

Define

$$\langle x, y \rangle = \sum_{i=1}^n a_i \bar{b}_i.$$

- (a) Prove that $\langle \cdot, \cdot \rangle$ is an inner product on V and that β is an orthonormal basis for V . Thus every real or complex vector space may be regarded as an inner product space.
 - (b) Prove that if $V = \mathbb{R}^n$ or $V = \mathbb{C}^n$ and β is the standard ordered basis, then the inner product defined above is the standard inner product.
23. Let $V = F^n$, and let $A \in M_{n \times n}(F)$.
- (a) Prove that $\langle x, Ay \rangle = \langle A^*x, y \rangle$ for all $x, y \in V$.
 - (b) Suppose that for some $B \in M_{n \times n}(F)$, we have $\langle x, Ay \rangle = \langle Bx, y \rangle$ for all $x, y \in V$. Prove that $B = A^*$.
 - (c) Let α be the standard ordered basis for V . For any orthonormal basis β for V , let Q be the $n \times n$ matrix whose columns are the vectors in β . Prove that $Q^* = Q^{-1}$.
 - (d) Define linear operators T and U on V by $T(x) = Ax$ and $U(x) = A^*x$. Show that $[U]_\beta = [T]_\beta^*$ for any orthonormal basis β for V .
24. Let V be a complex inner product space with an inner product $\langle \cdot, \cdot \rangle$. Let $[\cdot, \cdot]$ be the real-valued function such that $[x, y]$ is the real part of the complex number $\langle x, y \rangle$ for all $x, y \in V$. Prove that $[\cdot, \cdot]$ is an inner product for V , where V is regarded as a vector space over R . Prove, furthermore, that $[x, ix] = 0$ for all $x \in V$.
25. Let V be a vector space over C , and suppose that $[\cdot, \cdot]$ is a real inner product on V , where V is regarded as a vector space over R , such that $[x, ix] = 0$ for all $x \in V$. Let $\langle \cdot, \cdot \rangle$ be the complex-valued function defined by
- $$\langle x, y \rangle = [x, y] + i[x, iy] \quad \text{for } x, y \in V.$$

Prove that $\langle \cdot, \cdot \rangle$ is a complex inner product on V .

The following definition is used in Exercises 26–30.

Definition. Let V be a vector space over F , where F is either R or C . Regardless of whether V is or is not an inner product space, we may still define a norm $\|\cdot\|_V$ as a real-valued function on V satisfying the following three conditions for all $x, y \in V$ and $a \in F$:

- (1) $\|x\|_V \geq 0$, and $\|x\|_V = 0$ if and only if $x = 0$.
- (2) $\|ax\|_V = |a| \cdot \|x\|_V$.
- (3) $\|x + y\|_V \leq \|x\|_V + \|y\|_V$.

26. Prove that the following are norms on the given vector spaces V .
- $V = \mathbb{R}^2$; $\|(a, b)\|_v = |a| + |b|$ for all $(a, b) \in V$
 - $V = C([0, 1])$; $\|f\|_v = \max_{t \in [0, 1]} |f(t)|$ for all $f \in V$
 - $V = C([0, 1])$; $\|f\|_v = \int_0^1 |f(t)| dt$ for all $f \in V$
 - $V = M_{m \times n}(F)$; $\|A\|_v = \max_{i,j} |A_{ij}|$ for all $A \in V$
27. Use Exercise 11 to show that there is no inner product $\langle \cdot, \cdot \rangle$ on \mathbb{R}^2 such that $\|x\|_v^2 = \langle x, x \rangle$ for all $x \in \mathbb{R}^2$ if the norm is defined as in Exercise 26(a).
28. Let $\|\cdot\|_v$ be a norm on a vector space V , and define, for each ordered pair of vectors, the scalar $d(x, y) = \|x - y\|_v$, called the **distance** between x and y . Prove the following results for all $x, y, z \in V$.
- $d(x, y) \geq 0$.
 - $d(x, y) = d(y, x)$.
 - $d(x, y) \leq d(x, z) + d(z, y)$.
 - $d(x, x) = 0$ if and only if $x = 0$.
 - $d(x, y) \neq 0$ if $x \neq y$.
29. Let $\|\cdot\|_v$ be a norm on a real vector space V satisfying the parallelogram law given in Exercise 11. Define

$$\langle x, y \rangle = \frac{1}{4} [\|x + y\|_v^2 - \|x - y\|_v^2].$$

Prove that $\langle \cdot, \cdot \rangle$ defines an inner product on V such that $\|x\|_v^2 = \langle x, x \rangle$ for all $x \in V$. *Hints:*

