

Regression

Where substantial error is associated with data, polynomial interpolation is inappropriate and may yield unsatisfactory results when used to predict intermediate values.

Experimental data are often of this type. For example, Fig. 17.1a shows seven experimentally derived data points exhibiting significant variability. Visual inspection of these data suggests a positive relationship between y and x . That is, the overall trend indicates that higher values of y are associated with higher values of x . Now, if a sixth-order interpolating polynomial is fitted to these data (Fig. 17.1b), it will pass exactly through all of the points. However, because of the variability in these data, the curve oscillates widely in the interval between the points. In particular, the interpolated values at $x = 1.5$ and $x = 6.5$ appear to be well beyond the range suggested by these data.

A more appropriate strategy for such cases is to derive an approximating function that fits the shape or general trend of the data without necessarily matching the individual points. Figure 17.1c illustrates how a straight line can be used to generally characterize the trend of these data without passing through any particular point.

One way to determine the line in Fig. 17.1c is to visually inspect the plotted data and then sketch a “best” line through the points. Although such “eyeball” approaches have commonsense appeal and are valid for “back-of-the-envelope” calculations, they are deficient because they are arbitrary. That is, unless the points define a perfect straight line (in which case, interpolation would be appropriate), different analysts would draw different lines.

To remove this subjectivity, some criterion must be devised to establish a basis for the fit. One way to do this is to derive a curve that minimizes the discrepancy between the data points and the curve. A technique for accomplishing this objective, called *least-squares regression*

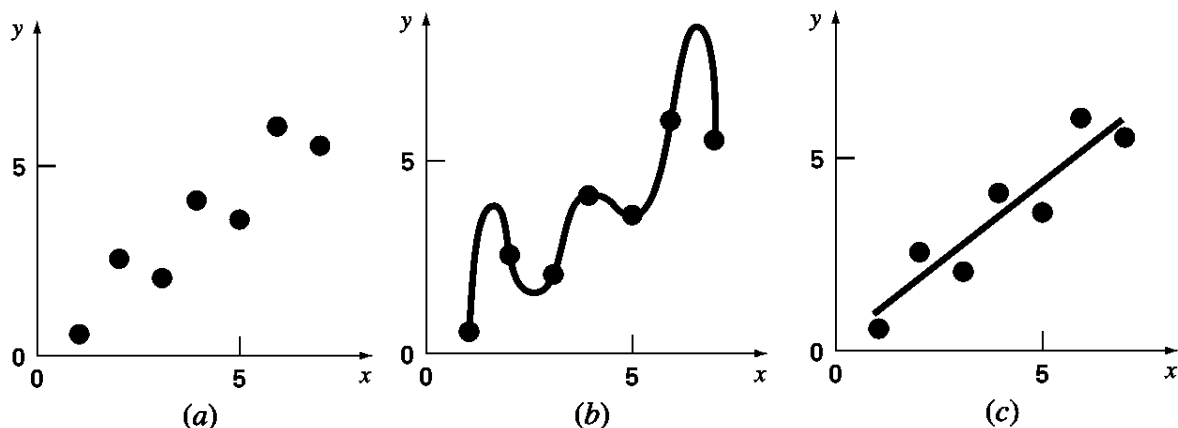


FIGURE 17.1

(a) Data exhibiting significant error. (b) Polynomial fit oscillating beyond the range of the data. (c) More satisfactory result using the least-squares fit.

LINEAR REGRESSION

The simplest example of a least-squares approximation is fitting a straight line to a set of paired observations: $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$. The mathematical expression for the straight line is

$$y = a_0 + a_1x + e \quad (17.1)$$

where a_0 and a_1 are coefficients representing the intercept and the slope, respectively, and e is the error, or residual, between the model and the observations, which can be represented by rearranging Eq. (17.1) as

$$e = y - a_0 - a_1x$$

Thus, the error, or *residual*, is the discrepancy between the true value of y and the approximate value, $a_0 + a_1x$, predicted by the linear equation.

One strategy for fitting a “best” line through the data would be to minimize the sum of the residual errors for all the available data, as in

$$\sum_{i=1}^n e_i = \sum_{i=1}^n (y_i - a_0 - a_1x_i) \quad (17.2)$$

where n = total number of points. However, this is an inadequate criterion, as illustrated by Fig. 17.2a which depicts the fit of a straight line to two points. Obviously, the best fit is the 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

