain the differences in quality and accuracy of the most common estimators.
- Show examples of estimators that you can use in your work.
- Show how to devise and evaluate estimators for new problems
- Show solutions to problems specific to data fusion issues

## 3 Accuracy and Confidence Limits

A *confidence limit* is a probability that is set as part of an experiment design as an *a priori* threshold on the significance of results of the experiment. For example, if we want to be 95% sure that our data is biased before we declare that our experiment has shown that it is biased, then we will devise an estimator for the bias, find the probability distribution function for the estimator under the condition that a bias is not present, and set thresholds at the 5% points in this distribution. Then, the results of an ensemble of experiments will declare the data unbiased 95% of the time when it is, indeed, unbiased.

Setting confidence limits requires that the requirements of the experiment be understood. The impact of incorrect interpretation of the experiment are part of the process of setting confidence limits.

## 4 Properties of Estimators

### 4.1 The Unbiased and Consistency Properties

An *unbiased estimator* is one whose ensemble mean is equal to the value of the unknown parameter.

A *consistent estimator* is one which approaches the value of the unknown parameter as the number of measurements increases without bound.

### 4.2 Sufficiency

A *sufficient estimator* is one that uses all the information in the available data that applies to estimating the value of the unknown parameters.

The *Neyman-Fisher Factorization Theorem* says that an estimator  $\hat{\phi}(\underline{x})$  is sufficient if, and only if, the conditional probability density of the data given the states can be factored in this way:

$$p(\underline{y}|\underline{x}) = g(\underline{x}, \hat{\phi}) \cdot h(\underline{y}) \quad (4.1)$$

where  $h(\underline{y})$  is independent of  $\underline{x}$ . This is important in mathematical proofs of sufficiency. We will not have the time in this course to delve into the mathematics of estimation

theory to the extent that we will use this but you should remember it as background, and try to understand what it means in terms of the property of sufficiency.

### 4.3 Invariance

An estimator  $\underline{\phi}(\underline{x})$  is *invariant* if, when applied to a nonlinear but single-valued function of the parameters  $\underline{f}(\underline{x})$ , the result is  $\underline{f}(\underline{\phi}(\underline{x}))$ . A sufficient estimator will always be invariant.

### 4.4 Efficiency

An estimator is *efficient* (or, *statistically efficient*) if the variance of the estimator is the smallest that is obtainable by any estimator. If an estimator is sufficient and unbiased, it will be efficient.

A term that I sometimes use that is not found in the literature as a standard term is *statistical efficiency*. Occasionally an *ad hoc* estimator that is commonly used has a variance that is demonstrably less than the Cramer-Rao bound, or that of an efficient estimator that is available. An example is averaging the absolute values of the errors to estimate standard deviation instead of averaging the squared errors to find the sample variance. I call the ratio of the variances the statistical efficiency of the estimator. When the statistical efficiency is 100%, the estimator is efficient in the classical usage of the term “efficient estimator.”

## 5 Types of Estimators

### 5.1 Ad Hoc

Estimators that seem “correct” but that have no provenance in estimation theory are common. An example is estimating the number of fish of various types that exist in a lake. The technique that is normally used is to band or mark a small number and return them, then look at the number of fish that are caught of each type. For each type, the ratio of fish that are caught that are banded to those that are not banded is considered to be the best estimate of the number that were banded and thrown back and the rest in the lake.

### 5.2 Averages

Simply averaging noisy quantities to obtain an estimate of their ensemble average is a common example. This method is easily analyzed in terms of the accuracy of the estimate.

