

## 17.1 LINEAR REGRESSION

457

line connecting the points. However, any straight line passing through the midpoint of the connecting line (except a perfectly vertical line) results in a minimum value of Eq. (17.2) equal to zero because the errors cancel.

Therefore, another logical criterion might be to minimize the sum of the absolute values of the discrepancies, as in

$$\sum_{i=1}^n |e_i| = \sum_{i=1}^n |y_i - a_0 - a_1 x_i|$$

Figure 17.2*b* demonstrates why this criterion is also inadequate. For the four points shown, any straight line falling within the dashed lines will minimize the sum of the absolute values. Thus, this criterion also does not yield a unique best fit.

A third strategy for fitting a best line is the *minimax* criterion. In this technique, the line is chosen that minimizes the maximum distance that an individual point falls from the line. As depicted in Fig. 17.2*c*, this strategy is ill-suited for regression because it gives undue influence to an outlier, that is, a single point with a large error. It should be noted that the minimax principle is sometimes well-suited for fitting a simple function to a complicated function (Carnahan, Luther, and Wilkes, 1969).

A strategy that overcomes the shortcomings of the aforementioned approaches is to minimize the sum of the squares of the residuals between the measured  $y$  and the  $y$  calculated with the linear model

$$S_r = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_{i,\text{measured}} - y_{i,\text{model}})^2 = \sum_{i=1}^n (y_i - a_0 - a_1 x_i)^2 \quad (17.3)$$

This criterion has a number of advantages, including the fact that it yields a unique line for a given set of data. Before discussing these properties, we will present a technique for determining the values of  $a_0$  and  $a_1$  that minimize Eq. (17.3).

### 17.1.2 Least-Squares Fit of a Straight Line

To determine values for  $a_0$  and  $a_1$ , Eq. (17.3) is differentiated with respect to each coefficient:

$$\begin{aligned} \frac{\partial S_r}{\partial a_0} &= -2 \sum (y_i - a_0 - a_1 x_i) \\ \frac{\partial S_r}{\partial a_1} &= -2 \sum [(y_i - a_0 - a_1 x_i) x_i] \end{aligned}$$

Note that we have simplified the summation symbols; unless otherwise indicated, all summations are from  $i = 1$  to  $n$ . Setting these derivatives equal to zero will result in a minimum  $S_r$ . If this is done, the equations can be expressed as

$$\begin{aligned} 0 &= \sum y_i - \sum a_0 - \sum a_1 x_i \\ 0 &= \sum y_i x_i - \sum a_0 x_i - \sum a_1 x_i^2 \end{aligned}$$

Now, realizing that  $\Sigma a_0 = na_0$ , we can express the equations as a set of two simultaneous linear equations with two unknowns ( $a_0$  and  $a_1$ ):

$$na_0 + \left(\sum x_i\right)a_1 = \sum y_i \quad (17.4)$$

$$\left(\sum x_i\right)a_0 + \left(\sum x_i^2\right)a_1 = \sum x_i y_i \quad (17.5)$$

These are called the *normal equations*. They can be solved simultaneously

$$a_1 = \frac{n \Sigma x_i y_i - \Sigma x_i \Sigma y_i}{n \Sigma x_i^2 - (\Sigma x_i)^2} \quad (17.6)$$

This result can then be used in conjunction with Eq. (17.4) to solve for

$$a_0 = \bar{y} - a_1 \bar{x} \quad (17.7)$$

where  $\bar{y}$  and  $\bar{x}$  are the means of  $y$  and  $x$ , respectively.

**EXAMPLE 17.1****Linear Regression**

**Problem Statement.** Fit a straight line to the  $x$  and  $y$  values in the first two columns of Table 17.1.

**Solution.** The following quantities can be computed:

$$n = 7 \quad \sum x_i y_i = 119.5 \quad \sum x_i^2 = 140$$

$$\sum x_i = 28 \quad \bar{x} = \frac{28}{7} = 4$$

$$\sum y_i = 24 \quad \bar{y} = \frac{24}{7} = 3.428571$$

Using Eqs. (17.6) and (17.7),

$$a_1 = \frac{7(119.5) - 28(24)}{7(140) - (28)^2} = 0.8392857$$

$$a_0 = 3.428571 - 0.8392857(4) = 0.07142857$$

**TABLE 17.1** Computations for an error analysis of the linear fit.

$x_i$	$y_i$	$(y_i - \bar{y})^2$	$(y_i - a_0 - a_1 x_i)^2$
1	0.5	8.5765	0.1687
2	2.5	0.8622	0.5625
3	2.0	2.0408	0.3473
4	4.0	0.3265	0.3265
5	3.5	0.0051	0.5896
6	6.0	6.6122	0.7972
7	5.5	4.2908	0.1993
$\Sigma$	24.0	22.7143	2.9911