- Prove $\langle x, 2y \rangle = 2\langle x, y \rangle$ for all $x, y \in V$.
- Prove $\langle x + u, y \rangle = \langle x, y \rangle + \langle u, y \rangle$ for all $x, u, y \in V$.
- Prove $\langle nx, y \rangle = n\langle x, y \rangle$ for every positive integer n and every $x, y \in V$.
- Prove $m\langle \frac{1}{m}x, y \rangle = \langle x, y \rangle$ for every positive integer m and every $x, y \in V$.
- Prove $\langle rx, y \rangle = r\langle x, y \rangle$ for every rational number r and every $x, y \in V$.
- Prove $|\langle x, y \rangle| \leq \|x\|_v\|y\|_v$ for every $x, y \in V$. *Hint:* Condition (3) in the definition of norm can be helpful.
- Prove that for every $c \in R$, every rational number r , and every $x, y \in V$,

$$\begin{aligned} |c\langle x, y \rangle - \langle cx, y \rangle| &= |(c - r)\langle x, y \rangle - \langle (c - r)x, y \rangle| \\ &\leq 2|c - r|\|x\|_v\|y\|_v. \end{aligned}$$

- (h) Use the fact that for any $c \in R$, $|c - r|$ can be made arbitrarily small, where r varies over the set of rational numbers, to establish item (b) of the definition of *inner product*.
30. Let $\|\cdot\|_V$ be a norm (as defined on page 337) on a complex vector space V satisfying the parallelogram law given in Exercise 11. Prove that there is an inner product $\langle \cdot, \cdot \rangle$ on V such that $\|x\|_V^2 = \langle x, x \rangle$ for all $x \in V$. *Hint:* Apply Exercise 29 to V regarded as a vector space over R . Then apply Exercise 25.

6.2 THE GRAM-SCHMIDT ORTHOGONALIZATION PROCESS AND ORTHOGONAL COMPLEMENTS

In previous chapters, we have seen the special role of the standard ordered bases for C^n and R^n . The special properties of these bases stem from the fact that the basis vectors form an orthonormal set. Just as bases are the building blocks of vector spaces, bases that are also orthonormal sets are the building blocks of inner product spaces. We now name such bases.

Definition. Let V be an inner product space. A subset of V is an **orthonormal basis** for V if it is an ordered basis that is orthonormal.

Example 1

The standard ordered basis for F^n is an orthonormal basis for F^n . ◆

Example 2

The set

$$\left\{ \left(\frac{1}{\sqrt{5}}, \frac{2}{\sqrt{5}} \right), \left(\frac{2}{\sqrt{5}}, \frac{-1}{\sqrt{5}} \right) \right\}$$

is an orthonormal basis for R^2 . ◆

The next theorem and its corollaries illustrate why orthonormal sets and, in particular, orthonormal bases are so important.

Theorem 6.3. Let V be an inner product space and $S = \{v_1, v_2, \dots, v_k\}$ be an orthogonal subset of V consisting of nonzero vectors. If $y \in \text{span}(S)$, then

$$y = \sum_{i=1}^k \frac{\langle y, v_i \rangle}{\|v_i\|^2} v_i.$$

Proof. Write $y = \sum_{i=1}^k a_i v_i$, where $a_1, a_2, \dots, a_k \in F$. Then, for $1 \leq j \leq k$,

we have

$$\langle y, v_j \rangle = \left\langle \sum_{i=1}^k a_i v_i, v_j \right\rangle = \sum_{i=1}^k a_i \langle v_i, v_j \rangle = a_j \langle v_j, v_j \rangle = a_j \|v_j\|^2.$$

So $a_j = \frac{\langle y, v_j \rangle}{\|v_j\|^2}$, and the result follows. ■

The next corollary follows immediately from Theorem 6.3.

Corollary 1. *If, in addition to the hypotheses of Theorem 6.3, S is orthonormal and $y \in \text{span}(S)$, then*

$$y = \sum_{i=1}^k \langle y, v_i \rangle v_i.$$

If V possesses a finite orthonormal basis, then Corollary 1 allows us to compute the coefficients in a linear combination very easily. (See Example 3.)