### 5.3 Least Squares

Least squares is a simple technique that was first used by Gauss to determine orbits of celestial bodies and the Earth from astronomical observations. An example we commonly use is a way to fit a polynomial to a set of data points. If we denote the set of  $M$  data points as

$$(x_i, y_i), \quad i = 1 \dots M \quad (5.1)$$

and wish to fit a polynomial  $p_N(x)$  of degree  $N$  to those points, we have

$$\begin{aligned} \tilde{y}_i &= p_N(x_i) \\ &= \sum_{j=0}^{N-1} a_j \cdot x_i^j \end{aligned} \quad (5.2)$$

We define a cost function as the total squared error:

$$J = \sum_{i=1}^M (\tilde{y}_i - y_i)^2 \quad (5.3)$$

and find the polynomial coefficients that minimize the total squared error  $J$ . We can pose the problem in vectors and matrices and write the solution immediately. We collect the polynomial coefficients  $a_j$ , the powers of the independent variables  $x_i$ , and the dependent variables  $y_i$  as vectors:

$$\underline{a} = \begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_{N-1} \end{bmatrix} \quad (5.4)$$

$$\underline{x}_i = \begin{bmatrix} 1 \\ x_i \\ \vdots \\ x_i^{N-1} \end{bmatrix} \quad (5.5)$$

$$\underline{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_M \end{bmatrix} \quad (5.6)$$

The data equations become

$$\underline{\tilde{y}} = X \cdot \underline{a} \quad (5.7)$$

where the rows of the matrix  $X$  are the vectors  $\underline{x}_i$ :

$$X = \begin{bmatrix} \underline{x}_1^T \\ \underline{x}_2^T \\ \vdots \\ \underline{x}_M^T \end{bmatrix} \quad (5.8)$$

The error vector  $\underline{y}_e$  is

$$\underline{y}_e = \underline{\tilde{y}} - \underline{y} \quad (5.9)$$

and the total squared error  $J$  is

$$\begin{aligned}
 J &= \underline{y}_e^T \cdot \underline{y}_e = (\tilde{\underline{y}} - \underline{y})^T \cdot (\tilde{\underline{y}} - \underline{y}) \\
 &= (X \cdot \underline{a} - \underline{y})^T \cdot (X \cdot \underline{a} - \underline{y})
 \end{aligned} \tag{5.10}$$

We find the value of  $\underline{a}$  that minimizes the total squared error  $J$  with simple calculus. We take the gradient of  $J$  with respect to  $\underline{a}$  to get  $N$  equations in  $N$  unknowns and solve:

$$\frac{\partial J}{\partial \underline{a}} = 2 \cdot X^T \cdot (X \cdot \underline{a} - \underline{y}) = \underline{0} \tag{5.11}$$

The solution is

$$\underline{a} = (X^T \cdot X)^{-1} \cdot X^T \cdot \underline{y} \tag{5.12}$$

The matrix  $X^T X$  is very interesting and, because it occurs in this very common problem, it has been studied quite a bit. It is of the form

$$X^T \cdot X = \begin{bmatrix} M & \sum_{i=1}^M x_i & \sum_{i=1}^M x_i^2 & \cdots & \sum_{i=1}^M x_i^{N-1} \\ \sum_{i=1}^M x_i & \sum_{i=1}^M x_i^2 & \sum_{i=1}^M x_i^3 & \cdots & \sum_{i=1}^M x_i^N \\ \sum_{i=1}^M x_i^2 & \sum_{i=1}^M x_i^3 & \sum_{i=1}^M x_i^4 & \cdots & \vdots \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \sum_{i=1}^M x_i^{N-1} & \sum_{i=1}^M x_i^N & \cdots & \sum_{i=1}^M x_i^{2N-2} \end{bmatrix} \tag{5.13}$$

When  $M=N$  then  $X$  is square and we can simply invert  $X$  to find a polynomial that passes through all the points; in this case,  $X$  is a Vandermonde matrix, which occurs also in other problems. The Vandermonde matrix is known to be ill-conditioned – that is, practical problems in accuracy are often encountered in inverting it. When  $M>N$ , we are doing a least-squares fit of data to a polynomial, which is a smoothing operation, not a polynomial fit. The matrix for this operation as given by (5.13) is less troublesome than the Vandermonde matrix, and problems are rarely encountered when  $M \gg N$ .

