

## Lecture 18: Maximum Likelihood Parameter Estimation

### Forms of Parameter Estimation [1, Ch.4]

We continue our discussion of parameter estimation for dynamical systems by first considering the simpler parameter estimation for a static nonlinear model of the form

$$\mathbf{y} = \varphi(\boldsymbol{\theta}) + \mathbf{v} \quad (1)$$

where  $\boldsymbol{\theta}$  is an unknown parameter vector,  $\varphi(\boldsymbol{\theta})$  is the (noise-free) output, and  $\mathbf{y}$  is a measurement that consists of the output corrupted by noise  $\mathbf{v}$ . As a special case of the above nonlinear function we have the linear model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta} + \mathbf{v} \quad (2)$$

Morreli and Klein [1] discuss three forms parameter estimation models:

#### 1. The Least-Squares Model

- (a)  $\boldsymbol{\theta}$  is a constant (but unknown) vector of parameters (i.e., it has no p.d.f.)
- (b)  $\mathbf{v}$  is a random vector of measurement noise which may or may not have a p.d.f. The p.d.f. is not needed for the least-square estimate but can be used to compute a confidence, if available.
- (c) The least-square estimate is

$$\hat{\boldsymbol{\theta}} = \min_{\boldsymbol{\theta}} \sum_{i=1}^N [\mathbf{y}_i - \varphi(\boldsymbol{\theta})]^2 \quad (3)$$

for  $\mathbf{y}_i$  data points.

#### 2. The Fisher Model

- (a)  $\boldsymbol{\theta}$  is a constant (but unknown) vector of parameters (i.e., it has no p.d.f.)
- (b)  $\mathbf{v}$  is a random vector with a probability distribution function  $p(\mathbf{v})$
- (c) The constant vector  $\boldsymbol{\theta}$  is viewed as parameter in the p.d.f. that describes the distribution of the  $N$  data points  $\mathbf{Y}_{1:N}$  according to the likelihood function

$$\mathcal{L}(\mathbf{Y}_{1:N}; \boldsymbol{\theta}) = p(\mathbf{Y}_{1:N} | \boldsymbol{\theta}) \quad (4)$$

- (d) The maximum likelihood estimator (MLE) is

$$\hat{\boldsymbol{\theta}} = \max_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{Y}_{1:N}; \boldsymbol{\theta}) \quad (5)$$

#### 3. The Bayesian Model

- (a)  $\boldsymbol{\theta}$  is a random vector with a probability distribution  $p(\boldsymbol{\theta})$
- (b)  $\mathbf{v}$  is a random vector with a probability distribution  $p(\mathbf{v})$
- (c) The p.d.f. of the parameter  $\boldsymbol{\theta}$  is related to the data  $\mathbf{y}$  by Bayes' Rule:

$$p(\boldsymbol{\theta} | \mathbf{Y}_{1:N}) = \frac{p(\mathbf{Y}_{1:N} | \boldsymbol{\theta}) p(\boldsymbol{\theta})}{p(\mathbf{Y}_{1:N})} \quad (6)$$

(d) The Bayesian estimator is

$$\hat{\theta} = \max_{\theta} p(\theta | \mathbf{Y}_{1:N}) \quad (7)$$

We've already discussed least-square regression and in this lecture we will introduce maximum likelihood estimation. A future lecture will describe Bayesian estimation.

## Maximum Likelihood Estimation of a Static Parameter

The *likelihood function*

$$\mathcal{L}(\mathbf{y}; \theta) = p(\mathbf{y} | \theta) \quad (8)$$

is the probability of the data/measurements  $\mathbf{y}$  for a particular choice of parameters  $\theta$ . Now, consider a system that produces a sequence of measurements  $\{\mathbf{y}_1, \dots, \mathbf{y}_N\}$  and let  $\mathbf{Y}_{1:N} = \{\mathbf{y}_1, \dots, \mathbf{y}_N\}$  denote the set of  $N$  measurements. The likelihood function that considers *all* these measurements is denoted

$$\mathcal{L}(\mathbf{Y}_{1:N}; \theta) = \mathcal{L}(\{\mathbf{y}_1, \dots, \mathbf{y}_N\}; \theta) = p(\mathbf{Y}_{1:N} | \theta) \quad (9)$$

Split the set  $\mathbf{Y}_{1:N}$  into two subsets:

$$\begin{aligned} \mathbf{Y}_{1:N} &= (\mathbf{y}_N, \{\mathbf{y}_1, \dots, \mathbf{y}_{N-1}\}) \\ &= (\mathbf{y}_N, \mathbf{Y}_{1:N-1}) \end{aligned}$$

**Aside:** Recall the relationship between joint and conditional probabilities

$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x} | \mathbf{y}) p(\mathbf{y}) \quad (10)$$

