

### Least Squares Estimation

The general objective of least squares estimation is that you want to minimize some measure of the deviation between the model estimate  $f(x_i)$  and the known values  $f_0(x_i)$  - thus hoping that you will have a “good fit” of the model to the data. The closeness of the fit will depend on data errors and the selected function model.

We will show at a later time that if these deviations  $\{ f(x_i) - f_0(x_i) \}$  are Gaussian (normally distributed), the “best estimate” is found by minimizing the sum of squares of the deviations; i.e., we minimize

$$E = \sum_{i=1}^N [\tilde{f}(x_i) - f_0(x_i)]^2, \quad (1)$$

where  $N$  is the number of observations.

Our first example will be a one-dimensional case, in which the model is a 2<sup>nd</sup>-order polynomial. Thus the model is

$$\tilde{f}(x) = a_0 + a_1x + a_2x^2, \quad (2)$$

and therefore we minimize

$$E = \sum_{i=1}^N [(a_0 + a_1x_i + a_2x_i^2) - f_0(x_i)]^2. \quad (3)$$

Note that for this quadratic model,  $N$  must be equal to or greater than 3.

Since the  $x_i$  and  $f_0(x_i)$  are known, we need to determine  $a_0$ ,  $a_1$  and  $a_2$ .

To find  $a_0$ ,  $a_1$  and  $a_2$  such that  $E$  is a minimum, we derive the **normal equations** by setting

$$\frac{\partial E}{\partial a_0} \rightarrow 0, \quad \frac{\partial E}{\partial a_1} \rightarrow 0, \quad \frac{\partial E}{\partial a_2} \rightarrow 0.$$

Note that setting  $\frac{\partial A}{\partial Q} \rightarrow 0$  usually finds an extrema only, but since  $E$  is positive definite, this extrema is always a minima in this case.

Thus, performing these operations on eq. (3) yields

$$\frac{\partial E}{\partial a_0} = 2 \sum_{i=1}^N [(a_0 + a_1 x_i + a_2 x_i^2) - f_0(x_i)] = 0, \quad (4)$$

$$\frac{\partial E}{\partial a_1} = 2 \sum_{i=1}^N [(a_0 + a_1 x_i + a_2 x_i^2) - f_0(x_i)] x_i = 0, \quad (5)$$

$$\frac{\partial E}{\partial a_2} = 2 \sum_{i=1}^N [(a_0 + a_1 x_i + a_2 x_i^2) - f_0(x_i)] x_i^2 = 0. \quad (6)$$

Taking the summation inside eq. (4) (and dividing out the 2), we have

$$\sum_{i=1}^N a_0 + \sum_{i=1}^N a_1 x_i + \sum_{i=1}^N a_2 x_i^2 - \sum_{i=1}^N f_0(x_i) = 0$$

or

$$Na_0 + \sum_{i=1}^N a_1 x_i + \sum_{i=1}^N a_2 x_i^2 - \sum_{i=1}^N f_0(x_i) = 0$$

and, dividing by  $N$  and defining the ensemble average  $\overline{(\quad)} = \frac{1}{N} \sum_{i=1}^N (\quad)$

we have  $a_0 + a_1 \overline{x_i} + a_2 \overline{x_i^2} = \overline{f_0(x_i)}$ .

Similarly, eqs. (5) and (6) can be rewritten as

$$a_0 \overline{x_i} + a_1 \overline{x_i^2} + a_2 \overline{x_i^3} = \overline{x_i f_0(x_i)},$$

$$a_0 \overline{x_i^2} + a_1 \overline{x_i^3} + a_2 \overline{x_i^4} = \overline{x_i^2 f_0(x_i)}.$$

These normal equations can be written in matrix form:

$$\begin{bmatrix} 1 & \overline{x_i} & \overline{x_i^2} \\ \overline{x_i} & \overline{x_i^2} & \overline{x_i^3} \\ \overline{x_i^2} & \overline{x_i^3} & \overline{x_i^4} \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ a_2 \end{bmatrix} = \begin{bmatrix} \overline{f_0(x_i)} \\ \overline{x_i f_0(x_i)} \\ \overline{x_i^2 f_0(x_i)} \end{bmatrix}. \quad (7)$$

$$\begin{matrix} \overline{x_i^2} & \overline{x_i^3} & \overline{x_i^4} & & \\ & \mathbf{a}_2 & & & \\ & & & & \overline{x_i^2 f_0(x_i)} \end{matrix}$$

or  $\mathbf{X} \mathbf{A} = \mathbf{F}$  (8)

Note that the coefficient matrix  $\mathbf{X}$  is symmetric.

