

## 6.6 EXERCISES

In Exercises 1–4, find the equation  $y = \beta_0 + \beta_1 x$  of the least-squares line that best fits the given data points.

- (0, 1), (1, 1), (2, 2), (3, 2)
- (1, 0), (2, 1), (4, 2), (5, 3)
- (−1, 0), (0, 1), (1, 2), (2, 4)
- (2, 3), (3, 2), (5, 1), (6, 0)
- Let  $X$  be the design matrix used to find the least-squares line to fit data  $(x_1, y_1), \dots, (x_n, y_n)$ . Use a theorem in Section 6.5 to show that the normal equations have a unique solution if and only if the data include at least two data points with different  $x$ -coordinates.
- Let  $X$  be the design matrix in Example 2 corresponding to a least-squares fit of a parabola to data  $(x_1, y_1), \dots, (x_n, y_n)$ . Suppose  $x_1, x_2$ , and  $x_3$  are distinct. Explain why there is only one parabola that fits the data best, in a least-squares sense. (See Exercise 5.)
- A certain experiment produces the data (1, 1.8), (2, 2.7), (3, 3.4), (4, 3.8), (5, 3.9). Describe the model that produces a least-squares fit of these points by a function of the form

$$y = \beta_1 x + \beta_2 x^2$$

Such a function might arise, for example, as the revenue from the sale of  $x$  units of a product, when the amount offered for sale affects the price to be set for the product.

- Give the design matrix, the observation vector, and the unknown parameter vector.
  - [M] Find the associated least-squares curve for the data.
- A simple curve that often makes a good model for the variable costs of a company, as a function of the sales level  $x$ , has the form  $y = \beta_1 x + \beta_2 x^2 + \beta_3 x^3$ . There is no constant term because fixed costs are not included.
    - Give the design matrix and the parameter vector for the linear model that leads to a least-squares fit of the equation above, with data  $(x_1, y_1), \dots, (x_n, y_n)$ .
    - [M] Find the least-squares curve of the form above to fit the data (4, 1.58), (6, 2.08), (8, 2.5), (10, 2.8), (12, 3.1), (14, 3.4), (16, 3.8), and (18, 4.32), with values in thousands. If possible, produce a graph that shows the data points and the graph of the cubic approximation.
  - A certain experiment produces the data (1, 7.9), (2, 5.4), and (3, −.9). Describe the model that produces a least-squares fit of these points by a function of the form

$$y = A \cos x + B \sin x$$

- Suppose radioactive substances A and B have decay constants of .02 and .07, respectively. If a mixture of these two substances at time  $t = 0$  contains  $M_A$  grams of A and  $M_B$  grams of B, then a model for the total amount  $y$  of the mixture present at time  $t$  is

$$y = M_A e^{-.02t} + M_B e^{-.07t} \tag{6}$$

Suppose the initial amounts  $M_A$  and  $M_B$  are unknown, but a scientist is able to measure the total amounts present at several times and records the following points  $(t_i, y_i)$ : (10, 21.34), (11, 20.68), (12, 20.05), (14, 18.87), and (15, 18.30).

- Describe a linear model that can be used to estimate  $M_A$  and  $M_B$ .
- [M] Find the least-squares curve based on (6).



Halley's Comet last appeared in 1986 and will reappear in 2061.

- [M] According to Kepler's first law, a comet should have an elliptic, parabolic, or hyperbolic orbit (with gravitational attractions from the planets ignored). In suitable polar coordinates, the position  $(r, \vartheta)$  of a comet satisfies an equation of the form

$$r = \beta + e(r \cdot \cos \vartheta)$$

where  $\beta$  is a constant and  $e$  is the *eccentricity* of the orbit, with  $0 \leq e < 1$  for an ellipse,  $e = 1$  for a parabola, and  $e > 1$  for a hyperbola. Suppose observations of a newly discovered comet provide the data below. Determine the type of orbit, and predict where the comet will be when  $\vartheta = 4.6$  (radians).<sup>3</sup>

$\vartheta$	.88	1.10	1.42	1.77	2.14
$r$	3.00	2.30	1.65	1.25	1.01

- [M] A healthy child's systolic blood pressure  $p$  (in millimeters of mercury) and weight  $w$  (in pounds) are approximately related by the equation

$$\beta_0 + \beta_1 \ln w = p$$

Use the following experimental data to estimate the systolic blood pressure of a healthy child weighing 100 pounds.

<sup>3</sup> The basic idea of least-squares fitting of data is due to K. F. Gauss (and, independently, to A. Legendre), whose initial rise to fame occurred in 1801 when he used the method to determine the path of the asteroid *Ceres*. Forty days after the asteroid was discovered, it disappeared behind the sun. Gauss predicted it would appear ten months later and gave its location. The accuracy of the prediction astonished the European scientific community.

