

- (c) Solve the systems using A and A' and determine the actual relative error.
- (d) Suppose \mathbf{b} is changed to $\mathbf{b}' = \begin{bmatrix} 100 \\ 101 \end{bmatrix}$. How large a relative change can this change produce in the solution to $A\mathbf{x} = \mathbf{b}$? [Hint: Use Exercise 42.]
- (e) Solve the systems using \mathbf{b} and \mathbf{b}' and determine the actual relative error.

44. Let $A = \begin{bmatrix} 1 & 1 & 1 \\ 2 & 5 & 0 \\ 1 & -1 & 2 \end{bmatrix}$ and $\mathbf{b} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$.

- (a) Compute $\text{cond}_1(A)$.
- (b) Suppose A is changed to $A' = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 5 & 0 \\ 1 & -1 & 2 \end{bmatrix}$. How large a relative change can this change produce in the solution to $A\mathbf{x} = \mathbf{b}$? [Hint: Use inequality (1) from this section.]
- (c) Solve the systems using A and A' and determine the actual relative error.
- (d) Suppose \mathbf{b} is changed to $\mathbf{b}' = \begin{bmatrix} 1 \\ 1 \\ 3 \end{bmatrix}$. How large a relative change can this change produce in the solution to $A\mathbf{x} = \mathbf{b}$? [Hint: Use Exercise 42.]
- (e) Solve the systems using \mathbf{b} and \mathbf{b}' and determine the actual relative error.

45. Show that if A is an invertible matrix, then $\text{cond}(A) \geq 1$ with respect to any matrix norm.

46. Show that if A and B are invertible matrices, then $\text{cond}(AB) \leq \text{cond}(A)\text{cond}(B)$ with respect to any matrix norm.
47. Let A be an invertible matrix and let λ_1 and λ_n be the eigenvalues with the largest and smallest absolute values, respectively. Show that

$$\text{cond}(A) \geq \frac{|\lambda_1|}{|\lambda_n|}$$

[Hint: See Exercise 34 and Theorem 4.18(b) in Section 4.3.]

CAS In Exercises 48–51, write the given system in the form of Equation (7). Then use the method of Example 7.22 to estimate the number of iterations of Jacobi's method that will be needed to approximate the solution to three-decimal-place accuracy. (Use $\mathbf{x}_0 = \mathbf{0}$.) Compare your answer with the solution computed in the given exercise from Section 2.5.

48. Exercise 1, Section 2.5 49. Exercise 3, Section 2.5
50. Exercise 4, Section 2.5 51. Exercise 5, Section 2.5

Exercise 52(c) refers to the Leontief model of an open economy, as discussed in Sections 2.4 and 3.7.

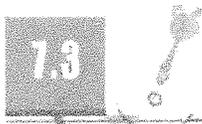
52. Let A be an $n \times n$ matrix such that $\|A\| < 1$, where the norm is either the sum norm or the max norm.

- (a) Prove that $A^n \rightarrow O$ as $n \rightarrow \infty$.
(b) Deduce from (a) that $I - A$ is invertible and

$$(I - A)^{-1} = I + A + A^2 + A^3 + \cdots$$

[Hint: See the proof of Theorem 3.34.]

- (c) Show that (b) can be used to prove Corollaries 3.35 and 3.36.



7.3 Least Squares Approximation

In many branches of science, experimental data are used to infer a mathematical relationship among the variables being measured. For example, we might measure the height of a tree at various points in time and try to deduce a function that expresses the tree's height h in terms of time t . Or, we might measure the size p of a population over time and try to find a rule that relates p to t . Relationships between variables are also of interest in business; for example, a company producing widgets may be interested in knowing the relationship between its total costs c and the number n of widgets produced.

In each of these examples, the data come in the form of two measurements: one for the independent variable and one for the (supposedly) dependent variable. Thus, we have a set of *data points* (x_i, y_i) , and we are looking for a function that best approximates the relationship between the independent variable x and the dependent variable y . Figure 7.9 shows examples in which experimental data points are plotted, along with a curve that approximately “fits” the data.