The joint probability  $p(\mathbf{y}_N, \mathbf{Y}_{1:N-1}; \theta)$  can be re-written as

$$p(\mathbf{y}_N, \mathbf{Y}_{1:N-1}; \theta) = p(\mathbf{y}_N | \mathbf{Y}_{1:N-1}; \theta) p(\mathbf{Y}_{1:N-1}; \theta) \quad (11)$$

or, in terms of likelihood functions, it may be factored further as:

$$\begin{aligned} \mathcal{L}(\mathbf{Y}_{1:N}; \theta) &= \mathcal{L}(\{\mathbf{y}_1, \dots, \mathbf{y}_N\}; \theta) \\ &= \mathcal{L}(\mathbf{y}_N | \mathbf{Y}_{1:N-1}; \theta) \mathcal{L}(\mathbf{Y}_{1:N-1}; \theta) \\ &= \mathcal{L}(\mathbf{y}_N | \mathbf{Y}_{1:N-1}; \theta) \mathcal{L}(\mathbf{y}_{N-1} | \mathbf{Y}_{1:N-2}; \theta) \mathcal{L}(\mathbf{Y}_{1:N-2}; \theta) \\ &\quad \vdots \\ &= \prod_{i=1}^N \mathcal{L}(\mathbf{y}_i | \mathbf{Y}_{1:i-1}; \theta) \end{aligned}$$

It is advantageous to consider minimizing the negative logarithm of this function, rather than

maximize the likelihood. The maximum likelihood (ML) estimator is

$$\begin{aligned}\hat{\theta} &= \max_{\theta} \mathcal{L}(\mathbf{Y}_{1:N}; \theta) \\ &= \max_{\theta} \prod_{i=1}^N \mathcal{L}(\mathbf{y}_i | \mathbf{Y}_{1:i-1}; \theta) \\ &= \min_{\theta} \underbrace{\sum_{i=1}^N -\ln[\mathcal{L}(\mathbf{y}_i | \mathbf{Y}_{1:i-1}; \theta)]}_{=J(\theta)}\end{aligned}$$

If we further assume that the noise is uncorrelated in time then the likelihood functions also become independent:

$$\mathcal{L}(\mathbf{y}_i | \mathbf{Y}_{1:i-1}; \theta) = \mathcal{L}(\mathbf{y}_i; \theta) \quad (12)$$

and the ML estimator is:

$$\hat{\theta}_{\text{MLE}} = \min_{\theta} \underbrace{\sum_{i=1}^N -\ln[\mathcal{L}(\mathbf{y}_i; \theta)]}_{=J(\theta)}$$

The above cost function can be minimized analytically if it is simple enough, otherwise we may resort to numerical methods. One advantage of MLE is that it allows us to consider non-Gaussian distributions, as in the following example.

**Example (Rayleigh distribution).** The distribution of wave heights in the ocean is commonly modeled using a *Rayleigh* probability distribution that has probability density function (p.d.f.):

$$f(y; \theta) = \frac{y}{\theta^2} e^{-y^2/(2\theta^2)} \quad (13)$$

where  $\theta$  is a scale parameter that describes the shape of the distribution. The distribution is restricted to the positive real numbers  $y \in [0, \infty]$ . Now suppose that we obtain a set of data

$$Y = \{y_1, y_2, \dots, y_N\} \quad (14)$$

that represents wave height measurements and we wish to determine estimate the parameter  $\theta$  that best matches the data. The negative log-likelihood is

$$\begin{aligned}-\ln[\mathcal{L}(y; \theta)] &= -\ln\left[\frac{y}{\theta^2} e^{-y^2/(2\theta^2)}\right] \\ &= -\ln[y] + 2\ln[\theta] + [y^2/(2\theta^2)]\end{aligned}$$

The MLE minimizes the cost function

$$\begin{aligned}J(\theta) &= \sum_{i=1}^N \{-\ln[y_i] + 2\ln[\theta] + [y_i^2/(2\theta^2)]\} \\ &= -\sum_{i=1}^N \ln[y_i] + 2N\ln[\theta] + \frac{1}{2\theta^2} \sum_{i=1}^N y_i^2\end{aligned}$$

which considers all of the data. To find the minimum we set the partial derivative to zero

$$\begin{aligned}
 0 &= \frac{\partial J(\theta)}{\partial \theta} = \frac{\partial}{\partial \theta} \left\{ - \sum_{i=1}^N \ln [y_i] \right\} + 2N \frac{\partial}{\partial \theta} \ln [\theta] + \frac{\partial}{\partial \theta} \frac{1}{2\theta^2} \sum_{i=1}^N y_i^2 \\
 0 &= 0 + \frac{2N}{\hat{\theta}_{MLE}} - \frac{1}{\hat{\theta}_{MLE}^3} \sum_{i=1}^N y_i^2 \\
 2N \hat{\theta}_{MLE}^2 &= \sum_{i=1}^N y_i^2 \\
 \Rightarrow \hat{\theta}_{MLE}^2 &= \sqrt{\frac{1}{2N} \sum_{i=1}^N y_i^2}
 \end{aligned}$$

In the above example the parameter of interest characterized the probability distribution, however the same approach can be taken for any parameters that appear in the model, or both the model parameters and noise parameters. Consider the following example.