Corollary 2. *Let V be an inner product space, and let S be an orthogonal subset of V consisting of nonzero vectors. Then S is linearly independent.*

Proof. Suppose that $v_1, v_2, \dots, v_k \in S$ and

$$\sum_{i=1}^k a_i v_i = 0.$$

As in the proof of Theorem 6.3 with $y = 0$, we have $a_j = \langle 0, v_j \rangle / \|v_j\|^2 = 0$ for all j . So S is linearly independent. ■

Example 3

By Corollary 2, the orthonormal set

$$\left\{ \frac{1}{\sqrt{2}}(1, 1, 0), \frac{1}{\sqrt{3}}(1, -1, 1), \frac{1}{\sqrt{6}}(-1, 1, 2) \right\}$$

obtained in Example 8 of Section 6.1 is an orthonormal basis for \mathbb{R}^3 . Let $x = (2, 1, 3)$. The coefficients given by Corollary 1 to Theorem 6.3 that express x as a linear combination of the basis vectors are

$$a_1 = \frac{1}{\sqrt{2}}(2 + 1) = \frac{3}{\sqrt{2}}, \quad a_2 = \frac{1}{\sqrt{3}}(2 - 1 + 3) = \frac{4}{\sqrt{3}},$$

and

$$a_3 = \frac{1}{\sqrt{6}}(-2 + 1 + 6) = \frac{5}{\sqrt{6}}.$$

As a check, we have

$$(2, 1, 3) = \frac{3}{2}(1, 1, 0) + \frac{4}{3}(1, -1, 1) + \frac{5}{6}(-1, 1, 2). \quad \blacklozenge$$

Corollary 2 tells us that the vector space H in Section 6.1 contains an infinite linearly independent set, and hence H is not a finite-dimensional vector space.

Of course, we have not yet shown that every finite-dimensional inner product space possesses an orthonormal basis. The next theorem takes us most of the way in obtaining this result. It tells us how to construct an orthogonal set from a linearly independent set of vectors in such a way that both sets generate the same subspace.

Before stating this theorem, let us consider a simple case. Suppose that $\{w_1, w_2\}$ is a linearly independent subset of an inner product space (and hence a basis for some two-dimensional subspace). We want to construct an orthogonal set from $\{w_1, w_2\}$ that spans the same subspace. Figure 6.1 suggests that the set $\{v_1, v_2\}$, where $v_1 = w_1$ and $v_2 = w_2 - cw_1$, has this property if c is chosen so that v_2 is orthogonal to w_1 .

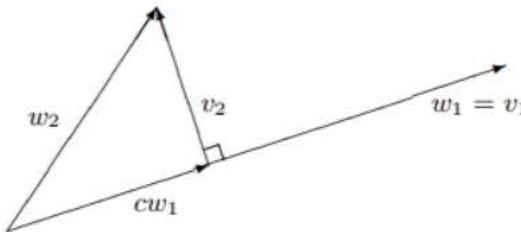


Figure 6.1

To find c , we need only solve the following equation:

$$0 = \langle v_2, w_1 \rangle = \langle w_2 - cw_1, w_1 \rangle = \langle w_2, w_1 \rangle - c \langle w_1, w_1 \rangle.$$

So

$$c = \frac{\langle w_2, w_1 \rangle}{\|w_1\|^2}.$$

Thus

$$v_2 = w_2 - \frac{\langle w_2, w_1 \rangle}{\|w_1\|^2} w_1.$$

The next theorem shows us that this process can be extended to any finite linearly independent subset.

Theorem 6.4. *Let V be an inner product space and $S = \{w_1, w_2, \dots, w_n\}$ be a linearly independent subset of V . Define $S' = \{v_1, v_2, \dots, v_n\}$, where $v_1 = w_1$ and*

$$v_k = w_k - \sum_{j=1}^{k-1} \frac{\langle w_k, v_j \rangle}{\|v_j\|^2} v_j \quad \text{for } 2 \leq k \leq n. \quad (1)$$

Then S' is an orthogonal set of nonzero vectors such that $\text{span}(S') = \text{span}(S)$.

Proof. The proof is by mathematical induction on n , the number of vectors in S . For $k = 1, 2, \dots, n$, let $S_k = \{w_1, w_2, \dots, w_k\}$. If $n = 1$, then the theorem is proved by taking $S'_1 = S_1$; i.e., $v_1 = w_1 \neq 0$. Assume then that the set $S'_{k-1} = \{v_1, v_2, \dots, v_{k-1}\}$ with the desired properties has been constructed by the repeated use of (1). We show that the set $S'_k = \{v_1, v_2, \dots, v_{k-1}, v_k\}$ also has the desired properties, where v_k is obtained from S'_{k-1} by (1). If $v_k = 0$, then (1) implies that $w_k \in \text{span}(S'_{k-1}) = \text{span}(S_{k-1})$, which contradicts the assumption that S_k is linearly independent. For $1 \leq i \leq k-1$, it follows from (1) that