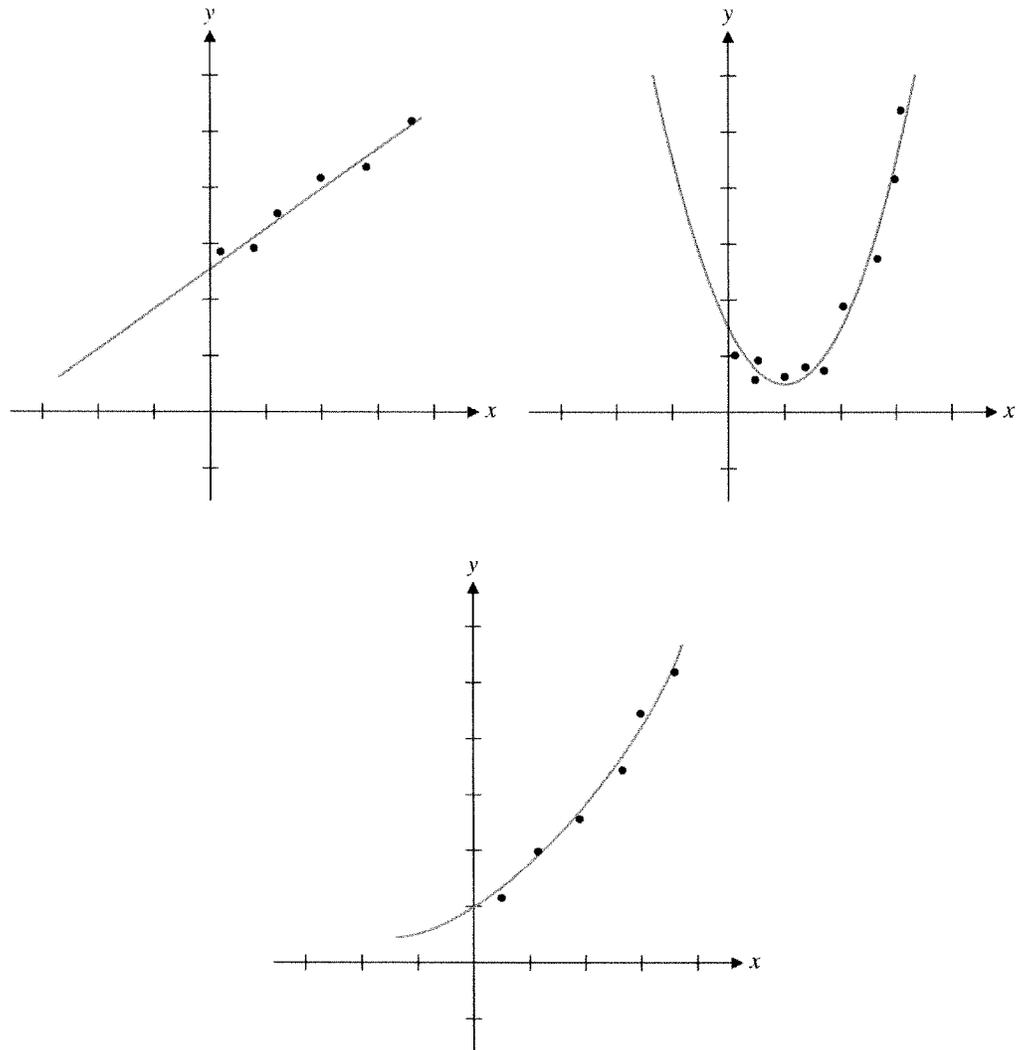


Figure 7.9
Curves of “best fit”

Roger Cotes (1682–1716) was an English mathematician who, while a fellow at Cambridge, edited the second edition of Newton’s *Principia*. Although he published little, he made important discoveries in the theory of logarithms, integral calculus, and numerical methods.

The method of least squares, which we are about to consider, is attributed to Gauss. A new asteroid, Ceres, was discovered on New Year’s Day, 1801, but it disappeared behind the sun shortly after it was observed. Astronomers predicted when and where Ceres would reappear, but their calculations differed greatly from those done, independently, by Gauss. Ceres reappeared on December 7, 1801, almost exactly where Gauss had predicted it would be. Although he did not disclose his methods at the time, Gauss had used his least squares approximation method, which he described in a paper in 1809. The same method was actually known earlier; Cotes anticipated the method in the early 18th century, and Legendre published a paper on it in 1806. Nevertheless, Gauss is generally given credit for the method of least squares approximation.

We begin our exploration of approximation with a more general result.

The Best Approximation Theorem

In the sciences, there are many problems that can be phrased generally as “What is the best approximation to X of type Y ?” X might be a set of data points, a function, a vector, or many other things, while Y might be a particular type of function, a vector belonging to a certain vector space, etc. A typical example of such a problem is finding the vector \mathbf{w} in a subspace W of a vector space V that best approximates (i.e., is closest to) a given vector \mathbf{v} in V . This problem gives rise to the following definition.

Definition If W is a subspace of a normed linear space V and if \mathbf{v} is a vector in V , then the **best approximation to \mathbf{v} in W** is the vector $\bar{\mathbf{v}}$ in W such that

$$\|\mathbf{v} - \bar{\mathbf{v}}\| < \|\mathbf{v} - \mathbf{w}\|$$

for every vector \mathbf{w} in W different from $\bar{\mathbf{v}}$.