**Example (MLE for Scalar Linear Model with Gaussian Noise).** Consider a scalar measurement model of the form:

$$\begin{aligned}
 y &= mx + b + v \\
 v &\sim \mathcal{N}(0, \sigma^2)
 \end{aligned} \tag{15}$$

where  $m$ ,  $b$ , and  $\sigma^2$  are unknown model parameters. A data set is collected consisting of pairs

$$\{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\} \tag{16}$$

The vector of parameters to be estimated is

$$\hat{\theta} = [\hat{m}, \hat{b}, \hat{\sigma}^2]^T \tag{17}$$

The parameters  $m$  and  $b$  are unknown but they are deterministic. Thus, the probability of measuring some  $y_i$  given an input  $x_i$  is equivalent to the probability that the noise  $v_i$  takes on the particular value

$$v_i = y_i - (mx_i + b) \tag{18}$$

Since the noise is normally distributed the probability can be expressed using the Gaussian p.d.f. function:

$$\mathcal{L}(y; m, b, \sigma^2) = p(v|m, b, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{-v^2}{2\sigma^2}\right) \tag{19}$$

$$= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{-(z - mx - b)^2}{2\sigma^2}\right) \tag{20}$$

The negative log-likelihood is:

$$-\ln [\mathcal{L}(y; m, b, \sigma^2)] = -\ln \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left( \frac{-(y - mx - b)^2}{2\sigma^2} \right) \right] \quad (21)$$

$$= - \left( \frac{-(y - mx - b)^2}{2\sigma^2} \right) + \ln \sqrt{2\pi\sigma^2} \quad (22)$$

$$= - \left( \frac{-(y - mx - b)^2}{2\sigma^2} \right) + \frac{1}{2} (\ln 2\pi + \ln \sigma^2) \quad (23)$$

$$= - \left( \frac{-(y - mx - b)^2}{2\sigma^2} \right) + \frac{1}{2} \ln 2\pi + \frac{1}{2} \ln \sigma^2 \quad (24)$$

The MLE minimizes the cost function

$$J(\theta) = \sum_{i=1}^N \left[ - \left( \frac{-(y_i - mx_i - b)^2}{2\sigma^2} \right) + \frac{1}{2} \ln 2\pi + \frac{1}{2} \ln \sigma^2 \right] \quad (25)$$

$$= - \sum_{i=1}^N \left( \frac{-(y_i - mx_i - b)^2}{2\sigma^2} \right) + \frac{N}{2} \ln 2\pi + \frac{N}{2} \ln \sigma^2 \quad (26)$$

which considers all of the data. To find the minimum we set the partial derivative with respect to each parameter to zero. Begin with the partial with respect to  $b$ :

$$0 = \frac{\partial J(\theta)}{\partial b} = - \frac{\partial}{\partial b} \sum_{i=1}^N \left( \frac{-(y_i - mx_i - b)^2}{2\sigma^2} \right) \quad (27)$$

$$0 = - \sum_{i=1}^N \left( \frac{-2(y_i - \hat{m}x_i - \hat{b})}{2\sigma^2} \right) (-1) \quad (28)$$

$$0 = \sum_{i=1}^N (y_i - \hat{m}x_i - \hat{b}) \quad (29)$$

$$\hat{b} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{m}x_i) \quad (30)$$

$$\implies \hat{b} = \bar{y} - \hat{m}\bar{x} \quad (31)$$

where  $\bar{x}$  and  $\bar{y}$  denote the mean values of the  $x$  and  $y$  data points (i.e., the *sample mean*). Next, evaluate the partial derivative with respect to  $m$ :

$$0 = \frac{\partial J(\theta)}{\partial m} = - \frac{\partial}{\partial m} \sum_{i=1}^N \left( \frac{-(y_i - \hat{m}x_i - \hat{b})^2}{2\sigma^2} \right) \quad (32)$$

$$0 = \sum_{i=1}^N \left( \frac{2(y_i - \hat{m}x_i - \hat{b})}{2\sigma^2} \right) (-x_i) \quad (33)$$

$$0 = \sum_{i=1}^N (y_i x_i - \hat{m}x_i^2 - \hat{b}x_i) \quad (34)$$

the substituting  $\hat{b}$

$$0 = \sum_{i=1}^N (y_i x_i - \hat{m} x_i^2 - (\bar{z} - \hat{m} \bar{x}) x_i) \quad (35)$$