The solution of (8) is

$$\mathbf{A} = \mathbf{X}^{-1} \mathbf{F} \quad (9)$$

provided the inverse exists (i.e.  $\mathbf{X}$  is nonsingular.)

Note: Moments in  $\mathbf{X}$  increase greatly for higher-order polynomials;  $\mathbf{X}$  may be singular or ill-conditioned beyond 5<sup>th</sup> or 6<sup>th</sup> order.

Now we are ready to examine Sec. 2.1 of Daley, which is an example of two-dimensional local polynomial fitting. This technique, developed by Panofsky (1949), produced the first objective analysis ever used in NWP. Local fitting is still occasionally used in regional data analysis studies.

The problem is illustrated by the diagram on the right, for which we need to determine an analysis value  $f_A(x, y)$  at the center point  $r_i$  of the circle with an influence radius  $R_i$ . Only observations within the influence radius are used to determine  $f_A(x, y)$ ; - i.e.

$$x_k^2 + y_k^2 \leq R_i^2 .$$

Denote the observations  $f_o(x_k, y_k)$ ,  $k = 1 \text{ K}$ , where  $K$  is the number of observations in the influence region.

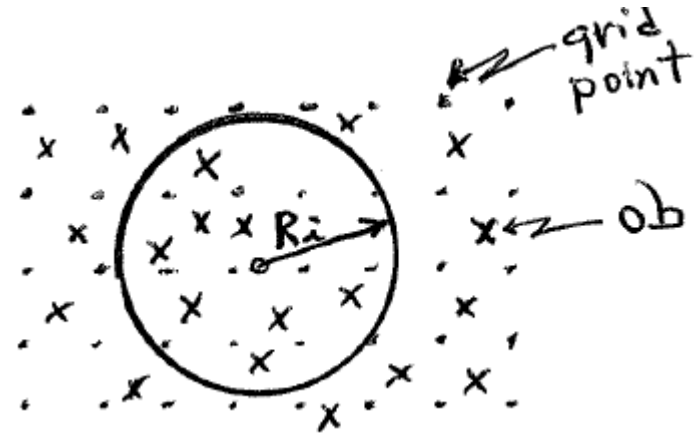
Define  $x = y = 0$  at analysis point  $r_i$  .

Assume  $f_A(x, y)$  can be represented by the following 2-D polynomial expansion:

$$f_A(x, y) = \sum_m \sum_n c_{m,n} x^m y^n \quad m, n \geq 0, \quad m + n \leq M . \quad (10)$$

Note that at the center point  $x = 0, y = 0$ , in eq. (10) yields  $f_A(0,0) = c_{00}$  .

At each  $r_i$  , we want to minimize the differences between the function model (10) and the observations; i.e., we minimize  $I$ , given by



$$I = \frac{1}{2} \sum_{k=1}^{K_i} \left[ \sum_m \sum_n c_{m,n} x_k^m y_k^n - f_0(x_k, y_k) \right]^2. \quad (11)$$

As usual, we minimize  $I$  by taking the first derivative of (11) w.r.t. the unknown coefficients  $c_{m,n}$  and setting the result = 0:

$$\frac{\partial I}{\partial c_{m,n}} I = \sum_{k=1}^{K_i} x_k^m y_k^n \left[ \sum_\mu \sum_\nu c_{\mu,\nu} x_k^\mu y_k^\nu - f_0(x_k, y_k) \right] \rightarrow 0,$$

which can be rewritten as (class exercise)

$$\sum_\mu \sum_\nu c_{\mu,\nu} \sum_{k=1}^{K_i} x_k^{m+\mu} y_k^{n+\nu} = \sum_{k=1}^{K_i} x_k^m y_k^n f_0(x_k, y_k). \quad (12)$$

Now, defining the ensemble average here as  $\overline{(\quad)} = \frac{1}{K_i} \sum_{k=1}^{K_i} (\quad)$ , also show that (12) can be written in matrix form given by eq. (2.1.4) in Daley when  $M = 2$ . (This permits the following coefficients:  $c_{00}$ ,  $c_{10}$ ,  $c_{20}$ ,  $c_{11}$ ,  $c_{01}$ ,  $c_{02}$ , and the matrices will be of order