In \mathbb{R}^2 or \mathbb{R}^3 , we are used to thinking of “shortest distance” as corresponding to “perpendicular distance.” In algebraic terminology, “shortest distance” relates to the notion of orthogonal projection: If W is a subspace of \mathbb{R}^n and \mathbf{v} is a vector in \mathbb{R}^n , then we expect $\text{proj}_W(\mathbf{v})$ to be the vector in W that is closest to \mathbf{v} (Figure 7.10).

Since orthogonal projection can be defined in any inner product space, we have the following theorem.

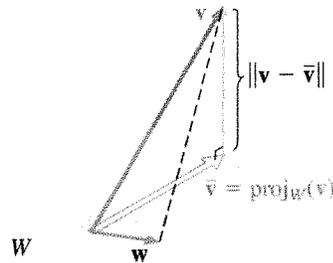


Figure 7.10

If $\bar{\mathbf{v}} = \text{proj}_W(\mathbf{v})$, then
 $\|\mathbf{v} - \bar{\mathbf{v}}\| < \|\mathbf{v} - \mathbf{w}\|$ for all $\mathbf{w} \neq \bar{\mathbf{v}}$

Theorem 7.8 The Best Approximation Theorem

If W is a finite-dimensional subspace of an inner product space V and if \mathbf{v} is a vector in V , then $\text{proj}_W(\mathbf{v})$ is the best approximation to \mathbf{v} in W .

Proof Let \mathbf{w} be a vector in W different from $\text{proj}_W(\mathbf{v})$. Then $\text{proj}_W(\mathbf{v}) - \mathbf{w}$ is also in W , so $\mathbf{v} - \text{proj}_W(\mathbf{v}) = \text{perp}_W(\mathbf{v})$ is orthogonal to $\text{proj}_W(\mathbf{v}) - \mathbf{w}$, by Exercise 43 in Section 7.1. Pythagoras' Theorem now implies that

$$\begin{aligned} \|\mathbf{v} - \text{proj}_W(\mathbf{v})\|^2 + \|\text{proj}_W(\mathbf{v}) - \mathbf{w}\|^2 &= \|(\mathbf{v} - \text{proj}_W(\mathbf{v})) + (\text{proj}_W(\mathbf{v}) - \mathbf{w})\|^2 \\ &= \|\mathbf{v} - \mathbf{w}\|^2 \end{aligned}$$

as Figure 7.10 illustrates. However, $\|\text{proj}_W(\mathbf{v}) - \mathbf{w}\|^2 > 0$, since $\mathbf{w} \neq \text{proj}_W(\mathbf{v})$, so

$$\|\mathbf{v} - \text{proj}_W(\mathbf{v})\|^2 < \|\mathbf{v} - \text{proj}_W(\mathbf{v})\|^2 + \|\text{proj}_W(\mathbf{v}) - \mathbf{w}\|^2 = \|\mathbf{v} - \mathbf{w}\|^2$$

or, equivalently,

$$\|\mathbf{v} - \text{proj}_W(\mathbf{v})\| < \|\mathbf{v} - \mathbf{w}\|$$

Example 7.23

Let $\mathbf{u}_1 = \begin{bmatrix} 1 \\ 2 \\ -1 \end{bmatrix}$, $\mathbf{u}_2 = \begin{bmatrix} 5 \\ -2 \\ 1 \end{bmatrix}$, and $\mathbf{v} = \begin{bmatrix} 3 \\ 2 \\ 5 \end{bmatrix}$. Find the best approximation to \mathbf{v} in the plane $W = \text{span}(\mathbf{u}_1, \mathbf{u}_2)$ and find the Euclidean distance from \mathbf{v} to W .

Solution The vector in W that best approximates \mathbf{v} is $\text{proj}_W(\mathbf{v})$. Since \mathbf{u}_1 and \mathbf{u}_2 are orthogonal,

$$\begin{aligned} \text{proj}_W(\mathbf{v}) &= \left(\frac{\mathbf{u}_1 \cdot \mathbf{v}}{\mathbf{u}_1 \cdot \mathbf{u}_1} \right) \mathbf{u}_1 + \left(\frac{\mathbf{u}_2 \cdot \mathbf{v}}{\mathbf{u}_2 \cdot \mathbf{u}_2} \right) \mathbf{u}_2 \\ &= \frac{2}{6} \begin{bmatrix} 1 \\ 2 \\ -1 \end{bmatrix} + \frac{16}{30} \begin{bmatrix} 5 \\ -2 \\ 1 \end{bmatrix} = \begin{bmatrix} \frac{3}{5} \\ -\frac{2}{5} \\ \frac{1}{5} \end{bmatrix} \end{aligned}$$

The distance from \mathbf{v} to W is the distance from \mathbf{v} to the point in W closest to \mathbf{v} . But this distance is just $\|\text{perp}_W(\mathbf{v})\| = \|\mathbf{v} - \text{proj}_W(\mathbf{v})\|$. We compute

$$\mathbf{v} - \text{proj}_W(\mathbf{v}) = \begin{bmatrix} 3 \\ 2 \\ 5 \end{bmatrix} - \begin{bmatrix} \frac{3}{5} \\ -\frac{2}{5} \\ \frac{1}{5} \end{bmatrix} = \begin{bmatrix} \frac{12}{5} \\ \frac{24}{5} \\ \frac{24}{5} \end{bmatrix}$$

so $\|\mathbf{v} - \text{proj}_W(\mathbf{v})\| = \sqrt{0^2 + \left(\frac{12}{5}\right)^2 + \left(\frac{24}{5}\right)^2} = \sqrt{\frac{720}{25}} = 12\sqrt{5}/5$

which is the distance from \mathbf{v} to W .