$$0 = \sum_{i=1}^N y_i x_i - \hat{m} \sum_{i=1}^N x_i^2 - \bar{z} \sum_{i=1}^N x_i + \hat{m} \bar{x} \sum_{i=1}^N x_i \quad (36)$$

$$0 = \sum_{i=1}^N y_i x_i - \hat{m} \sum_{i=1}^N x_i^2 - N \bar{z} \bar{x} + N \hat{m} \bar{x}^2 \quad (37)$$

$$\hat{m} = \frac{\sum_{i=1}^N y_i x_i - N \bar{z} \bar{x}}{\sum_{i=1}^N x_i^2 - N \bar{x}^2} \quad (38)$$

$$\hat{m} = \frac{\sum_{i=1}^N y_i x_i - \sum_{i=1}^N \bar{z} \bar{x}}{\sum_{i=1}^N x_i^2 - \sum_{i=1}^N \bar{x}^2} \quad (39)$$

$$\hat{m} = \frac{\sum_{i=1}^N (y_i x_i - \bar{z} \bar{x})}{\sum_{i=1}^N (x_i^2 - \bar{x}^2)} \quad (40)$$

$$\Rightarrow \hat{m} = \frac{\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{z})}{\sum_{i=1}^N (x_i - \bar{x})^2} \quad (41)$$

Notice tha the numerator of this expression is equivalent to the *sample covariance* between  $x$  and  $y$  and the denominator is the sample variance of  $x$ . Lastly, take the partial derivative with respect to  $\sigma^2$ :

$$0 = \frac{\partial J(\theta)}{\partial \sigma^2} = -\frac{\partial}{\partial \sigma^2} \sum_{i=1}^N \left( \frac{-(y_i - m x_i - b)^2}{2\sigma^2} \right) + \frac{\partial}{\partial \sigma^2} \frac{N}{2} \ln \sigma^2 \quad (42)$$

$$0 = \sum_{i=1}^N \left( \frac{-(y_i - m x_i - b)^2}{2\sigma^4} \right) + \frac{N}{2\sigma^2} \quad (43)$$

$$0 = -\sum_{i=1}^N (y_i - m x_i - b)^2 + N \sigma^2 \quad (44)$$

$$\Rightarrow \hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{m} x_i - \hat{b})^2 \quad (45)$$

Together the three equations (31), (41), and (45) form the maximum likelihood estimate.

**Additive Gaussian Measurement Noise.** The example above which considers scalar Gaussian measurement noise can be extended to a more general vector-valued and nonlinear measurement model of the form

$$\begin{aligned} \mathbf{y} &= \varphi(\mathbf{x}; \boldsymbol{\theta}) + \mathbf{v} \\ \mathbf{v} &\sim \mathcal{N}(\mathbf{0}, \mathbf{R}) \end{aligned} \quad (46)$$

As before, the probability of observing some data point  $\mathbf{y}$  is equivalent to the probability that the noise takes on the value

$$\mathbf{v} = \mathbf{y} - \varphi(\mathbf{x}; \boldsymbol{\theta}) \quad (47)$$

for a given set of parameters  $\theta$ . The likelihood function is given by the standard Gaussian multivariate p.d.f. for the probability that  $\mathbf{v} = \mathbf{y} - \varphi(\mathbf{x}; \theta)$ . That is,

$$\mathcal{L}(\mathbf{y}; \theta) = p(\mathbf{y}|\theta) = \frac{\exp\left(-\frac{1}{2}(\mathbf{y} - \varphi(\mathbf{x}; \theta))^T \mathbf{R}^{-1}(\mathbf{y} - \varphi(\mathbf{x}; \theta))\right)}{(2\pi)^{n/2} \sqrt{|\mathbf{R}|}} \quad (48)$$

If we compute the negative log likelihood (and ignore parts of the cost function that are not dependent on  $\theta$ ) we end up with the cost function:

$$J(\theta) = \sum_{i=1}^N \frac{1}{2} (\mathbf{y}_i - \varphi(\mathbf{x}; \theta))^T \mathbf{R}^{-1} (\mathbf{y}_i - \varphi(\mathbf{x}; \theta)) \quad (49)$$