Also instructive is the product  $X^T \underline{y}$ : that we use when  $M>N$ :

$$X^T \cdot \underline{y} = \begin{bmatrix} \sum_{i=1}^M y_i \\ \sum_{i=1}^M x_i \cdot y_i \\ \sum_{i=1}^M x_i^2 \cdot y_i \\ \vdots \\ \sum_{i=1}^M x_i^{N-1} \cdot y_i \end{bmatrix} \tag{5.14}$$

Note that these sums can be accumulated as data is accrued, so that only the partial sums need be kept in memory as the data is collected. Since the element at row  $p$  and column  $q$  of  $X^T X$  is given by

$$\left( X^T \cdot X \right)_{p,q} = \sum_{i=1}^M x_i^{p+q-2} \quad (5.15)$$

This allows these terms to be accumulated as data is accrued also.

## 5.4 Maximum Likelihood

The method of maximum likelihood is based on a simple principle. In general, we can usually write the conditional probability density function of the data, given the states:

$$\text{Likelihood}(\underline{x}) = p(\underline{y}|\underline{x}) \quad (5.16)$$

This is called the *likelihood equation* because at the time that R. A. Fisher devised it, the common term for the probability density function was the likelihood function. This method is simple to use because the likelihood function can be written easily whenever the probability density function of the errors in the measurements can be characterized, which is nearly all the time. Estimators found in this way are called *maximum likelihood estimators* (MLEs). We will present the method of maximum likelihood in greater detail later today.

## 5.5 Maximum a Posteriori

We can write a conditional probability distribution of the states  $\underline{x}$ , given the measurements  $\underline{y}$ , as a *posterior probability* of  $\underline{x}$  given  $\underline{y}$ :

$$PP(\underline{x}, \underline{y}) = p(\underline{x}|\underline{y}) \quad (5.17)$$

If we find the states  $\underline{x}$  that maximize this as a function of the observed data  $\underline{y}$ , this is the *maximum a posteriori* (MAP) estimator.

From Bayes' rule, we know that conditional probabilities are given in terms of joint probabilities according to the rule

$$p(\underline{x}, \underline{y}) = p(\underline{x}|\underline{y}) \cdot p(\underline{y}) = p(\underline{y}|\underline{x}) \cdot p(\underline{x}) \quad (5.18)$$

where the probability densities of separate sets of variables are given by

$$\begin{aligned} p(\underline{x}) &= \int p(\underline{x}, \underline{y}) \cdot d\underline{y} \\ p(\underline{y}) &= \int p(\underline{x}, \underline{y}) \cdot d\underline{x} \end{aligned} \quad (5.19)$$

Thus the posterior probability function can be written in three forms:

$$PP(\underline{x}, \underline{y}) = p(\underline{x}|\underline{y}) = \frac{p(\underline{y}|\underline{x}) \cdot p(\underline{x})}{p(\underline{y})} = \frac{p(\underline{x}, \underline{y})}{p(\underline{y})} \quad (5.20)$$

Thus the MAP estimator can be applied whenever we can write any one of these.

When the *a priori* probability density function of  $\underline{x}$  is unknown without other information, as is often the case, and  $p(\underline{x})$  is taken as uniform, the MAP estimator reduces to the MLE.

## 5.6 Bayesian Mean

If we look at the probability distribution of the states  $\underline{x}$  given the measurements  $\underline{y}$ , and we use it to compute the mean of the states given the measurements  $\underline{y}$ ,

$$\hat{\underline{x}} = \int \underline{x} \cdot p(\underline{x}|\underline{y}) \cdot d\underline{x} \quad (5.21)$$

is called the *Bayesian mean* estimator. This estimator can be shown to be minimum variance, and thus is sufficient, and often it produces an efficient estimator. In addition, it is unbiased. These properties make it a method of choice in mathematical developments but in real-world cases the difficulty in writing  $p(\underline{x}|\underline{y})$  and then performing the integral make its use less practical than the MLE.