In Section 7.5, we will look at other examples of the Best Approximation Theorem when we explore the problem of approximating functions.

Remark The orthogonal projection of a vector \mathbf{v} onto a subspace W is defined in terms of an orthogonal basis for W . The Best Approximation Theorem gives us an alternative proof that $\text{proj}_W(\mathbf{v})$ does not depend on the choice of this basis, since there can be only one vector in W that is closest to \mathbf{v} —namely, $\text{proj}_W(\mathbf{v})$.

Least Squares Approximation

We now turn to the problem of finding a curve that “best fits” a set of data points. Before we can proceed, however, we need to define what we mean by “best fit.” Suppose the data points (1, 2), (2, 2), and (3, 4) have arisen from measurements taken during some experiment. Also suppose we have reason to believe that the x and y values are related by a linear function; that is, we expect the points to lie on some line with equation $y = a + bx$. If our measurements were accurate, all three points would satisfy this equation and we would have

$$2 = a + b \cdot 1 \quad 2 = a + b \cdot 2 \quad 4 = a + b \cdot 3$$

This is a system of three linear equations in two variables:

$$\begin{array}{l} a + b = 2 \\ a + 2b = 2 \\ a + 3b = 4 \end{array} \quad \text{or} \quad \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} 2 \\ 2 \\ 4 \end{bmatrix}$$

Unfortunately, this system is inconsistent (since the three points do not lie on a straight line). So we will settle for a line that comes “as close as possible” to passing through our points. For any line, we will measure the vertical distance from each data point to the line (representing the *errors* in the y -direction), and then we will try to choose the line that minimizes the *total error*. Figure 7.11 illustrates the situation.

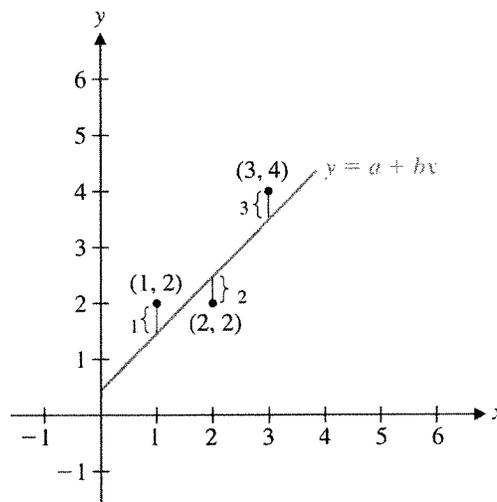


Figure 7.11

Finding the line that minimizes $\varepsilon_1^2 + \varepsilon_2^2 + \varepsilon_3^2$

If the errors are denoted by ε_1 , ε_2 , and ε_3 , then we can form the *error vector*

$$\mathbf{e} = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \varepsilon_3 \end{bmatrix}$$

We want \mathbf{e} to be as small as possible, so $\|\mathbf{e}\|$ must be as close to zero as possible. Which norm should we use? It turns out that the familiar Euclidean norm is the best choice. (The sum norm would also be a sensible choice, since $\|\mathbf{e}\|_s = |\varepsilon_1| + |\varepsilon_2| + |\varepsilon_3|$ is the actual sum of the errors in Figure 7.11. However, the absolute value signs are hard to work with, and, as you will soon see, the choice of the Euclidean norm leads to some very nice formulas.) So we are going to minimize

$$\|\mathbf{e}\| = \sqrt{\varepsilon_1^2 + \varepsilon_2^2 + \varepsilon_3^2} \quad \text{or, equivalently,} \quad \|\mathbf{e}\|^2 = \varepsilon_1^2 + \varepsilon_2^2 + \varepsilon_3^2$$

This is where the term “least squares” comes from: We need to find the smallest sum of squares, in the sense of the foregoing equation. The number $\|\mathbf{e}\|$ is called the *least squares error* of the approximation.

From Figure 7.11, we also obtain the following formulas for ϵ_1 , ϵ_2 , and ϵ_3 in our example:

$$\epsilon_1 = 2 - (a + b \cdot 1) \quad \epsilon_2 = 2 - (a + b \cdot 2) \quad \epsilon_3 = 4 - (a + b \cdot 3)$$

Example 7.24 Which of the following lines gives the smallest least squares error for the data points (1, 2), (2, 2), and (3, 4)?