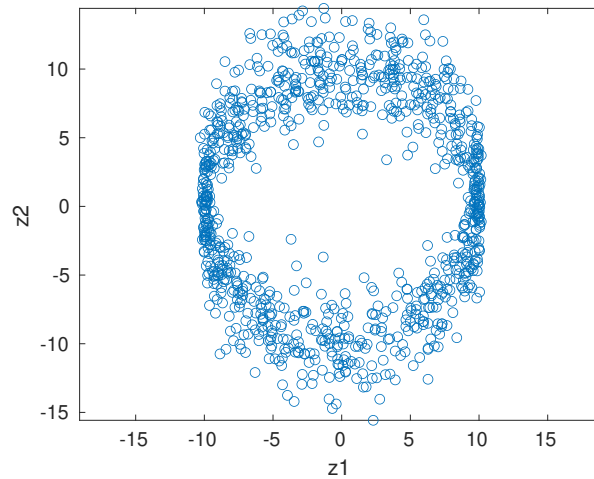
that appears within the exponent in the p.d.f. The optimal  $\hat{\theta}$  can be solved for using various optimization procedures such as gradient descent. For low dimensional  $\theta$  a coarse optimization may be possible by brute force (i.e., evaluating  $J(\theta)$  for a grid of  $\theta$  values and identifying the minimum).

**Example.** Consider a nonlinear vector-valued measurement model of the form:

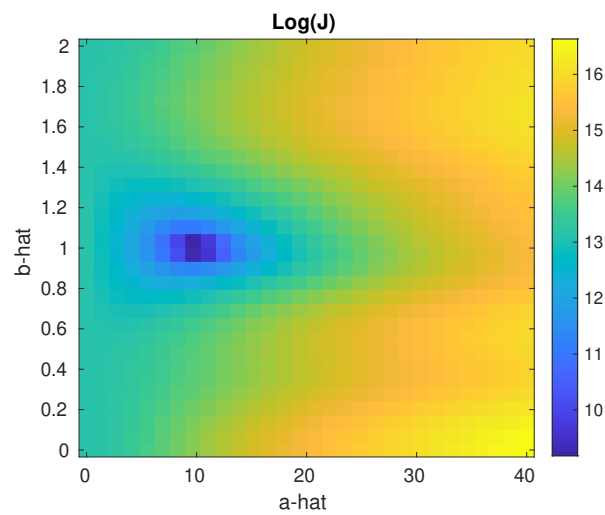
$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} a \cos bx_1 \\ a \sin bx_1 \end{bmatrix} + \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} \quad (50)$$

$$\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{R})$$

where  $a$  and  $b$  are unknown model parameters (with true values  $a = 10$  and  $b = 1$ ) to be estimated and  $\mathbf{R} = \text{diag}([\sigma_1^2, \sigma_2^2])$  with known variances  $\sigma_1^2 = 0.1$  and  $\sigma_2^2 = 2$ . An example data set of  $N = 1000$  points was generated by varying  $x_1 \in [0, 2\pi]$  as shown below:



The data resembles a noisy circle that is “spread out” along the  $z_2$  axis since  $\sigma_2 > \sigma_1$ . The cost function is defined by computing  $J(\theta)$  with all  $N$  data points and is plotted below for a  $M \times M$  search grid with  $\hat{a} \in [0, 40]$  and  $\hat{b} \in [0, 2]$  with  $M = 30$ . To emphasize changes in the cost function the logarithm of  $J$  is plotted below.



As shown, the cost function has a minimum near the correct values of  $a = 10$  and  $b = 1$ . In this example our actual estimate  $\hat{a}$  and  $\hat{b}$  will depend on the resolution of the search grid.



## Maximum Likelihood Estimation in A Dynamical System

The above examples discuss maximum likelihood estimation for a static system model. How can we extend this to find the parameters of a dynamical system using the output data? The parameters of a dynamic system can appear anywhere in the state equation, measurement equation, or the noise covariance matrices. Maximum likelihood estimation provides an approach for estimating these parameters. The most common methods described in [1, Ch.4] methods are the (a) *filter-error method*, (b) the *output-error method*, and (c) the *equation-error method*. A block diagram of each methods is shown below. (Note: the block diagram is for a particular system “aircraft” but the approach is more general.)