The MLE produces estimators that are the same as those produced as the Bayesian mean except in cases where nonlinearities exist between the measurements and the states. When this occurs, the MLE is biased while the Bayesian mean is not, while both are minimum variance estimators. The difficulty in removing the mean from the MLE is typically less than that in implementing a Bayesian mean estimator and produces an equivalent result.

## 6 The Method of Maximum Likelihood

### 6.1 Definition and properties

The method of maximum likelihood is based on a simple principle. In general, we can usually write the conditional probability density function of the data, given the states:

$$\text{Likelihood}(\underline{x}) = p(\underline{y}|\underline{x}) \quad (6.1)$$

This is called the *likelihood equation* because at the time that R. A. Fisher devised it, the common term for the probability density function was the likelihood function.

The method of maximum likelihood is based on the idea that if you find the states  $\underline{x}$  that maximize the value of the likelihood function for the observed values of the measurements  $\underline{y}$ , then you would have a “good” estimator. The importance of maximum likelihood estimators (MLEs) lies in the fact that these properties have been proven for it:

- MLEs are sufficient and consistent.
- MLEs are asymptotically unbiased.
- The errors in MLEs are asymptotically Gaussian.
- MLEs are asymptotically efficient.
- If an unbiased, efficient estimator exists, the MLE will be that estimator.

The MLE is useful in our work is because

- Nearly all common estimators can be shown to be found by the method of maximum likelihood, including the examples given
- The method of maximum likelihood is easily stated and applied for nearly all practical problems.

- When the measurement errors are distributed by a rule other than Gaussian random variables, the method of maximum likelihood can be easily applied without approximation.
- MLEs are usually simple to write, implement, and analyze.

## 6.2 The Likelihood Equations and the MLE

Most probability density functions are given primarily in terms of an exponential. Examples are Gaussian, Rayleigh, and Poisson. An exception is the Cauchy distribution. As such, we usually proceed with the log likelihood function,

$$L(\underline{x}) = \ln\left(p\left(\underline{y}|\underline{x}\right)\right) \quad (6.2)$$

We then find the maximum of this function as the states  $\underline{x}$  are varied through ordinary calculus. We take the gradient of the log likelihood function  $L(\underline{x})$  and set it equal to zero to form the *likelihood equation*:

$$l(\hat{\underline{x}}) = \frac{\partial L(\underline{x})}{\partial \underline{x}} = \underline{0} \quad (6.3)$$

This is  $N$  equations in  $N$  unknowns. Much is said by some that the generality of this equation makes it suspect, particularly in the fact that such an equation can have many roots. In practical cases this happens only in cases where the result is predictable and easily handled, such as in phase ambiguities of  $2\pi$  radians. When this equation does have multiple solutions, then the solution or solutions with the largest value of the log likelihood function are the solution. If there are multiple such solutions, the ambiguous results are part of the solution to the estimation problem.

The solution to the likelihood equation (6.3) is the MLE. Now, we need to know the covariance of the estimate. The Fisher Information Matrix is the second gradient of the log likelihood function:

$$F = - \left. \frac{\partial^2 L(\underline{x})}{\partial \underline{x}^2} \right|_{\underline{x}=\hat{\underline{x}}} \quad (6.4)$$

Fisher showed in his landmark papers (see the Library) that the asymptotic value of the covariance of the estimate is

$$P = F^{-1} \quad (6.5)$$

This value of the covariance matrix is called the *Cramer-Rao Bound*.

If the likelihood equation (6.3) is linear in the states, the MLE is unbiased and statistically efficient, and its covariance is given by (6.5). If the likelihood equation is nonlinear, the MLE is sufficient, asymptotically unbiased, and asymptotically efficient.