$$\sum_{i=1}^n |e_i| = \sum_{i=1}^n |y_i - a_0 - a_1x_i|$$

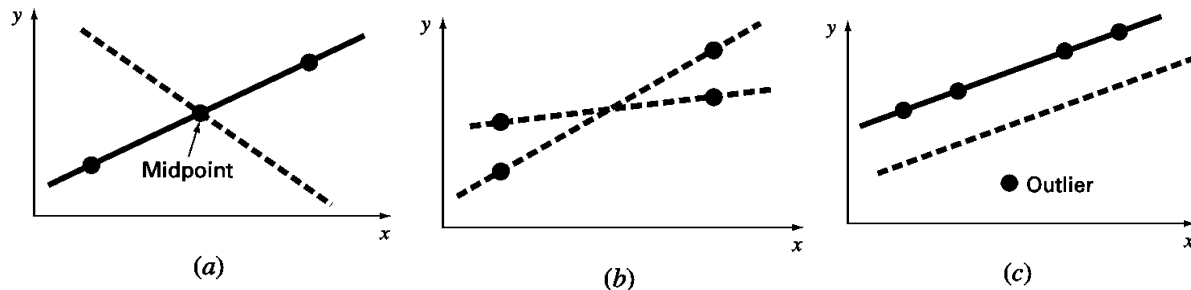


FIGURE 17.2

Examples of some criteria for “best fit” that are inadequate for regression: (a) minimizes the sum of the residuals, (b) minimizes the sum of the absolute values of the residuals, and (c) minimizes the maximum error of any individual point.

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).

Least-Squares Fit of a Straight Line

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

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.

17.1.3 Quantification of Error of Linear Regression

Any line other than the one computed in Example 17.1 results in a larger sum of the squares of the residuals. Thus, the line is unique and in terms of our chosen criterion is a “best” line through the points. A number of additional properties of this fit can be elucidated by examining more closely the way in which residuals were computed. Recall that the sum of the squares is defined as [Eq. (17.3)]

$$S_r = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - a_0 - a_1 x_i)^2 \quad (17.8)$$

Notice the similarity between Eqs. (PT5.3) and (17.8). In the former case, the square of the residual represented the square of the discrepancy between the data and a single estimate of the measure of central tendency—the mean. In Eq. (17.8), the square of the residual represents the square of the vertical distance between the data and another measure of central tendency—the straight line (Fig. 17.3).

The analogy can be extended further for cases where (1) the spread of the points around the line is of similar magnitude along the entire range of the data and (2) the distribution of these points about the line is normal. It can be demonstrated that if these criteria are met, least-squares regression will provide the best (that is, the most likely) estimates of a_0 and a_1 (Draper and Smith, 1981). This is called the *maximum likelihood principle* in statistics. In addition, if these criteria are met, a “standard deviation” for the regression line can be determined as [compare with Eq. (PT5.2)]

$$s_{y/x} = \sqrt{\frac{S_r}{n - 2}} \quad (17.9)$$

where $s_{y/x}$ is called the *standard error of the estimate*. The subscript notation “y/x” designates that the error is for a predicted value of y corresponding to a particular value of x. Also, notice that we now divide by $n - 2$ because two data-derived estimates— a_0 and a_1 —were used to compute S_r ; thus, we have lost two degrees of freedom. As with our discussion of the standard deviation in PT5.2.1, another justification for dividing by $n - 2$ is that there is no such thing as the “spread of data” around a straight line connecting two points. Thus, for the case where $n = 2$, Eq. (17.9) yields a meaningless result of infinity.

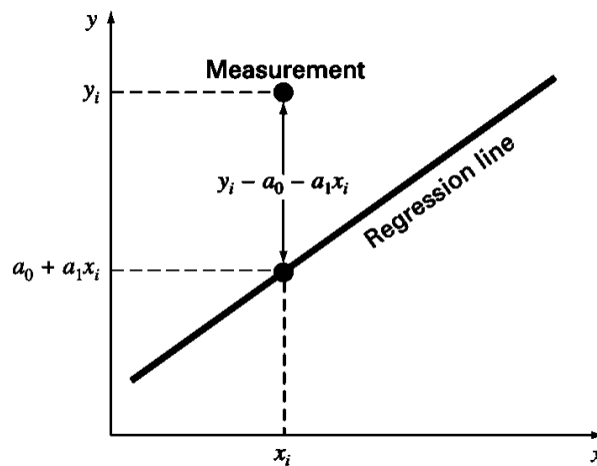


FIGURE 17.3

The residual in linear regression represents the vertical distance between a data point and the straight line

Just as was the case with the standard deviation, the standard error of the estimate quantifies the spread of the data. However, $s_{y/x}$ quantifies the spread *around the regression line* as shown in Fig. 17.4b in contrast to the original standard deviation s_y that quantified the spread *around the mean* (Fig. 17.4a).

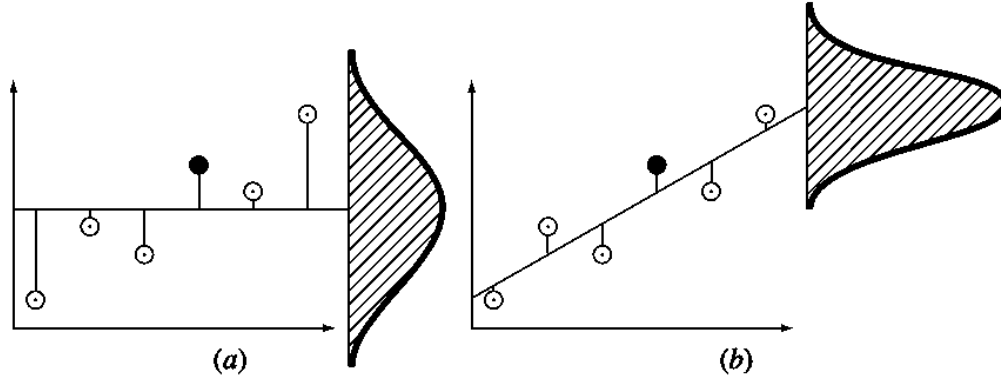


FIGURE 17.4

Regression data showing (a) the spread of the data around the mean of the dependent variable and (b) the spread of the data around the best-fit line. The reduction in the spread in going from (a) to (b), as indicated by the bell-shaped curves at the right, represents the improvement due to linear regression.