$$\langle v_k, v_i \rangle = \langle w_k, v_i \rangle - \sum_{j=1}^{k-1} \frac{\langle w_k, v_j \rangle}{\|v_j\|^2} \langle v_j, v_i \rangle = \langle w_k, v_i \rangle - \frac{\langle w_k, v_i \rangle}{\|v_i\|^2} \|v_i\|^2 = 0,$$

since $\langle v_j, v_i \rangle = 0$ if $i \neq j$ by the induction assumption that S'_{k-1} is orthogonal. Hence S'_k is an orthogonal set of nonzero vectors. Now, by (1), we have that $\text{span}(S'_k) \subseteq \text{span}(S_k)$. But by Corollary 2 to Theorem 6.3, S'_k is linearly independent; so $\dim(\text{span}(S'_k)) = \dim(\text{span}(S_k)) = k$. Therefore $\text{span}(S'_k) = \text{span}(S_k)$. ■

The construction of $\{v_1, v_2, \dots, v_n\}$ by the use of Theorem 6.4 is called the **Gram–Schmidt process**.

Example 4

In \mathbb{R}^4 , let $w_1 = (1, 0, 1, 0)$, $w_2 = (1, 1, 1, 1)$, and $w_3 = (0, 1, 2, 1)$. Then $\{w_1, w_2, w_3\}$ is linearly independent. We use the Gram–Schmidt process to compute the orthogonal vectors v_1 , v_2 , and v_3 , and then we normalize these vectors to obtain an orthonormal set.

Take $v_1 = w_1 = (1, 0, 1, 0)$. Then

$$\begin{aligned} v_2 &= w_2 - \frac{\langle w_2, v_1 \rangle}{\|v_1\|^2} v_1 \\ &= (1, 1, 1, 1) - \frac{2}{2} (1, 0, 1, 0) \\ &= (0, 1, 0, 1). \end{aligned}$$

Finally,

$$\begin{aligned} v_3 &= w_3 - \frac{\langle w_3, v_1 \rangle}{\|v_1\|^2} v_1 - \frac{\langle w_3, v_2 \rangle}{\|v_2\|^2} v_2 \\ &= (0, 1, 2, 1) - \frac{2}{2} (1, 0, 1, 0) - \frac{2}{2} (0, 1, 0, 1) \\ &= (-1, 0, 1, 0). \end{aligned}$$

These vectors can be normalized to obtain the orthonormal basis $\{u_1, u_2, u_3\}$, where

$$u_1 = \frac{1}{\|v_1\|} v_1 = \frac{1}{\sqrt{2}}(1, 0, 1, 0),$$

$$u_2 = \frac{1}{\|v_2\|} v_2 = \frac{1}{\sqrt{2}}(0, 1, 0, 1),$$

and

$$u_3 = \frac{v_3}{\|v_3\|} = \frac{1}{\sqrt{2}}(-1, 0, 1, 0). \quad \blacklozenge$$

Example 5

Let $V = \mathbb{P}(R)$ with the inner product $\langle f(x), g(x) \rangle = \int_{-1}^1 f(t)g(t) dt$, and consider the subspace $\mathbb{P}_2(R)$ with the standard ordered basis β . We use the Gram-Schmidt process to replace β by an orthogonal basis $\{v_1, v_2, v_3\}$ for $\mathbb{P}_2(R)$, and then use this orthogonal basis to obtain an orthonormal basis for $\mathbb{P}_2(R)$.

Take $v_1 = 1$. Then $\|v_1\|^2 = \int_{-1}^1 1^2 dt = 2$, and $\langle x, v_1 \rangle = \int_{-1}^1 t \cdot 1 dt = 0$. Thus

$$v_2 = x - \frac{\langle v_1, x \rangle}{\|v_1\|^2} v_1 = x - \frac{0}{2} = x.$$

Furthermore,

$$\langle x^2, v_1 \rangle = \int_{-1}^1 t^2 \cdot 1 dt = \frac{2}{3} \quad \text{and} \quad \langle x^2, v_2 \rangle = \int_{-1}^1 t^2 \cdot t dt = 0.$$

Therefore

$$\begin{aligned} v_3 &= x^2 - \frac{\langle x^2, v_1 \rangle}{\|v_1\|^2} v_1 - \frac{\langle x^2, v_2 \rangle}{\|v_2\|^2} v_2 \\ &= x^2 - \frac{1}{3} \cdot 1 - 0 \cdot x \\ &= x^2 - \frac{1}{3}. \end{aligned}$$

We conclude that $\{1, x, x^2 - \frac{1}{3}\}$ is an orthogonal basis for $\mathbb{P}_2(R)$.