(a) $y = 1 + x$
 (b) $y = -2 + 2x$
 (c) $y = \frac{2}{3} + x$

Solution Table 7.1 shows the necessary calculations.

Table 7.1

	$y = 1 + x$	$y = -2 + 2x$	$y = \frac{2}{3} + x$
ϵ_1	$2 - (1 + 1) = 0$	$2 - (-2 + 2) = 2$	$2 - (\frac{2}{3} + 1) = \frac{1}{3}$
ϵ_2	$2 - (1 + 2) = -1$	$2 - (-2 + 4) = 0$	$2 - (\frac{2}{3} + 2) = -\frac{2}{3}$
ϵ_3	$4 - (1 + 3) = 0$	$4 - (-2 + 6) = 0$	$4 - (\frac{2}{3} + 3) = \frac{1}{3}$
$\epsilon_1^2 + \epsilon_2^2 + \epsilon_3^2$	$0^2 + (-1)^2 + 0^2 = 1$	$2^2 + 0^2 + 0^2 = 4$	$(\frac{1}{3})^2 + (-\frac{2}{3})^2 + (\frac{1}{3})^2 = \frac{2}{3}$
$\ e\ $	1	2	$\sqrt{\frac{2}{3}} \approx 0.816$

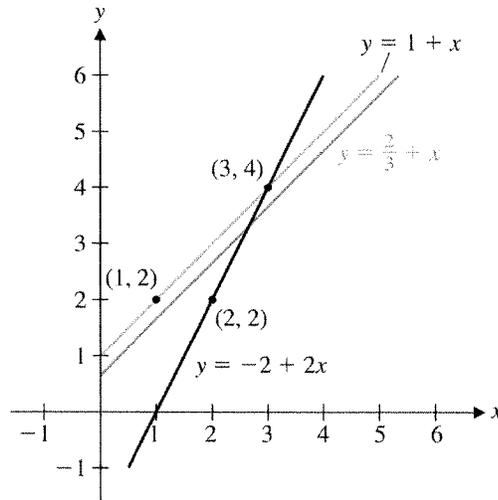


Figure 7.12

We see that the line $y = \frac{2}{3} + x$ produces the smallest least squares error among these three lines. Figure 7.12 shows the data points and all three lines.



It turns out that the line $y = \frac{2}{3} + x$ in Example 7.24 gives the smallest least squares error of *any* line, even though it passes through *none* of the given points. The rest of this section is devoted to illustrating why this is so.

In general, suppose we have n data points $(x_1, y_1), \dots, (x_n, y_n)$ and a line $y = a + bx$. Our error vector is

$$\mathbf{e} = \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

where $\varepsilon_i = y_i - (a + bx_i)$. The line $y = a + bx$ that minimizes $\varepsilon_1^2 + \dots + \varepsilon_n^2$ is called the **least squares approximating line** (or the **line of best fit**) for the points $(x_1, y_1), \dots, (x_n, y_n)$. As noted prior to Example 7.24, we can express this problem in matrix form. If the given points were actually on the line $y = a + bx$, then the n linear equations

$$\begin{aligned} a + bx_1 &= y_1 \\ &\vdots \\ a + bx_n &= y_n \end{aligned}$$

would all be true (i.e., the system would be consistent). Our interest is in the case where the points are *not* collinear, in which case the system is *inconsistent*. In matrix form, we have

$$\begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

which is of the form $A\mathbf{x} = \mathbf{b}$, where

$$A = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} a \\ b \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$



The error vector \mathbf{e} is just $\mathbf{b} - A\mathbf{x}$ (check this), and we want to minimize $\|\mathbf{e}\|^2$ or, equivalently, $\|\mathbf{e}\|$. We can therefore rephrase our problem in terms of matrices as follows.

Definition If A is an $m \times n$ matrix and \mathbf{b} is in \mathbb{R}^m , a **least squares solution** of $A\mathbf{x} = \mathbf{b}$ is a vector $\bar{\mathbf{x}}$ in \mathbb{R}^n such that

$$\|\mathbf{b} - A\bar{\mathbf{x}}\| \leq \|\mathbf{b} - A\mathbf{x}\|$$

for all \mathbf{x} in \mathbb{R}^n .