## 6.3 The MLE for Mean and Covariance of Correlated Gaussian Variables

Suppose that we have a series of samples of a distribution that is Gaussian, but we do not know the mean or covariance of the distribution. Thus, our measurements are



$$\underline{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_M \end{bmatrix}, \quad \underline{y}_i = \underline{m} + \underline{v}, \quad \text{Cov}\{\underline{v}\} = R \quad (6.6)$$

and our states are

$$\underline{x} = \begin{bmatrix} \underline{m} \\ \text{(Elements of } R) \end{bmatrix} \quad (6.7)$$

The likelihood function is

$$p(\underline{y}|\underline{x}) = \prod_{j=1}^M p(y_j|\underline{x}) \quad (6.8)$$

where each set of measurements  $y_j$  has the probability distribution

$$p(y_j|\underline{x}) = \frac{1}{(2\pi)^{N/2} \cdot |R|^{M/2}} \cdot \exp\left(-\frac{1}{2} \cdot (\underline{y}_j - \underline{m})^T \cdot R^{-1} \cdot (\underline{y}_j - \underline{m})\right) \quad (6.9)$$

The log likelihood function is

$$L(\underline{x}) = -\frac{M \cdot N}{2} \cdot \ln(2\pi) - \frac{M}{2} \cdot \ln(|R|) - \frac{1}{2} \cdot \sum_{j=1}^M (\underline{y}_j - \underline{m})^T \cdot R^{-1} \cdot (\underline{y}_j - \underline{m}) \quad (6.10)$$

We find separate likelihood functions for  $\underline{m}$  and  $R$ , solving for  $\underline{m}$  first:

$$l(\hat{\underline{m}}) = \frac{\partial L(\underline{x})}{\partial \underline{m}} = \sum_{j=1}^M R^{-1} \cdot (\underline{y}_j - \hat{\underline{m}}) = \underline{0} \quad (6.11)$$

The solution to this one is simple:

$$\hat{\underline{m}} = \frac{1}{M} \cdot \sum_{j=1}^M y_j \quad (6.12)$$

We find the covariance  $R$  by using the property of consistency, and solving for  $R^{-1}$  instead. We do this taking the gradient of the log likelihood equation with respect to the matrix  $R^{-1}$ . We need a couple of identities from our work with gradients first, though.

See Gelb's Equation (2.1-75) page 23 for

$$\frac{\partial \ln(|R^{-1}|)}{\partial R^{-1}} = \frac{1}{|R^{-1}|} \cdot \frac{\partial |R^{-1}|}{\partial R^{-1}} = R \quad (6.13)$$

and Gelb's Equation (2.1-72) page 23 for

$$\frac{\partial (\underline{y}_i - \underline{m})^T \cdot R^{-1} \cdot (\underline{y}_i - \underline{m})}{\partial R^{-1}} = (\underline{y}_i - \underline{m}) \cdot (\underline{y}_i - \underline{m})^T \quad (6.14)$$

With these two equations, we have the likelihood equation for the covariance as

$$l(R^{-1}) = \frac{\partial L(\underline{x})}{\partial R^{-1}} = \frac{M}{2} \cdot \hat{R} - \frac{1}{2} \cdot \sum_{j=1}^M (\underline{y}_j - \hat{\underline{m}}) \cdot (\underline{y}_j - \hat{\underline{m}}) = 0 \quad (6.15)$$

This has the solution

$$\hat{R} = \frac{1}{M} \cdot \sum_{j=1}^M (\underline{y}_j - \hat{\underline{m}}) \cdot (\underline{y}_j - \hat{\underline{m}})^T \quad (6.16)$$

Note that the estimator for  $\hat{m}$  is linear, so it is statically efficient and unbiased. However the estimator for  $\hat{R}$  involves  $\hat{m}$  so that this part of the estimator is nonlinear. The inclusion of  $\hat{m}$  can be seen to make the sum have one less degree of freedom, and to make the variation in each sample smaller by a factor of  $(M-1)/M$ . A more rigorous analysis will bear out that an unbiased estimator of the covariance is

$$\hat{R}_{UNBIASED} = \frac{1}{M-1} \cdot \sum_{j=1}^M (\underline{y}_i - \hat{m}) \cdot (\underline{y}_i - \hat{m})^T \quad (6.17)$$

## 6.4 Covariance of Estimates

We can use (6.4) and (6.11) to obtain the Cramer-Rao bound for the covariance of  $\hat{m}$ .