To obtain an orthonormal basis, we normalize v_1 , v_2 , and v_3 to obtain

$$u_1 = \frac{1}{\sqrt{\int_{-1}^1 1^2 dt}} = \frac{1}{\sqrt{2}},$$

$$u_2 = \frac{x}{\sqrt{\int_{-1}^1 t^2 dt}} = \sqrt{\frac{3}{2}} x,$$

and similarly,

$$u_3 = \frac{v_3}{\|v_3\|} = \sqrt{\frac{5}{8}} (3x^2 - 1).$$

Thus $\{u_1, u_2, u_3\}$ is the desired orthonormal basis for $P_2(R)$. \blacklozenge

Continuing to apply the Gram–Schmidt orthogonalization process to the basis $\{1, x, x^2, \dots\}$ for $P(R)$, we obtain an orthogonal basis $\{v_1, v_2, v_3, \dots\}$. For each n , the polynomial $(1/v_k(1))v_k$ is called the k th *Legendre polynomial*. The first three Legendre polynomials are 1 , x and $\frac{1}{2}(3x^2 - 1)$. The set of Legendre polynomials is also an orthogonal basis for $P(R)$.

The following result gives us a simple method of representing a vector as a linear combination of the vectors in an orthonormal basis.

Theorem 6.5. *Let V be a nonzero finite-dimensional inner product space. Then V has an orthonormal basis β . Furthermore, if $\beta = \{v_1, v_2, \dots, v_n\}$ and $x \in V$, then*

$$x = \sum_{i=1}^n \langle x, v_i \rangle v_i.$$

Proof. Let β_0 be an ordered basis for V . Apply Theorem 6.4 to obtain an orthogonal set β' of nonzero vectors with $\text{span}(\beta') = \text{span}(\beta_0) = V$. By normalizing each vector in β' , we obtain an orthonormal set β that generates V . By Corollary 2 to Theorem 6.3, β is linearly independent; therefore β is an orthonormal basis for V . The remainder of the theorem follows from Corollary 1 to Theorem 6.3. \blacksquare

Example 6

We use Theorem 6.5 to represent the polynomial $f(x) = 1 + 2x + 3x^2$ as a linear combination of the vectors in the orthonormal basis $\{u_1, u_2, u_3\}$ for $P_2(R)$ obtained in Example 5. Observe that

$$\langle f(x), u_1 \rangle = \int_{-1}^1 \frac{1}{\sqrt{2}} (1 + 2t + 3t^2) dt = 2\sqrt{2},$$

$$\langle f(x), u_2 \rangle = \int_{-1}^1 \sqrt{\frac{3}{2}} t (1 + 2t + 3t^2) dt = \frac{2\sqrt{6}}{3},$$

and

$$\langle f(x), u_3 \rangle = \int_{-1}^1 \sqrt{\frac{5}{8}} (3t^2 - 1) (1 + 2t + 3t^2) dt = \frac{2\sqrt{10}}{5}.$$

Therefore $f(x) = 2\sqrt{2} u_1 + \frac{2\sqrt{6}}{3} u_2 + \frac{2\sqrt{10}}{5} u_3$. \blacklozenge

Theorem 6.5 gives us a simple method for computing the entries of the matrix representation of a linear operator with respect to an orthonormal basis.

Corollary. *Let V be a finite-dimensional inner product space with an orthonormal basis $\beta = \{v_1, v_2, \dots, v_n\}$. Let T be a linear operator on V , and let $A = [T]_\beta$. Then for any i and j , $A_{ij} = \langle T(v_j), v_i \rangle$.*

Proof. From Theorem 6.5, we have

$$T(v_j) = \sum_{i=1}^n \langle T(v_j), v_i \rangle v_i.$$

Hence $A_{ij} = \langle T(v_j), v_i \rangle$. \blacksquare

The scalars $\langle x, v_i \rangle$ given in Theorem 6.5 have been studied extensively for special inner product spaces. Although the vectors v_1, v_2, \dots, v_n were chosen from an orthonormal basis, we introduce a terminology associated with orthonormal sets β in more general inner product spaces.

