# Appendix 2: Least squares analysis

# 11

## CHAPTER OUTLINE HEAD

## 11.1 The least squares criterion

The *least squares criterion* is one of the foundations of *estimation theory*. This is the theory that concerns extracting the **true** value of signals from **noisy** measurements. Estimation theory techniques have been used to guide Exocet missiles and astronauts on moon missions (where navigation data was derived using sextants!), all based on techniques which employ the least squares criterion. The least squares criterion was originally developed by Gauss when he was confronted by the problem of measuring the six parameters of the orbits of planets, given astronomical measurements. These measurements were naturally subject to error, and Gauss realized that they could be combined together in some way in order to reduce a best estimate of the six parameters of interest.

Gauss assumed that the noise corrupting the measurements would have a *normal distribution*; indeed such distributions are often now called Gaussian to honor his great insight. As a consequence of the *central limit theorem*, it may be assumed that many real random noise sources are normally distributed. In cases where this assumption is not valid, the mathematical advantages that accrue from its use generally offset any resulting loss of accuracy. Also, the assumption of normality is particularly invaluable in view of the fact that the output of a system excited by Gaussian-distributed noise is also Gaussian-distributed (as seen in Fourier analysis, Chapter 2). A Gaussian probability distribution of a variable $x$ is defined by

$$p(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{\frac{-(x-\bar{x})^2}{\sigma^2}} \qquad (11.1)$$

where $\bar{x}$ is the mean (loosely the average) of the distribution and $\sigma^2$ is the second moment or variance of the distribution. Given many measurements of a single unknown quantity, when that quantity is subject to errors of a zero-mean (symmetric) normal distribution, it is well known that the best estimate of the

unknown quantity is the average of the measurements. In the case of two or more unknown quantities, the requirement is to combine the measurements in such a way that the error in the estimates of the unknown quantities is minimized. Clearly, direct averaging will not suffice when measurements are a function of two or more unknown quantities.

Consider the case where $N$ equally precise measurements, $f_1, f_2, \ldots, f_N$, are made on a linear function $f(a)$ of a single parameter $a$. The measurements are subject to zero-mean additive Gaussian noise $v_i(t)$, as such the measurements are given by

$$f_i = f(a) + v_i(t) \quad \forall i \in 1, N \tag{11.2}$$

The differences $\tilde{f}$ between the true value of the function and the noisy measurements of it are

$$\tilde{f}_i = f(a) - f_i \quad \forall i \in 1, N \tag{11.3}$$

By Eq. (11.1), the probability distribution of these errors is

$$p(\tilde{f}_i) = \frac{1}{\sigma\sqrt{2\pi}} e^{\frac{-(\tilde{f}_i)^2}{\sigma^2}} \quad \forall i \in 1, N \tag{11.4}$$

Since the errors are **independent**, the **compound** distribution of these errors is the product of their distributions and is given by

$$p(\tilde{f}) = \frac{1}{\sigma\sqrt{2\pi}} e^{\frac{-((\tilde{f}_1)^2 + (\tilde{f}_2)^2 + (\tilde{f}_3)^2 + \cdots + (\tilde{f}_N)^2)}{\sigma^2}} \tag{11.5}$$

Each of the errors is a function of the unknown quantity, $a$, which is to be estimated. Different estimates of $a$ will give different values for $p(\tilde{f})$. The most probable system of errors will be that for which $p(\tilde{f})$ is a **maximum** and this corresponds to the **best** estimate of the unknown quantity. Thus, to maximize $p(\tilde{f})$

$$\begin{aligned}
\max\{p(\tilde{f})\} &= \max\left\{ \frac{1}{\sigma\sqrt{2\pi}} e^{\frac{-((\tilde{f}_1)^2 + (\tilde{f}_2)^2 + (\tilde{f}_3)^2 + \cdots + (\tilde{f}_N)^2)}{\sigma^2}} \right\} \\
&= \max\left\{ e^{\frac{-((\tilde{f}_1)^2 + (\tilde{f}_2)^2 + (\tilde{f}_3)^2 + \cdots + (\tilde{f}_N)^2)}{\sigma^2}} \right\} \\
&= \max\{-((\tilde{f}_1)^2 + (\tilde{f}_2)^2 + (\tilde{f}_3)^2 + \cdots + (\tilde{f}_N)^2)\} \\
&= \min\{-((\tilde{f}_1)^2 + (\tilde{f}_2)^2 + (\tilde{f}_3)^2 + \cdots + (\tilde{f}_N)^2)\}
\end{aligned} \tag{11.6}$$

Thus, the required estimate is that which **minimizes** the sum of the differences squared, and this estimate is the one that is **optimal** by the least squares criterion.

This criterion leads on to the method of least squares which follows in the next section. This is a method commonly used to fit curves to measured data. This concerns estimating the values of parameters from a complete set of measurements.

There are also techniques that provide estimate of parameters at time instants, based on a set of previous measurements. These techniques include the Weiner filter and the Kalman filter. The Kalman filter was the algorithm chosen for guiding Exocet missiles and moon missions (an extended square root Kalman filter, no less).