The Fisher information matrix for  $\hat{m}$  is

$$F_M = -\frac{\partial^2 L(\underline{x})}{\partial \underline{m}^2} = M \cdot R^{-1} \quad (6.18)$$

from which we have the covariance of the estimate  $\hat{m}$ :

$$P_m = \frac{1}{M} \cdot R \approx \frac{1}{M} \cdot \hat{R} \quad (6.19)$$

Since the estimate of  $\underline{m}$  is linear, this covariance, the Cramer-Rao bound, is applicable because the estimator is efficient.

The remainder of the covariance is more complex. The estimator for  $\hat{R}^{-1}$  is nonlinear, and in fact the state vector contains both the elements of  $\underline{m}$  and  $R$ , a total of  $M^2+M$  elements, not just  $M$ .

## 7 The Kalman Update as an MLE

We can easily pose the Kalman update as an MLE by focusing on the measurement data and the extrapolated state vector. To avoid a notation clash here, we will denote the measurement for the Kalman filter update as  $\underline{yy}$  and reserve the standard notation for measurements  $\underline{y}$  as our measurements for the MLE. The data is then

$$\underline{y} = \begin{bmatrix} \underline{yy} \\ \tilde{x} \end{bmatrix} \quad (7.1)$$

and the covariance of the measurements is

$$\text{Cov}\{\underline{y}\} = \begin{bmatrix} R & 0 \\ 0 & \tilde{P} \end{bmatrix} \quad (7.2)$$

The mean values of the measurements are

$$E(\underline{y}) = \begin{bmatrix} \underline{h}(\underline{x}) \\ \underline{x} \end{bmatrix} \quad (7.3)$$

We can write the conditional probability as

$$p(\underline{y}, \tilde{\underline{x}} | \underline{x}) = \frac{1}{(2\pi)^{(M+N)/2} \cdot |R|^{1/2} \cdot |\tilde{P}|^{1/2}} \cdot \exp \left( -\frac{1}{2} \cdot \left( \begin{array}{l} (\underline{yy} - \underline{h}(\underline{x})) \cdot R^{-1} \cdot (\underline{yy} - \underline{h}(\underline{x})) \\ + (\tilde{\underline{x}} - \underline{x})^T \cdot \tilde{P}^{-1} \cdot (\tilde{\underline{x}} - \underline{x}) \end{array} \right) \right) \quad (7.4)$$

The log likelihood equation is

$$L(\underline{x}) = -\frac{M+N}{2} \cdot \ln(\pi) - \frac{1}{2} \cdot \ln(|R| \cdot |\tilde{P}|) - \frac{1}{2} \cdot \left( (\underline{yy} - \underline{h}(\underline{x})) \cdot R^{-1} \cdot (\underline{yy} - \underline{h}(\underline{x})) + (\tilde{\underline{x}} - \underline{x})^T \cdot \tilde{P}^{-1} \cdot (\tilde{\underline{x}} - \underline{x}) \right) \quad (7.5)$$

The likelihood equation is

$$l(\hat{\underline{x}}) = \frac{\partial L(\hat{\underline{x}})}{\partial \hat{\underline{x}}} = H^T \cdot R^{-1} \cdot (\underline{yy} - \underline{h}(\hat{\underline{x}})) + \tilde{P} \cdot (\tilde{\underline{x}} - \hat{\underline{x}}) = \underline{0} \quad (7.6)$$