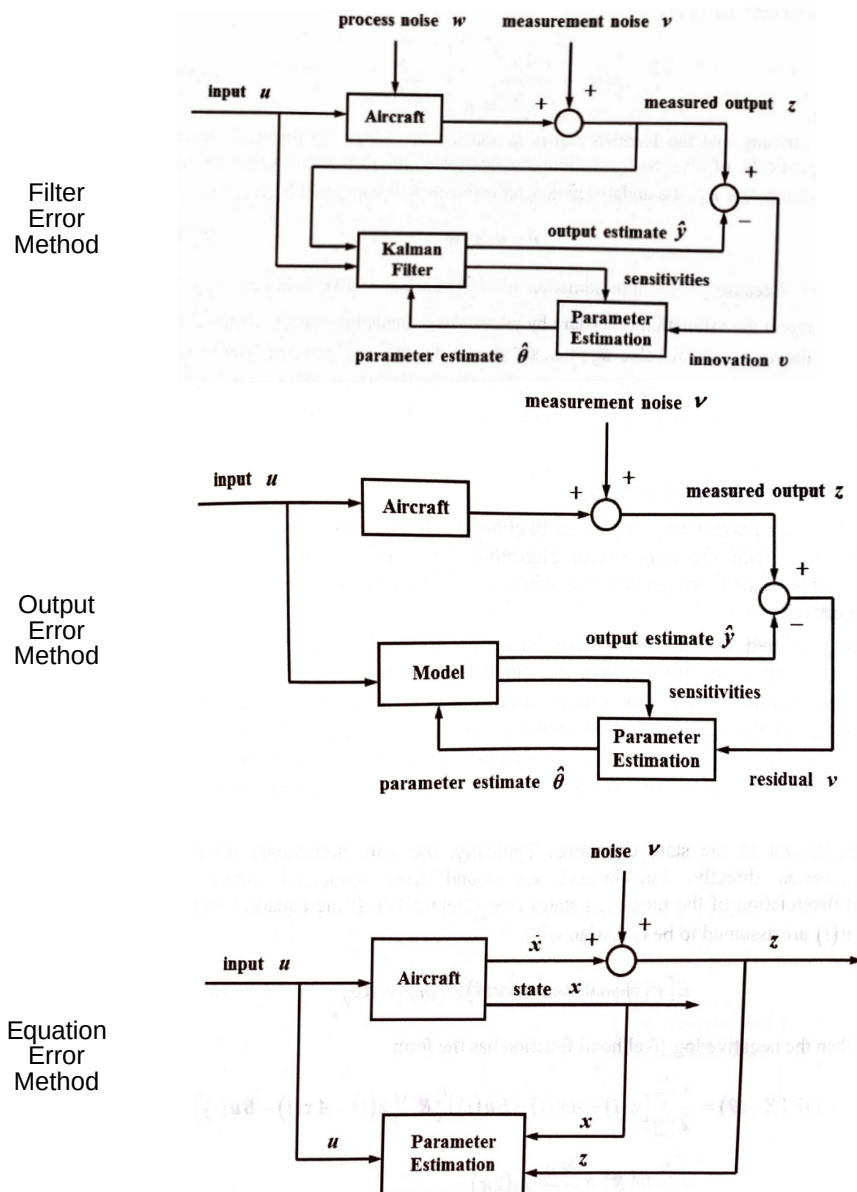


Image source: [1, Ch.4]

The filter-error is the most general form and works in tandem with a Kalman filter, as we will show below. The output-error considers a simplified system model in which the process noise is

ignored. Thus the state equation is entirely deterministic. Lastly, notice that the equation error model has no output/measurement equation. That is, this method assumes all of the states and all of the state derivatives are available. In this lecture we will focus on the filter-error method as described next.

### Filter-error Parameter Estimation

Consider a LTI discrete system in the form

$$\mathbf{x}_k = \mathbf{F}\mathbf{x}_{k-1} + \mathbf{G}\mathbf{u}_{k-1} + \mathbf{w}_{k-1} \quad (51)$$

$$\mathbf{y}_k = \mathbf{H}\mathbf{x}_k + \mathbf{v}_k \quad (52)$$

$$\mathbf{w}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}) \quad (53)$$

$$\mathbf{v}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{R}) \quad (54)$$

where the matrices  $\mathbf{F}$ ,  $\mathbf{G}$ ,  $\mathbf{Q}$ ,  $\mathbf{H}$ , and  $\mathbf{R}$  may be a function of some unknown parameters. Following our approach for static system maximum likelihood parameter estimation, we consider a sequence of noisy measurements  $\{\mathbf{y}_1, \dots, \mathbf{y}_N\}$  and let  $\mathbf{Y}_{1:N} = [\mathbf{y}_1^T, \dots, \mathbf{y}_N^T]^T$  denote the a vector that stacks these  $N$  measurements into a single column. The likelihood function that considers *all* these measurements is denoted

$$\mathcal{L}(\mathbf{Y}_{1:N}; \boldsymbol{\theta}) = \mathcal{L}([\mathbf{y}_1^T, \dots, \mathbf{y}_N^T]^T; \boldsymbol{\theta}) \quad (55)$$

As before, the MLE  $\hat{\boldsymbol{\theta}}$  minimizes the cost function of the negative log-likelihood (assuming uncorrelated measurements):

$$J(\boldsymbol{\theta}) = -\ln[\mathcal{L}(\mathbf{Y}_{1:N}; \boldsymbol{\theta})] = \sum_{i=1}^N -\ln[\mathcal{L}(\mathbf{y}_i; \boldsymbol{\theta})]$$

That is, the maximum likelihood (ML) estimator is

$$\hat{\boldsymbol{\theta}} = \min_{\boldsymbol{\theta}} \underbrace{\sum_{i=1}^N -\ln[\mathcal{L}(\mathbf{y}_i; \boldsymbol{\theta})]}_{=J(\boldsymbol{\theta})}$$