$w$	44	61	81	113	131
$\ln w$	3.78	4.11	4.39	4.73	4.88
$p$	91	98	103	110	112

13. [M] To measure the takeoff performance of an airplane, the horizontal position of the plane was measured every second, from  $t = 0$  to  $t = 12$ . The positions (in feet) were: 0, 8.8, 29.9, 62.0, 104.7, 159.1, 222.0, 294.5, 380.4, 471.1, 571.7, 686.8, and 809.2.
- Find the least-squares cubic curve  $y = \beta_0 + \beta_1 t + \beta_2 t^2 + \beta_3 t^3$  for these data.
  - Use the result of part (a) to estimate the velocity of the plane when  $t = 4.5$  seconds.
14. Let  $\bar{x} = \frac{1}{n}(x_1 + \cdots + x_n)$  and  $\bar{y} = \frac{1}{n}(y_1 + \cdots + y_n)$ . Show that the least-squares line for the data  $(x_1, y_1), \dots, (x_n, y_n)$  must pass through  $(\bar{x}, \bar{y})$ . That is, show that  $\bar{x}$  and  $\bar{y}$  satisfy the linear equation  $\bar{y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{x}$ . [Hint: Derive this equation from the vector equation  $\mathbf{y} = X\hat{\boldsymbol{\beta}} + \boldsymbol{\epsilon}$ . Denote the first column of  $X$  by  $\mathbf{1}$ . Use the fact that the residual vector  $\boldsymbol{\epsilon}$  is orthogonal to the column space of  $X$  and hence is orthogonal to  $\mathbf{1}$ .]

Given data for a least-squares problem,  $(x_1, y_1), \dots, (x_n, y_n)$ , the following abbreviations are helpful:

$$\sum x = \sum_{i=1}^n x_i, \quad \sum x^2 = \sum_{i=1}^n x_i^2,$$

$$\sum y = \sum_{i=1}^n y_i, \quad \sum xy = \sum_{i=1}^n x_i y_i$$

The normal equations for a least-squares line  $y = \hat{\beta}_0 + \hat{\beta}_1 x$  may be written in the form

$$\begin{aligned} n\hat{\beta}_0 + \hat{\beta}_1 \sum x &= \sum y \\ \hat{\beta}_0 \sum x + \hat{\beta}_1 \sum x^2 &= \sum xy \end{aligned} \quad (7)$$

15. Derive the normal equations (7) from the matrix form given in this section.
16. Use a matrix inverse to solve the system of equations in (7) and thereby obtain formulas for  $\hat{\beta}_0$  and  $\hat{\beta}_1$  that appear in many statistics texts.

17. a. Rewrite the data in Example 1 with new  $x$ -coordinates in mean deviation form. Let  $X$  be the associated design matrix. Why are the columns of  $X$  orthogonal?  
b. Write the normal equations for the data in part (a), and solve them to find the least-squares line,  $y = \beta_0 + \beta_1 x^*$ , where  $x^* = x - 5.5$ .
18. Suppose the  $x$ -coordinates of the data  $(x_1, y_1), \dots, (x_n, y_n)$  are in mean deviation form, so that  $\sum x_i = 0$ . Show that if  $X$  is the design matrix for the least-squares line in this case, then  $X^T X$  is a diagonal matrix.

Exercises 19 and 20 involve a design matrix  $X$  with two or more columns and a least-squares solution  $\hat{\boldsymbol{\beta}}$  of  $\mathbf{y} = X\boldsymbol{\beta}$ . Consider the following numbers.

- $\|X\hat{\boldsymbol{\beta}}\|^2$ —the sum of the squares of the “regression term.” Denote this number by SS(R).
- $\|\mathbf{y} - X\hat{\boldsymbol{\beta}}\|^2$ —the sum of the squares for error term. Denote this number by SS(E).
- $\|\mathbf{y}\|^2$ —the “total” sum of the squares of the  $y$ -values. Denote this number by SS(T).

Every statistics text that discusses regression and the linear model  $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon}$  introduces these numbers, though terminology and notation vary somewhat. To simplify matters, assume that the mean of the  $y$ -values is zero. In this case, SS(T) is proportional to what is called the *variance* of the set of  $y$ -values.

19. Justify the equation  $\text{SS(T)} = \text{SS(R)} + \text{SS(E)}$ . [Hint: Use a theorem, and explain why the hypotheses of the theorem are satisfied.] This equation is extremely important in statistics, both in regression theory and in the analysis of variance.
20. Show that  $\|X\hat{\boldsymbol{\beta}}\|^2 = \hat{\boldsymbol{\beta}}^T X^T \mathbf{y}$ . [Hint: Rewrite the left side and use the fact that  $\hat{\boldsymbol{\beta}}$  satisfies the normal equations.] This formula for SS(R) is used in statistics. From this and from Exercise 19, obtain the standard formula for SS(E):  
$$\text{SS(E)} = \mathbf{y}^T \mathbf{y} - \hat{\boldsymbol{\beta}}^T X^T \mathbf{y}$$

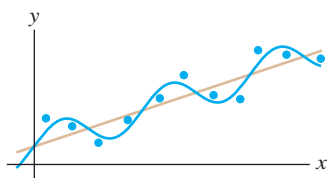
### SOLUTION TO PRACTICE PROBLEM

Construct  $X$  and  $\boldsymbol{\beta}$  so that the  $k$ th row of  $X\boldsymbol{\beta}$  is the predicted  $y$ -value that corresponds to the data point  $(x_k, y_k)$ , namely,

$$\beta_0 + \beta_1 x_k + \beta_2 \sin(2\pi x_k/12)$$

It should be clear that

$$X = \begin{bmatrix} 1 & x_1 & \sin(2\pi x_1/12) \\ \vdots & \vdots & \vdots \\ 1 & x_n & \sin(2\pi x_n/12) \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix}$$



Sales trend with seasonal fluctuations.