We note that we can use the first two terms of a Taylor series,

$$\underline{h}(\hat{\underline{x}}) = \underline{h}(\tilde{\underline{x}}) + H \cdot (\hat{\underline{x}} - \tilde{\underline{x}}) \quad (7.7)$$

and the likelihood equation becomes

$$l(\hat{\underline{x}}) = -(\tilde{P}^{-1} + H^T \cdot R^{-1} \cdot H) \cdot (\hat{\underline{x}} - \tilde{\underline{x}}) + H^T \cdot R^{-1} \cdot (\underline{yy} - \underline{h}(\tilde{\underline{x}})) \quad (7.8)$$

or,

$$(H \cdot \tilde{P} \cdot H^T + \tilde{R}) \cdot (\hat{\underline{x}} - \tilde{\underline{x}}) = H^T \cdot R^{-1} \cdot (\underline{yy} - \underline{h}(\tilde{\underline{x}})) \quad (7.9)$$

from which we get the Kalman update:

$$\hat{\underline{x}} = \tilde{\underline{x}} + K \cdot (\underline{yy} - \underline{h}(\tilde{\underline{x}})) \quad (7.10)$$

where the Kalman gain  $K$  is

$$K = (\tilde{P}^{-1} + H^T \cdot R^{-1} \cdot H)^{-1} \cdot H^T \cdot R^{-1} \quad (7.11)$$

The covariance  $P$  of the estimate is found from

$$\frac{\partial^2 l(\underline{x})}{\partial \underline{x}^2} = -P^{-1} = -\tilde{P}^{-1} - H^T \cdot R^{-1} \cdot H \quad (7.12)$$

The update, Kalman gain, and covariance of the estimate are as given in the Kalman filter equations.

Note that the Kalman filter extrapolation is an application of the linear covariance propagation equations. Thus, between those and the development of the update equations that we have here, we have a complete derivation of the Kalman filter equations.

## 8 Estimation Theory Library

### 8.1 Foundations of Modern Estimation Theory

- 1) R. A. Fisher, *Theory of Statistical Estimation*, Proc. Cambridge Phil. Soc., v. 22 p. 700 (1925).
- 2) R. A. Fisher, *On the Mathematical Foundations of Theoretical Statistics*, Phil. Trans. Roy. Soc., London, v. 222, p. 309 (1922).
- 3) R. A. Fisher, *Two New Properties of Mathematical Likelihood*, Proc. Roy. Soc., London, v. 144, p. 285 (1934).

- 4) R. A. Fisher, *The Logic of Inductive Inference*, J. Roy. Statist. Soc., v. 98, p. 39 (1935).
- 5) C. R. Rao, *Information and Accuracy Attainable in the Estimation of Statistical Parameters*, Bull. Calcutta Math. Soc., v. 37, pp 81-91 (1945).
- 6) Harold Cramer, *Mathematical Methods of Statistics*, Princeton University Press (1946).
- 7) A. Bhattacharyya, "On Some Analogues of the Amount of Information and their Use in Statistical Estimation, *Sankhya*, v. 8, pp, 1, 201, 314 (1946, 1947, 1948).

## 8.2 Major References

- 8) Richard Bellman, *Introduction to Matrix Analysis*, Second Edition, SIAM Press (1997) (reprint from 1984; first edition was in 1960).
- 9) M. G. Kendall and A. Stuart, *The Advanced Theory of Statistics, Vol. 2, Inference and Relationship*, Hafner, New York (1961).
- 10) Harry L. Van Trees, *Detection, Estimation, and Modulation Theory Part I*, John Wiley (1968).
- 11) Harold W. Sorenson, *Parameter Estimation*, Marcel Dekker (1980).
- 12) M. B. Priestley, *Spectral Analysis and Time Series*, Academic Press (1981).
- 13) Louis L. Scharf, *Statistical Signal Processing*, Addison-Wesley (1991)