Definition. *Let β be an orthonormal subset (possibly infinite) of an inner product space V , and let $x \in V$. We define the **Fourier coefficients** of x relative to β to be the scalars $\langle x, y \rangle$, where $y \in \beta$.*

In the first half of the 19th century, the French mathematician Jean Baptiste Fourier was associated with the study of the scalars

$$\int_0^{2\pi} f(t) \sin nt dt \quad \text{and} \quad \int_0^{2\pi} f(t) \cos nt dt,$$

or in the complex case,

$$c_n = \frac{1}{2\pi} \int_0^{2\pi} f(t) e^{-int} dt,$$

for a function f . In the context of Example 9 of Section 6.1, we see that $c_n = \langle f, f_n \rangle$, where $f_n(t) = e^{int}$; that is, c_n is the n th Fourier coefficient for a continuous function $f \in V$ relative to S . The coefficients c_n are the “classical” Fourier coefficients of a function, and the literature concerning their behavior is extensive. We learn more about Fourier coefficients in the remainder of this chapter.

Example 7

Let $S = \{e^{int} : n \text{ is an integer}\}$. In Example 9 of Section 6.1, S was shown to be an orthonormal set in \mathbb{H} . We compute the Fourier coefficients of $f(t) = t$ relative to S . Using integration by parts, we have, for $n \neq 0$,

$$\langle f, f_n \rangle = \frac{1}{2\pi} \int_0^{2\pi} te^{\overline{int}} dt = \frac{1}{2\pi} \int_0^{2\pi} te^{-int} dt = \frac{-1}{in},$$

and, for $n = 0$,

$$\langle f, 1 \rangle = \frac{1}{2\pi} \int_0^{2\pi} t(1) dt = \pi.$$

As a result of these computations, and using Exercise 16 of this section, we obtain an upper bound for the sum of a special infinite series as follows:

$$\begin{aligned} \|f\|^2 &\geq \sum_{n=-k}^{-1} |\langle f, f_n \rangle|^2 + |\langle f, 1 \rangle|^2 + \sum_{n=1}^k |\langle f, f_n \rangle|^2 \\ &= \sum_{n=-k}^{-1} \frac{1}{n^2} + \pi^2 + \sum_{n=1}^k \frac{1}{n^2} \\ &= 2 \sum_{n=1}^k \frac{1}{n^2} + \pi^2 \end{aligned}$$

for every k . Now, using the fact that $\|f\|^2 = \frac{4}{3}\pi^2$, we obtain

$$\frac{4}{3}\pi^2 \geq 2 \sum_{n=1}^k \frac{1}{n^2} + \pi^2,$$

or

$$\frac{\pi^2}{6} \geq \sum_{n=1}^k \frac{1}{n^2}.$$

Because this inequality holds for all k , we may let $k \rightarrow \infty$ to obtain

$$\frac{\pi^2}{6} \geq \sum_{n=1}^{\infty} \frac{1}{n^2}.$$

Additional results may be produced by replacing f by other functions. ♦

We are now ready to proceed with the concept of an *orthogonal complement*.

Definition. Let S be a nonempty subset of an inner product space V . We define S^\perp (read “ S perp”) to be the set of all vectors in V that are orthogonal to every vector in S ; that is, $S^\perp = \{x \in V : \langle x, y \rangle = 0 \text{ for all } y \in S\}$. The set S^\perp is called the **orthogonal complement** of S .

It is easily seen that S^\perp is a subspace of V for any subset S of V .

Example 8

The reader should verify that $\{\theta\}^\perp = V$ and $V^\perp = \{\theta\}$ for any inner product space V . ♦

Example 9

If $V = \mathbb{R}^3$ and $S = \{e_3\}$, then S^\perp equals the xy -plane (see Exercise 5). ♦

Exercise 18 provides an interesting example of an orthogonal complement in an infinite-dimensional inner product space.

Consider the problem in \mathbb{R}^3 of finding the distance from a point P to a plane W . (See Figure 6.2.) Problems of this type arise in many settings. If we let y be the vector determined by θ and P , we may restate the problem as follows: Determine the vector u in W that is “closest” to y . The desired distance is clearly given by $\|y - u\|$. Notice from the figure that the vector $z = y - u$ is orthogonal to every vector in W , and so $z \in W^\perp$.

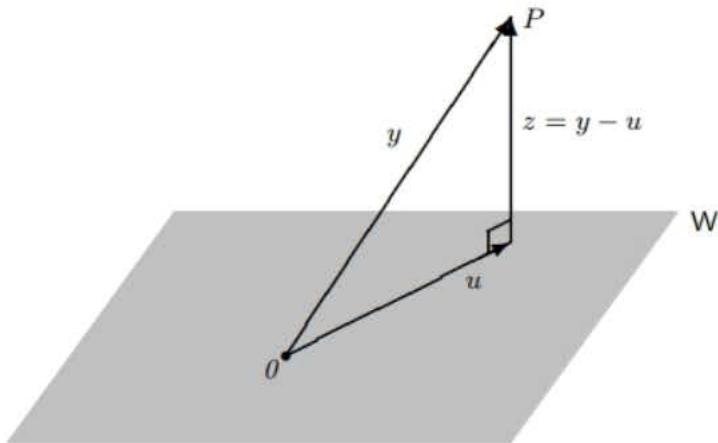


Figure 6.2

The next result presents a practical method of finding u in the case that W is a finite-dimensional subspace of an inner product space.