We can view each measurement  $\mathbf{y}_k$  as a random vector

$$\mathbf{y}_k \sim \mathcal{N}(\mathbf{y}_{k|k-1}, \mathbf{S}_k) \quad (56)$$

with mean given by the measurement  $\mathbf{y}_{k|k-1}$  predicted from a Kalman filter and the variance given as the innovation  $\mathbf{S}_k$ . That is,

$$\mathbf{y}_{k|k-1} = E[\mathbf{y}_k] \quad (57)$$

$$\mathbf{S}_k = E[(\mathbf{y}_k - \mathbf{y}_{k|k-1})(\mathbf{y}_k - \mathbf{y}_{k|k-1})^T] \quad (58)$$

$$= E[\boldsymbol{\epsilon}_k \boldsymbol{\epsilon}_k^T] \quad (59)$$

where

$$\boldsymbol{\epsilon}_k = \mathbf{y}_k - \mathbf{y}_{k|k-1} \quad (60)$$

is called the *innovation* vector and  $\mathbf{S}_k$  is the innovation covariance. As the number of measurements increases this assumption is justified and the p.d.f. of the innovation approaches a Gaussian distribution. The quantities (57) and (58) that describe this likelihood are provided by the Kalman filter. Recall that in the Kalman filter we perform the following steps at each iteration:

$$\text{State Prior : } \hat{\mathbf{x}}_{k-1|k-1} \quad (61)$$

$$\text{Covariance Prior : } \mathbf{P}_{k-1|k-1} \quad (62)$$

$$\text{Current Measurement : } \mathbf{y}_k \quad (63)$$

$$\text{State Prediction : } \hat{\mathbf{x}}_{k|k-1} = \mathbf{F}\hat{\mathbf{x}}_{k-1|k-1} + \mathbf{G}\mathbf{u}_{k-1} \quad (64)$$

$$\text{Covariance Prediction : } \mathbf{P}_{k|k-1} = \mathbf{F}\mathbf{P}_{k-1|k-1}\mathbf{F}^T + \mathbf{Q} \quad (65)$$

$$\text{Measurement Prediction : } \mathbf{y}_{k|k-1} = \mathbf{H}\hat{\mathbf{x}}_{k|k-1} \quad (66)$$

$$\text{Innovation Covariance : } \mathbf{S}_k = \mathbf{H}_k\mathbf{P}_{k|k-1}\mathbf{H}_k^T + \mathbf{R}_k \quad (67)$$

$$\text{Kalman Gain : } \mathbf{K}_k = \mathbf{P}_{k|k-1}\mathbf{H}_k^T\mathbf{S}_k^{-1} \quad (68)$$

$$\text{State Posterior : } \hat{\mathbf{x}}_{k|k} = \hat{\mathbf{x}}_{k|k-1} + \mathbf{K}_k(\mathbf{y}_k - \mathbf{y}_{k|k-1}) \quad (69)$$

$$\text{Covariance Posterior : } \mathbf{P}_{k|k} = (\mathbf{I} - \mathbf{K}_k\mathbf{H}_k)\mathbf{P}_{k|k-1} \quad (70)$$

Thus (57) and (58) are found from (66) and (67). Having characterized the mean and covariance of the Gaussian random vector (56) we may compute the likelihood:

$$\mathcal{L}(\mathbf{y}_k; \boldsymbol{\theta}) = p(\mathbf{y}|\boldsymbol{\theta}) = \frac{\exp\left(-\frac{1}{2}(\mathbf{y}_k - \mathbf{y}_{k|k-1})^T \mathbf{S}_k^{-1}(\mathbf{y}_k - \mathbf{y}_{k|k-1})\right)}{(2\pi)^{m/2} \sqrt{|\mathbf{S}_k|}} \quad (71)$$

where  $m$  is the size of the output. The negative log likelihood (across all  $N$  data points) is:

$$-\ln[\mathcal{L}(\mathbf{Y}_N; \boldsymbol{\theta})] = -\sum_{i=1}^N \left[ \left( -\frac{1}{2}(\mathbf{y}_i - \mathbf{y}_{i|i-1})^T \mathbf{S}_i^{-1}(\mathbf{y}_i - \mathbf{y}_{i|i-1}) \right) - \ln \sqrt{|\mathbf{S}_i|} - \ln(2\pi)^{m/2} \right] \quad (72)$$

$$= -\sum_{i=1}^N \left[ \left( -\frac{1}{2}(\mathbf{y}_i - \mathbf{y}_{i|i-1})^T \mathbf{S}_i^{-1}(\mathbf{y}_i - \mathbf{y}_{i|i-1}) \right) - \frac{1}{2} \ln |\mathbf{S}_i| - \frac{m}{2} \ln(2\pi) \right] \quad (73)$$

$$= \frac{1}{2} \sum_{i=1}^N \left[ \boldsymbol{\epsilon}_i^T \mathbf{S}_i^{-1} \boldsymbol{\epsilon}_i + \ln |\mathbf{S}_i| \right] + \frac{mN}{2} \ln(2\pi) \quad (74)$$