The above concepts can be used to quantify the “goodness” of our fit. This is particularly useful for comparison of several regressions (Fig. 17.5). To do this, we return to the original data and determine the *total sum of the squares* around the mean for the dependent variable (in our case, y). As was the case for Eq. (PT5.3), this quantity is designated S_t . This is the magnitude of the residual error associated with the dependent variable prior to regression. After performing the regression, we can compute S_r , the sum of the squares of the residuals around the regression line. This characterizes the residual error that remains after the regression. It is, therefore, sometimes called the unexplained sum of the squares. The difference between the two quantities, $S_t - S_r$, quantifies the improvement or error reduction due to describing the data in terms of a straight line rather than as an average value. Because the magnitude of this quantity is scale-dependent, the difference is normalized to S_t to yield

$$r^2 = \frac{S_t - S_r}{S_t} \quad (17.10)$$

where r^2 is called the *coefficient of determination* and r is the *correlation coefficient* ($=\sqrt{r^2}$). For a perfect fit, $S_r = 0$ and $r = r^2 = 1$, signifying that the line explains 100 percent of the variability of the data. For $r = r^2 = 0$, $S_r = S_t$ and the fit represents no improvement. An alternative formulation for r that is more convenient for computer implementation is

$$r = \frac{n\sum x_i y_i - (\sum x_i)(\sum y_i)}{\sqrt{n\sum x_i^2 - (\sum x_i)^2} \sqrt{n\sum y_i^2 - (\sum y_i)^2}} \quad (17.11)$$

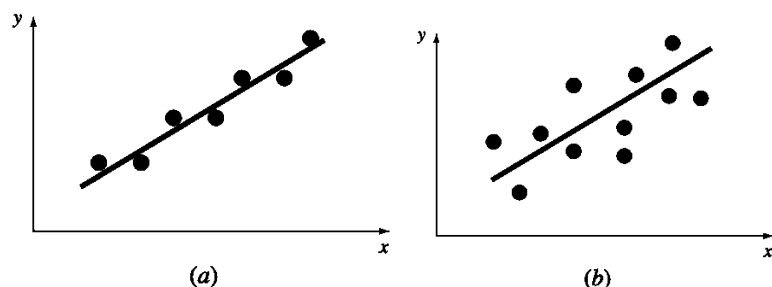


FIGURE 17.5

Examples of linear regression with (a) small and (b) large residual errors.

Example 4.1 Find the best values of a_0 and a_1 if the straight line $Y = a_0 + a_1x$ is fitted to the data (x_i, y_i) :

$(1, 0.6), (2, 2.4), (3, 3.5), (4, 4.8), (5, 5.7)$

Find also the correlation coefficient.

From the table of values given below, we find $\bar{x} = 3$, $\bar{y} = 3.4$, and

$$a_1 = \frac{5(63.6) - 15(17)}{5(55) - 225} = 1.26$$

Therefore,

$$a_0 = \bar{y} - a_1\bar{x} = -0.38.$$

x_i	y_i	x_i^2	$x_i y_i$	$(y_i - \bar{y})^2$	$(y_i - a_0 - a_1 x_i)^2$
1	0.6	1	0.6	7.84	0.0784
2	2.4	4	4.8	1.00	0.0676
3	3.5	9	10.5	0.01	0.0100
4	4.8	16	19.2	1.96	0.0196
5	5.7	25	28.5	5.29	0.0484
15	17.0	55	63.6	16.10	0.2240

$$\text{The correlation coefficient} = \sqrt{\frac{16.10 - 0.2240}{16.10}} = 0.9930.$$

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	$(y_i - \bar{y})$	$(y_i - a_0 - a_1 x_i)^2$
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
Σ	24.0	22.7143	2.9911

Therefore, the least-squares fit is

$$y = 0.07142857 + 0.8392857x$$

Problem Statement. Compute the total standard deviation, the standard error of the estimate, and the correlation coefficient for the data in Example 17.1.

Solution. The summations are performed and presented in Table 17.1. The standard deviation is [Eq. (PT5.2)]

$$s_y = \sqrt{\frac{22.7143}{7 - 1}} = 1.9457$$

and the standard error of the estimate is [Eq. (17.9)]

$$s_{y/x} = \sqrt{\frac{2.9911}{7 - 2}} = 0.7735$$

Thus, because $s_{y/x} < s_y$, the linear regression model has merit. The extent of the improvement is quantified by [Eq. (17.10)]

$$r^2 = \frac{22.7143 - 2.9911}{22.7143} = 0.868$$

or

$$r = \sqrt{0.868} = 0.932$$

These results indicate that 86.8 percent of the original uncertainty has been explained by the linear model.