## 11.2 Curve fitting by least squares

*Curve fitting* by the method of least squares concerns combining a set of measurements to derive **estimates** of the parameters which specify the curve that best fits the data. By the least squares criterion, given a set of $N$ (noisy) measurements $f_i, i \in 1, N$, which are to be fitted to a curve $f(\mathbf{a})$, where $\mathbf{a}$ is a vector of parameter values, we seek to minimize the square of the difference between the measurements and the values of the curve to give an estimate of the parameters $\hat{\mathbf{a}}$ according to

$$\hat{\mathbf{a}} = \min \sum_{i=1}^{N} (f_i - f(x_i, y_i, \mathbf{a}))^2 \tag{11.7}$$

Since we seek a minimum, by differentiation we obtain

$$\frac{\partial \sum_{i=1}^{N} (f_i - f(x_i, y_i, \mathbf{a}))^2}{\partial \mathbf{a}} = 0 \tag{11.8}$$

which implies that

$$2 \sum_{i=1}^{N} (f_i - f(x_i, y_i, \mathbf{a})) \frac{\partial f(\mathbf{a})}{\partial \mathbf{a}} = 0 \tag{11.9}$$

The solution is usually of the form

$$\mathbf{Ma} = \mathbf{F} \tag{11.10}$$

where $\mathbf{M}$ is a matrix of summations of products of the index $i$ and $\mathbf{F}$ is a vector of summations of products of the measurements and $i$. The solution, the best estimate of the values of $\mathbf{a}$, is then given by

$$\hat{\mathbf{a}} = \mathbf{M}^{-1} \mathbf{F} \tag{11.11}$$

For example, let us consider the problem of fitting a 2D surface to a set of data points. The surface is given by

$$f(x, y, \mathbf{a}) = a + bx + cy + dxy \tag{11.12}$$

where the vector of parameters $\mathbf{a} = [a\ b\ c\ d]^\mathrm{T}$ controls the shape of the surface and $(x,y)$ are the coordinates of a point on the surface. Given a set of (noisy) measurements of the value of the surface at points with coordinates $(x,y)$, $f_i = f(x,y) + v_i$,

we seek to estimate values for the parameters using the method of least squares. By Eq. (11.7), we seek

$$\hat{\mathbf{a}} = [\hat{a} \ \hat{b} \ \hat{c} \ \hat{d}]^{\mathbf{T}} = \min \sum_{i=1}^{N} (f_i - f(x_i, y_i, \mathbf{a}))^2 \qquad (11.13)$$

By Eq. (11.9), we require

$$2 \sum_{i=1}^{N} (f_i - (a + bx_i + cy_i + dx_iy_i)) \frac{\partial f(x_i, y_i, \mathbf{a})}{\partial \mathbf{a}} = 0 \qquad (11.14)$$

By differentiating $f(x, y, \mathbf{a})$ with respect to each parameter, we have

$$\frac{\partial f(x_i, y_i)}{\partial a} = 1 \qquad (11.15)$$

$$\frac{\partial f(x_i, y_i)}{\partial b} = x \qquad (11.16)$$

$$\frac{\partial f(x_i, y_i)}{\partial c} = y \qquad (11.17)$$

and

$$\frac{\partial f(x_i, y_i)}{\partial d} = xy \qquad (11.18)$$

and by substituting Eqs (11.15)−(11.18) in Eq. (11.14), we obtain four simultaneous equations:

$$\sum_{i=1}^{N} (f_i - (a + bx_i + cy_i + dx_iy_i)) \times 1 = 0 \qquad (11.19)$$

$$\sum_{i=1}^{N} (f_i - (a + bx_i + cy_i + dx_iy_i)) \times x_i = 0 \qquad (11.20)$$

$$\sum_{i=1}^{N} (f_i - (a + bx_i + cy_i + dx_iy_i)) \times y_i = 0 \qquad (11.21)$$

and

$$\sum_{i=1}^{N} (f_i - (a + bx_i + cy_i + dx_iy_i)) \times x_iy_i = 0 \qquad (11.22)$$

Since $\sum_{i=1}^{N} a = Na$, Eq. (11.19) can be reformulated as

$$\sum_{i=1}^{N} f_i - Na - b \sum_{i=1}^{N} x_i - c \sum_{i=1}^{N} y_i - d \sum_{i=1}^{N} x_iy_i = 0 \qquad (11.23)$$

and Eqs (11.20)−(11.22) can be reformulated likewise. By expressing the simultaneous equations in matrix form, we get

$$
\begin{bmatrix}
N & \sum\limits_{i=1}^{N} x_i & \sum\limits_{i=1}^{N} y_i & \sum\limits_{i=1}^{N} x_i y_i \\
\sum\limits_{i=1}^{N} x_i & \sum\limits_{i=1}^{N} (x_i)^2 & \sum\limits_{i=1}^{N} x_i y_i & \sum\limits_{i=1}^{N} (x_i)^2 y_i \\
\sum\limits_{i=1}^{N} y_i & \sum\limits_{i=1}^{N} x_i y_i & \sum\limits_{i=1}^{N} (y_i)^2 & \sum\limits_{i=1}^{N} x_i (y_i)^2 \\
\sum\limits_{i=1}^{N} x_i y_i & \sum\limits_{i=1}^{N} (x_i)^2 y_i & \sum\limits_{i=1}^{N} x_i (y_i)^2 & \sum\limits_{i=1}^{N} (x_i)^2 (y_i)^2
\end{bmatrix}
\begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix}
=
\begin{bmatrix}
\sum\limits_{i=1}^{N} f_i \\
\sum\limits_{i=1}^{N} f_i x_i \\
\sum\limits_{i=1}^{N} f_i y_i \\
\sum\limits_{i=1}^{N} f_i x_i y_i
\end{bmatrix}
\tag{11.24}
$$

and this is the same form as Eq. (11.10) and can be solved by inversion, as in Eq. (11.11). Note that the matrix is **symmetric** and its inversion, or solution, does not impose such a great computational penalty as appears. Given a set of data points, the values need to be entered in the summations, thus completing the matrices from which the solution is found. This technique can replace the one used in the zero-crossing detector within the Marr−Hildreth edge detection operator (Section 4.3.3) but appeared to offer no significant advantage over the (much simpler) function implemented there.