$$(6 \times 6) (6 \times 1) = (6 \times 1) \quad (2.1.4 - Daley)$$

$$\begin{pmatrix}
 \overline{1} & \overline{x_k} & \overline{y_k} & \overline{x_k^2} & \overline{y_k^2} & \overline{x_k y_k} \\
 \overline{x_k} & \overline{x_k^2} & \overline{x_k y_k} & \overline{x_k^3} & \overline{x_k y_k^2} & \overline{x_k^2 y_k} \\
 \overline{y_k} & \overline{x_k y_k} & \overline{y_k^2} & \overline{x_k^2 y_k} & \overline{y_k^3} & \overline{x_k y_k^2} \\
 \overline{x_k^2} & \overline{x_k^3} & \overline{x_k^2 y_k} & \overline{x_k^4} & \overline{x_k^2 y_k^2} & \overline{x_k^3 y_k} \\
 \overline{y_k^2} & \overline{x_k y_k^2} & \overline{y_k^3} & \overline{x_k^2 y_k^2} & \overline{y_k^4} & \overline{x_k y_k^3} \\
 \overline{x_k y_k} & \overline{x_k^2 y_k} & \overline{x_k y_k^2} & \overline{x_k^3 y_k} & \overline{x_k y_k^3} & \overline{x_k^2 y_k^2}
 \end{pmatrix}
 \begin{pmatrix}
 c_{00} \\
 c_{10} \\
 c_{01} \\
 c_{20} \\
 c_{02} \\
 c_{11}
 \end{pmatrix}
 =
 \begin{pmatrix}
 \overline{f_0(x_k, y_k)} \\
 \overline{x_k f_0(x_k, y_k)} \\
 \overline{y_k f_0(x_k, y_k)} \\
 \overline{x_k^2 f_0(x_k, y_k)} \\
 \overline{y_k^2 f_0(x_k, y_k)} \\
 \overline{x_k y_k f_0(x_k, y_k)}
 \end{pmatrix}
 \quad (2.1.4)$$

where the elements of the matrix on the left-hand side and of the column vector on the right-hand side are all known.

Equation (2.1.4) is then solved for the 6 coefficients to obtain one analysis value. We then move to the next grid point and repeat the process; i.e. – local fitting.

### **Least squares estimation derivation**

Since least squares minimization is used or implied throughout this course, we will examine it more closely. Here we will show that the least squares minimization value,  $S_a$ , is a maximum likelihood estimate of  $S$ .

Consider  $N$  observations  $s_1, s_2, s_3, \dots, s_N$  of a variable  $s$  - at one point.

Assume each observation is taken with a different instrument and the observational error of each ob. is  $\epsilon_n = s_n - s$ ,  $n = 1, N$ , where  $s$  is the (unknown) true value.

Assume errors are random, unbiased, and normally distributed, and recall the formula for the normal (Gaussian) probability density function  $p_G$  given by

$$p_G = \frac{1}{\sigma_s \sqrt{2\pi}} \exp \left[ -\frac{(s - \langle s \rangle)^2}{2\sigma_s^2} \right], \quad (13)$$

where  $\sigma_s$  is the standard deviation of  $s$  and



$\sigma_s^2$  is the variance (2<sup>nd</sup> moment) of  $s$ , denoted as  $\langle (s - \langle s \rangle)^2 \rangle$

where  $\langle s \rangle$  is the expected value or mean of  $s$  given by  $\langle s \rangle = \int_{-\infty}^{\infty} s \cdot p(s) ds$ ,

i.e., the sum over all values of  $s$ , each one weighted by its probability  $p$ .

$p(s)$  is any probability density function (pdf) that describes the distribution of  $s$ . A pdf gives the probability that  $s$  lies between  $s$  and  $s + ds$ , and satisfies the properties

$$p(s) \geq 0; \quad \int_{-\infty}^{\infty} p(s) ds = 1.$$

[For a review of statistics, see Appendix D in Daley.]

So, if the errors are Gaussian, we use eq. (13) and the definition of  $\varepsilon_n$  to write that the probability that the error in the  $n^{\text{th}}$  observation lies between  $\varepsilon_n$  and  $\varepsilon_n + d\varepsilon_n$  is given by

$$p(\varepsilon_n) = \frac{1}{\sigma_n \sqrt{2\pi}} \exp\left[-\frac{\varepsilon_n^2}{2\sigma_n^2}\right], \quad (14)$$

where

$$\sigma_n^2 = \langle (s_n - s)^2 \rangle = \langle \varepsilon_n^2 \rangle = \int_{-\infty}^{\infty} \varepsilon_n^2 \cdot p(\varepsilon_n) d\varepsilon_n . \quad (15)$$

Note that in (14) we have used the fact that  $\langle \varepsilon_n \rangle = 0$ , since errors are assumed unbiased.