Theorem 6.6. Let W be a finite-dimensional subspace of an inner product space V , and let $y \in V$. Then there exist unique vectors $u \in W$ and $z \in W^\perp$

such that $y = u + z$. Furthermore, if $\{v_1, v_2, \dots, v_k\}$ is an orthonormal basis for W , then

$$u = \sum_{i=1}^k \langle y, v_i \rangle v_i.$$

Proof. Let $\{v_1, v_2, \dots, v_k\}$ be an orthonormal basis for W , let u be as defined in the preceding equation, and let $z = y - u$. Clearly $u \in W$ and $y = u + z$.

To show that $z \in W^\perp$, it suffices to show, by Exercise 7, that z is orthogonal to each v_j . For any j , we have

$$\begin{aligned} \langle z, v_j \rangle &= \left\langle \left(y - \sum_{i=1}^k \langle y, v_i \rangle v_i \right), v_j \right\rangle = \langle y, v_j \rangle - \sum_{i=1}^k \langle y, v_i \rangle \langle v_i, v_j \rangle \\ &= \langle y, v_j \rangle - \langle y, v_j \rangle = 0. \end{aligned}$$

To show uniqueness of u and z , suppose that $y = u + z = u' + z'$, where $u' \in W$ and $z' \in W^\perp$. Then $u - u' = z' - z \in W \cap W^\perp = \{0\}$. Therefore $u = u'$ and $z = z'$. ■

Corollary. In the notation of Theorem 6.6, the vector u is the unique vector in W that is “closest” to y ; that is, for any $x \in W$, $\|y - x\| \geq \|y - u\|$, and this inequality is an equality if and only if $x = u$.

Proof. As in Theorem 6.6, we have that $y = u + z$, where $z \in W^\perp$. Let $x \in W$. Then $u - x$ is orthogonal to z , so, by Exercise 10 of Section 6.1, we have

$$\begin{aligned} \|y - x\|^2 &= \|u + z - x\|^2 = \|(u - x) + z\|^2 = \|u - x\|^2 + \|z\|^2 \\ &\geq \|z\|^2 = \|y - u\|^2. \end{aligned}$$

Now suppose that $\|y - x\| = \|y - u\|$. Then the inequality above becomes an equality, and therefore $\|u - x\|^2 + \|z\|^2 = \|z\|^2$. It follows that $\|u - x\| = 0$, and hence $x = u$. The proof of the converse is obvious. ■

The vector u in the corollary is called the **orthogonal projection** of y on W . We will see the importance of orthogonal projections of vectors in the application to least squares in Section 6.3.

Example 10

Let $V = P_3(R)$ with the inner product

$$\langle f(x), g(x) \rangle = \int_{-1}^1 f(t)g(t) dt \quad \text{for all } f(x), g(x) \in V.$$

We compute the orthogonal projection $f_1(x)$ of $f(x) = x^3$ on $\mathbb{P}_2(R)$.

By Example 5,

$$\{u_1, u_2, u_3\} = \left\{ \frac{1}{\sqrt{2}}, \sqrt{\frac{3}{2}}x, \sqrt{\frac{5}{8}}(3x^2 - 1) \right\}$$

is an orthonormal basis for $\mathbb{P}_2(R)$. For these vectors, we have

$$\langle f(x), u_1 \rangle = \int_{-1}^1 t^3 \frac{1}{\sqrt{2}} dt = 0, \quad \langle f(x), u_2 \rangle = \int_{-1}^1 t^3 \sqrt{\frac{3}{2}} t dt = \frac{\sqrt{6}}{5},$$

and

$$\langle f(x), u_3 \rangle = \int_{-1}^1 t^3 \sqrt{\frac{5}{8}} (3t^2 - 1) dt = 0.$$

Hence

$$f_1(x) = \langle f(x), u_1 \rangle u_1 + \langle f(x), u_2 \rangle u_2 + \langle f(x), u_3 \rangle u_3 = \frac{3}{5}x. \quad \blacklozenge$$

It was shown (Corollary 2 to the replacement theorem, p. 48) that any linearly independent set in a finite-dimensional vector space can be extended to a basis. The next theorem provides an interesting analog for an orthonormal subset of a finite-dimensional inner product space.