Solution of the Least Squares Problem

Any vector of the form $A\mathbf{x}$ is in the column space of A , and as \mathbf{x} varies over all vectors in \mathbb{R}^n , $A\mathbf{x}$ varies over all vectors in $\text{col}(A)$. A least squares solution of $A\mathbf{x} = \mathbf{b}$ is therefore equivalent to a vector $\bar{\mathbf{y}}$ in $\text{col}(A)$ such that

$$\|\mathbf{b} - \bar{\mathbf{y}}\| \leq \|\mathbf{b} - \mathbf{y}\|$$

for all \mathbf{y} in $\text{col}(A)$. In other words, we need the closest vector in $\text{col}(A)$ to \mathbf{b} . By the Best Approximation Theorem, the vector we want is the orthogonal projection of \mathbf{b} onto $\text{col}(A)$. Thus, if $\bar{\mathbf{x}}$ is a least squares solution of $A\mathbf{x} = \mathbf{b}$, we have

$$A\bar{\mathbf{x}} = \text{proj}_{\text{col}(A)}(\mathbf{b}) \quad (1)$$

In order to find $\bar{\mathbf{x}}$, it would appear that we need to first compute $\text{proj}_{\text{col}(A)}(\mathbf{b})$ and then solve the system (1). However, there is a better way to proceed.

We know that

$$\mathbf{b} - A\bar{\mathbf{x}} = \mathbf{b} - \text{proj}_{\text{col}(A)}(\mathbf{b}) = \text{perp}_{\text{col}(A)}(\mathbf{b})$$

is orthogonal to $\text{col}(A)$. So $\mathbf{b} - A\bar{\mathbf{x}}$ is in $(\text{col}(A))^\perp = \text{null}(A^T)$. Therefore $A^T(\mathbf{b} - A\bar{\mathbf{x}}) = \mathbf{0}$, which, in turn, is equivalent to $A^T\mathbf{b} - A^TA\bar{\mathbf{x}} = \mathbf{0}$ or

$$A^TA\bar{\mathbf{x}} = A^T\mathbf{b}$$

This represents a system of equations known as the *normal equations* for $\bar{\mathbf{x}}$.

We have just established that the solutions of the normal equations for $\bar{\mathbf{x}}$ are precisely the least squares solutions of $A\mathbf{x} = \mathbf{b}$. This proves the first part of the following theorem.

Theorem 7.9 The Least Squares Theorem

Let A be an $m \times n$ matrix and let \mathbf{b} be in \mathbb{R}^m . Then $A\mathbf{x} = \mathbf{b}$ always has at least one least squares solution $\bar{\mathbf{x}}$. Moreover:

- $\bar{\mathbf{x}}$ is a least squares solution of $A\mathbf{x} = \mathbf{b}$ if and only if $\bar{\mathbf{x}}$ is a solution of the normal equations $A^TA\bar{\mathbf{x}} = A^T\mathbf{b}$.
- A has linearly independent columns if and only if A^TA is invertible. In this case, the least squares solution of $A\mathbf{x} = \mathbf{b}$ is unique and is given by

$$\bar{\mathbf{x}} = (A^TA)^{-1}A^T\mathbf{b}$$

Proof We have already established property (a). For property (b), we note that the n columns of A are linearly independent if and only if $\text{rank}(A) = n$. But this is true if and only if A^TA is invertible, by Theorem 3.28. If A^TA is invertible, then the unique solution of $A^TA\bar{\mathbf{x}} = A^T\mathbf{b}$ is clearly $\bar{\mathbf{x}} = (A^TA)^{-1}A^T\mathbf{b}$.

Example 7.30

Use the QR factorization to find a least squares solution of $A\mathbf{x} = \mathbf{b}$, where

$$A = \begin{bmatrix} 1 & 2 & 2 \\ -1 & 1 & 2 \\ -1 & 0 & 1 \\ 1 & 1 & 2 \end{bmatrix} \quad \text{and} \quad \mathbf{b} = \begin{bmatrix} 2 \\ -3 \\ -2 \\ 0 \end{bmatrix}$$

Solution From Example 5.15,

$$A = QR = \begin{bmatrix} 1/2 & 3\sqrt{5}/10 & -\sqrt{6}/6 \\ -1/2 & 3\sqrt{5}/10 & 0 \\ -1/2 & \sqrt{5}/10 & \sqrt{6}/6 \\ 1/2 & \sqrt{5}/10 & \sqrt{6}/3 \end{bmatrix} \begin{bmatrix} 2 & 1 & 1/2 \\ 0 & \sqrt{5} & 3\sqrt{5}/2 \\ 0 & 0 & \sqrt{6}/2 \end{bmatrix}$$

We have

$$Q^T \mathbf{b} = \begin{bmatrix} 1/2 & -1/2 & -1/2 & 1/2 \\ 3\sqrt{5}/10 & 3\sqrt{5}/10 & \sqrt{5}/10 & \sqrt{5}/10 \\ -\sqrt{6}/6 & 0 & \sqrt{6}/6 & \sqrt{6}/3 \end{bmatrix} \begin{bmatrix} 2 \\ -3 \\ -2 \\ 0 \end{bmatrix} = \begin{bmatrix} 7/2 \\ -\sqrt{5}/2 \\ -2\sqrt{6}/3 \\ 0 \end{bmatrix}$$

so we require the solution to $R\bar{\mathbf{x}} = Q^T \mathbf{b}$, or

$$\begin{bmatrix} 2 & 1 & 1/2 \\ 0 & \sqrt{5} & 3\sqrt{5}/2 \\ 0 & 0 & \sqrt{6}/2 \end{bmatrix} \bar{\mathbf{x}} = \begin{bmatrix} 7/2 \\ -\sqrt{5}/2 \\ -2\sqrt{6}/3 \end{bmatrix}$$