The last term is a constant that does not affect the optimization. Thus the cost function to be minimized is

$$J(\boldsymbol{\theta}) = \sum_{i=1}^N \left[ \boldsymbol{\epsilon}_i^T \mathbf{S}_i^{-1} \boldsymbol{\epsilon}_i + \ln |\mathbf{S}_i| \right] \quad (75)$$

where  $\mathbf{S}_i$  and  $\boldsymbol{\epsilon}_i$  depend on the unknown parameters  $\boldsymbol{\theta}$ . To minimize this cost function we may use several techniques, for example, gradient descent. Implementing gradient descent requires computing the partial derivative of the cost with respect to each parameter. The cost depends on the Kalman filter equations given above. A review of how these partial derivatives can be computed (analytically or numerically) is presented in [2, 11.1].

### Maximum Likelihood Estimate via Gradient Descent

One approach is to approximate  $J(\theta)$  as a Taylor series and use a Newton-Raphson scheme (i.e., gradient descent). Assume that  $\theta$  can be expressed as a small perturbation from a nominal estimate (i.e.,  $\theta = \theta_0 + \Delta\theta$ ). Then,

$$J(\theta_0 + \Delta\theta) = J(\theta_0) + \Delta\theta^T \frac{\partial J}{\partial \theta} \Big|_{\theta=\theta_0} + \frac{1}{2} \Delta\theta^T \frac{\partial^2 J}{\partial \theta \partial \theta^T} \Big|_{\theta=\theta_0} \Delta\theta + \text{h.o.t.} \quad (76)$$

$$\approx J(\theta_0) + \Delta\theta^T \frac{\partial J}{\partial \theta} \Big|_{\theta=\theta_0} + \frac{1}{2} \Delta\theta^T \frac{\partial^2 J}{\partial \theta \partial \theta^T} \Big|_{\theta=\theta_0} \Delta\theta \quad (77)$$

A necessary condition for  $J(\theta_0 + \Delta\theta)$  to be a minimum is

$$\frac{\partial}{\partial \Delta\theta} [J(\theta_0 + \Delta\theta)] = 0 \quad (78)$$

Applying this necessary condition to (77):

$$0 = \frac{\partial}{\partial \Delta\theta} [J(\theta_0 + \Delta\theta)] = \frac{\partial}{\partial \Delta\theta} \left( J(\theta_0) + \Delta\theta^T \frac{\partial J}{\partial \theta} \Big|_{\theta=\theta_0} + \frac{1}{2} \Delta\theta^T \frac{\partial^2 J}{\partial \theta \partial \theta^T} \Big|_{\theta=\theta_0} \Delta\theta \right) \quad (79)$$

$$0 = 0 + \frac{\partial}{\partial \Delta\theta} \left( \Delta\theta^T \frac{\partial J}{\partial \theta} \Big|_{\theta=\theta_0} \right) + \frac{\partial}{\partial \Delta\theta} \left( \frac{1}{2} \Delta\theta^T \frac{\partial^2 J}{\partial \theta \partial \theta^T} \Big|_{\theta=\theta_0} \Delta\theta \right) \quad (80)$$

$$0 = \frac{\partial J}{\partial \theta} \Big|_{\theta=\theta_0} + \frac{\partial^2 J}{\partial \theta \partial \theta^T} \Big|_{\theta=\theta_0} \Delta\theta \quad (81)$$

Solving this last equation gives a vector of parameter changes

$$-\frac{\partial J}{\partial \theta} \Big|_{\theta=\theta_0} = \frac{\partial^2 J}{\partial \theta \partial \theta^T} \Big|_{\theta=\theta_0} \Delta\theta \quad (82)$$

$$\Delta\theta = - \left[ \frac{\partial^2 J}{\partial \theta \partial \theta^T} \Big|_{\theta=\theta_0} \right]^{-1} \frac{\partial J}{\partial \theta} \Big|_{\theta=\theta_0} \quad (83)$$

which can be implemented as the iterative gradient descent algorithm:

$$\hat{\theta} = \theta_0 + \Delta\hat{\theta} \quad (84)$$

### References

- [1] Eugene A Morelli and Vladislav Klein. *Aircraft System Identification: Theory and Practice*.
- [2] Bruce P Gibbs. *Advanced Kalman Filtering, Least-squares and Modeling*. John Wiley & Sons, 2011.
- [3] Stephen Boyd, Stephen P Boyd, and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.