Now, for  $N = 1$ , we have one observation  $s_1$  with error  $\varepsilon_1$ . Thus the most probable value of  $s$  is that value for which  $p(\varepsilon_1)$  is a maximum.

Since  $\langle \varepsilon_n \rangle = 0$ , this occurs at  $\varepsilon_1 = 0$ , or, via the left side of (15), when  $s_1 = s$ . (Note that these results come from the properties of a Gaussian distribution.)

When there are  $N > 1$  obs., we need to compute the joint probability that

$\varepsilon_1$  lies between  $\varepsilon_1 + d\varepsilon_1$ ,  $\varepsilon_2$  lies between  $\varepsilon_2 + d\varepsilon_2$ , .....,  $\varepsilon_N$  lies between  $\varepsilon_N + d\varepsilon_N$ , which is the product of all the probabilities – which we write as

$p(\varepsilon_1, \varepsilon_2, \dots, \varepsilon_N) = p(\varepsilon_1) p(\varepsilon_2) \dots p(\varepsilon_N)$ , which, via eq. (14) is

$$p(\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n) = \prod_{n=1}^N \frac{1}{\sigma_n \sqrt{2\pi}} \exp \left[ -\frac{\varepsilon_n^2}{2\sigma_n^2} \right],$$

which can also be written

$$p(\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n) = \prod_{n=1}^N \left[ \frac{1}{\sigma_n \sqrt{2\pi}} \right] \exp \left[ -\sum_{n=1}^N \frac{(s_n - s)^2}{2\sigma_n^2} \right], \quad (16)$$

since exponents are additive.

In this formulation, the most probable value of  $s$  ( $s_a$ ) is that value for which the probability  $p(\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n)$  is a maximum. (That is, if we know the most likely  $\varepsilon_n$ , that that will direct us to the most likely  $s_n$ .)

This will occur when the quantity  $\sum_{n=1}^N \frac{(s_n - s)^2}{2\sigma_n^2}$  in (16) is a minimum (so that the exponent is large).

Thus the most probable value,  $s_a$ , often called the maximum likelihood estimate of  $s$ , must minimize

$$I = \frac{1}{2} \sum_{n=1}^N \frac{(s_a - s_n)^2}{\sigma_n^2} \quad (17)$$

As before, to minimize a quantity, we differentiate  $I$  w.r.t. the unknown quantity  $s_a$  and set to zero:

$$\sum_{n=1}^N \sigma_n^{-2} (s_a - s_n) = \frac{s_a - s_1}{\sigma_1^2} + \frac{s_a - s_2}{\sigma_2^2} + \dots + \frac{s_a - s_N}{\sigma_N^2} = 0$$

or

$$s_a \sum_{n=1}^N \sigma_n^{-2} - \sum_{n=1}^N \sigma_n^{-2} s_n = 0$$

and, solving for  $s_a$ , we have

$$s_a = \frac{\sum_{n=1}^N \sigma_n^{-2} s_n}{\sum_{n=1}^N \sigma_n^{-2}} . \quad (18)$$

Therefore, the most probable value, or most likelihood estimate is given by a weighted sum of the observations, with weights inversely proportional to the observational error variance.

Once we have an analysis estimate  $s_a$ , it is important to examine the error  $\varepsilon_a$  of the analysis - we will usually want to know the expected error variance of the analysis.

Thus we define  $\varepsilon_a = s_a - s$  (19)

And, using eqs. (18) and (19), we can write the expected error variance of the estimate as

$$\langle \varepsilon_a^2 \rangle = \left\langle \left[ \frac{\sum_{n=1}^N \sigma_n^{-2} (s_n - s)}{\sum_{n=1}^N \sigma_n^{-2}} \right]^2 \right\rangle. \quad (20)$$

**Class exercise:** Show that (20) can be simplified to

$$\langle \varepsilon_a^2 \rangle = \left[ \sum_{n=1}^N \sigma_n^{-2} \right]^{-1}. \quad (21)$$

(Hint: Do  $N = 2$  case first before do general case.)

Now suppose all the obs. were taken with the same instrument. Thus,  $\sigma_n^2 = \sigma^2$  for all  $n$ , and eqs. (18) and (21) reduce to

$$s_a = \frac{1}{N} \sum_{n=1}^N s_n \quad \text{and} \quad \langle \varepsilon_a^2 \rangle = \sigma^2 / N. \quad (22a,b)$$

Thus  $s_a$  (22a) becomes a simple ensemble average. It is important to note from (22b) that when  $N > 1$ ,  $\langle \varepsilon_a^2 \rangle < \sigma^2$  which means that the analysis value is better than any one observation.