Back substitution quickly yields

$$\bar{\mathbf{x}} = \begin{bmatrix} 4/3 \\ 3/2 \\ -4/3 \end{bmatrix}$$

Orthogonal Projection Revisited

One of the nice byproducts of the least squares method is a new formula for the orthogonal projection of a vector onto a subspace of \mathbb{R}^m .

Theorem 7.11

Let W be a subspace of \mathbb{R}^m and let A be an $m \times n$ matrix whose columns form a basis for W . If \mathbf{v} is any vector in \mathbb{R}^m , then the orthogonal projection of \mathbf{v} onto W is the vector

$$\text{proj}_W(\mathbf{v}) = A(A^T A)^{-1} A^T \mathbf{v}$$

The linear transformation $P: \mathbb{R}^m \rightarrow \mathbb{R}^m$ that projects \mathbb{R}^m onto W has $A(A^T A)^{-1} A^T$ as its standard matrix.

Proof Given the way we have constructed A , its column space is W . Since the columns of A are linearly independent, the Least Squares Theorem guarantees that there is a unique least squares solution to $A\mathbf{x} = \mathbf{v}$ given by

$$\bar{\mathbf{x}} = (A^T A)^{-1} A^T \mathbf{v}$$

By Equation (1),

$$A\bar{\mathbf{x}} = \text{proj}_{\text{col}(A)}(\mathbf{v}) = \text{proj}_W(\mathbf{v})$$

Therefore, $\text{proj}_W(\mathbf{v}) = A((A^T A)^{-1} A^T \mathbf{v}) = (A(A^T A)^{-1} A^T) \mathbf{v}$

as required. Since this equation holds for all \mathbf{v} in \mathbb{R}^m , the last statement of the theorem follows immediately. ▬

We will illustrate Theorem 7.11 by revisiting Example 5.11.

Example 7.31

Find the orthogonal projection of $\mathbf{v} = \begin{bmatrix} 3 \\ -1 \\ 2 \end{bmatrix}$ onto the plane W in \mathbb{R}^3 with equation $x - y + 2z = 0$, and give the standard matrix of the orthogonal projection transformation onto W .

Solution As in Example 5.11, we will take as a basis for W the set

$$\left\{ \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} -1 \\ 1 \\ 1 \end{bmatrix} \right\}$$

We form the matrix

$$A = \begin{bmatrix} 1 & -1 \\ 1 & 1 \\ 0 & 1 \end{bmatrix}$$

with these basis vectors as its columns. Then

$$A^T A = \begin{bmatrix} 1 & 1 & 0 \\ -1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & -1 \\ 1 & 1 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 2 & 0 \\ 0 & 3 \end{bmatrix}$$

so $(A^T A)^{-1} = \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & \frac{1}{3} \end{bmatrix}$

By Theorem 7.11, the standard matrix of the orthogonal projection transformation onto W is

$$A(A^T A)^{-1} A^T = A \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & \frac{1}{3} \end{bmatrix} \begin{bmatrix} 1 & 1 & 0 \\ -1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} \frac{5}{6} & \frac{1}{6} & -\frac{1}{3} \\ \frac{1}{6} & \frac{5}{6} & \frac{1}{3} \\ -\frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{bmatrix}$$

so the orthogonal projection of \mathbf{v} onto W is

$$\text{proj}_W(\mathbf{v}) = A(A^T A)^{-1} A^T \mathbf{v} = \begin{bmatrix} \frac{5}{6} & \frac{1}{6} & -\frac{1}{3} \\ \frac{1}{6} & \frac{5}{6} & \frac{1}{3} \\ -\frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{bmatrix} \begin{bmatrix} 3 \\ -1 \\ 2 \end{bmatrix} = \begin{bmatrix} \frac{5}{3} \\ \frac{1}{3} \\ -\frac{2}{3} \end{bmatrix}$$

which agrees with our solution to Example 5.11. ↖

Remark Since the projection of a vector onto a subspace W is unique, the standard matrix of this linear transformation (as given by Theorem 7.11) cannot depend on the choice of basis for W . In other words, with a different basis for W , we have a different matrix A , but the matrix $A(A^T A)^{-1} A^T$ will be the same! (You are asked to verify this in Exercise 43.)

The Pseudoinverse of a Matrix

If A is an $n \times n$ matrix with linearly independent columns, then it is invertible, and the unique solution to $A\mathbf{x} = \mathbf{b}$ is $\mathbf{x} = A^{-1}\mathbf{b}$. If $m > n$ and A is $m \times n$ with linearly independent columns, then $A\mathbf{x} = \mathbf{b}$ has no exact solution, but the best approximation is given by the unique least squares solution $\bar{\mathbf{x}} = (A^T A)^{-1} A^T \mathbf{b}$. The matrix $(A^T A)^{-1} A^T$ therefore plays the role of an “inverse of A ” in this situation.

Definition If A is a matrix with linearly independent columns, then the **pseudoinverse** of A is the matrix A^+ defined by

$$A^+ = (A^T A)^{-1} A^T$$

Observe that if A is $m \times n$, then A^+ is $n \times m$.

Example 7.32

Find the pseudoinverse of $A = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \end{bmatrix}$.

Solution We have already done most of the calculations in Example 7.26. Using our previous work, we have

$$A^+ = (A^T A)^{-1} A^T = \begin{bmatrix} \frac{7}{3} & -1 \\ -1 & \frac{1}{2} \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 3 \end{bmatrix} = \begin{bmatrix} \frac{4}{3} & \frac{1}{3} & -\frac{2}{3} \\ -\frac{1}{2} & 0 & \frac{1}{2} \end{bmatrix}$$

The pseudoinverse is a convenient shorthand notation for some of the concepts we have been exploring. For example, if A is $m \times n$ with linearly independent columns, the least squares solution of $A\mathbf{x} = \mathbf{b}$ is given by

$$\bar{\mathbf{x}} = A^+ \mathbf{b}$$

and the standard matrix of the orthogonal projection P from \mathbb{R}^m onto $\text{col}(A)$ is

$$[P] = AA^+$$

If A is actually a square matrix, then it is easy to show that $A^+ = A^{-1}$ (see Exercise 53). In this case, the least squares solution of $A\mathbf{x} = \mathbf{b}$ is the *exact* solution, since

$$\bar{\mathbf{x}} = A^+ \mathbf{b} = A^{-1} \mathbf{b} = \mathbf{x}$$

→ The projection matrix becomes $[P] = AA^+ = AA^{-1} = I$. (What is the geometric interpretation of this equality?)

→ Theorem 7.12 summarizes the key properties of the pseudoinverse of a matrix. (Before reading the proof of this theorem, verify these properties for the matrix in Example 7.32.)

Theorem 7.12

Let A be a matrix with linearly independent columns. Then the pseudoinverse A^+ of A satisfies the following properties, called the **Penrose conditions** for A :

- $AA^+A = A$
- $A^+AA^+ = A^+$
- AA^+ and A^+A are symmetric.

Proof We prove condition (a) and half of condition (c) and leave the proofs of the remaining conditions as Exercises 54 and 55.

(a) We compute

$$\begin{aligned} AA^+A &= A((A^T A)^{-1} A^T)A \\ &= A(A^T A)^{-1}(A^T A) \\ &= AI = A \end{aligned}$$

(c) By Theorem 3.4, $A^T A$ is symmetric. Therefore, $(A^T A)^{-1}$ is also symmetric, by Exercise 46 in Section 3.3. Taking the transpose of AA^+ , we have

$$\begin{aligned} (AA^+)^T &= (A(A^T A)^{-1} A^T)^T \\ &= (A^T)^T ((A^T A)^{-1})^T A^T \\ &= A(A^T A)^{-1} A^T \\ &= AA^+ \end{aligned}$$

Exercise 56 explores further properties of the pseudoinverse. In the next section, we will see how to extend the definition of A^+ to handle *all* matrices, whether or not the columns of A are linearly independent.

Exercises 7.3

CAS

In Exercises 1–3, consider the data points $(1, 0)$, $(2, 1)$, and $(3, 5)$. Compute the least squares error for the given line. In each case, plot the points and the line.

1. $y = -2 + 2x$ 2. $y = x$ 3. $y = -3 + \frac{5}{2}x$

In Exercises 4–6, consider the data points $(-5, 3)$, $(0, 3)$, $(5, 2)$, and $(10, 0)$. Compute the least squares error for the given line. In each case, plot the points and the line.

4. $y = 3 - \frac{1}{3}x$ 5. $y = \frac{5}{2}$ 6. $y = 2 - \frac{1}{5}x$

In Exercises 7–14, find the least squares approximating line for the given points and compute the corresponding least squares error.

- $(1, 0)$, $(2, 1)$, $(3, 5)$
- $(1, 6)$, $(2, 3)$, $(3, 1)$
- $(0, 4)$, $(1, 1)$, $(2, 0)$
- $(0, 3)$, $(1, 3)$, $(2, 5)$
- $(-5, -1)$, $(0, 1)$, $(5, 2)$, $(10, 4)$