Theorem 6.7. Suppose that $S = \{v_1, v_2, \dots, v_k\}$ is an orthonormal set in an n -dimensional inner product space V . Then

- (a) S can be extended to an orthonormal basis $\{v_1, v_2, \dots, v_k, v_{k+1}, \dots, v_n\}$ for V .
- (b) If $W = \text{span}(S)$, then $S_1 = \{v_{k+1}, v_{k+2}, \dots, v_n\}$ is an orthonormal basis for W^\perp (using the preceding notation).
- (c) If W is any subspace of V , then $\dim(V) = \dim(W) + \dim(W^\perp)$.

Proof. (a) By Corollary 2 to the replacement theorem (p. 48), S can be extended to an ordered basis $S' = \{v_1, v_2, \dots, v_k, w_{k+1}, \dots, w_n\}$ for V . Now apply the Gram-Schmidt process to S' . The first k vectors resulting from this process are the vectors in S by Exercise 8, and this new set spans V . Normalizing the last $n - k$ vectors of this set produces an orthonormal set that spans V . The result now follows.

(b) Because S_1 is a subset of a basis, it is linearly independent. Since S_1 is clearly a subset of W^\perp , we need only show that it spans W^\perp . Note that, for any $x \in V$, we have

$$x = \sum_{i=1}^n \langle x, v_i \rangle v_i.$$

If $x \in W^\perp$, then $\langle x, v_i \rangle = 0$ for $1 \leq i \leq k$. Therefore

$$x = \sum_{i=k+1}^n \langle x, v_i \rangle v_i \in \text{span}(S_1).$$

(c) Let W be a subspace of V . It is a finite-dimensional inner product space because V is, and so it has an orthonormal basis $\{v_1, v_2, \dots, v_k\}$. By (a) and (b), we have

$$\dim(V) = n = k + (n - k) = \dim(W) + \dim(W^\perp). \quad \blacksquare$$

Example 11

Let $W = \text{span}(\{e_1, e_2\})$ in \mathbb{F}^3 . Then $x = (a, b, c) \in W^\perp$ if and only if $0 = \langle x, e_1 \rangle = a$ and $0 = \langle x, e_2 \rangle = b$. So $x = (0, 0, c)$, and therefore $W^\perp = \text{span}(\{e_3\})$. One can deduce the same result by noting that $e_3 \in W^\perp$ and, from (c), that $\dim(W^\perp) = 3 - 2 = 1$. \blacklozenge

EXERCISES

1. Label the following statements as true or false.
 - (a) The Gram–Schmidt orthogonalization process produces an orthonormal set from an arbitrary linearly independent set.
 - (b) Every nonzero finite-dimensional inner product space has an orthonormal basis.
 - (c) The orthogonal complement of any set is a subspace.
 - (d) If $\{v_1, v_2, \dots, v_n\}$ is a basis for an inner product space V , then for any $x \in V$ the scalars $\langle x, v_i \rangle$ are the Fourier coefficients of x .
 - (e) An orthonormal basis must be an ordered basis.
 - (f) Every orthogonal set is linearly independent.
 - (g) Every orthonormal set is linearly independent.
2. In each part, apply the Gram–Schmidt process to the given subset S of the inner product space V to obtain an orthogonal basis for $\text{span}(S)$. Then normalize the vectors in this basis to obtain an orthonormal basis β for $\text{span}(S)$, and compute the Fourier coefficients of the given vector relative to β . Finally, use Theorem 6.5 to verify your result.
 - (a) $V = \mathbb{R}^3$, $S = \{(1, 0, 1), (0, 1, 1), (1, 3, 3)\}$, and $x = (1, 1, 2)$
 - (b) $V = \mathbb{R}^3$, $S = \{(1, 1, 1), (0, 1, 1), (0, 0, 1)\}$, and $x = (1, 0, 1)$
 - (c) $V = P_2(R)$ with the inner product $\langle f(x), g(x) \rangle = \int_0^1 f(t)g(t) dt$, $S = \{1, x, x^2\}$, and $h(x) = 1 + x$
 - (d) $V = \text{span}(S)$, where $S = \{(1, i, 0), (1 - i, 2, 4i)\}$, and $x = (3 + i, 4i, -4)$

- (e) $V = \mathbb{R}^4$, $S = \{(2, -1, -2, 4), (-2, 1, -5, 5), (-1, 3, 7, 11)\}$, and $x = (-11, 8, -4, 18)$
- (f) $V = \mathbb{R}^4$, $S = \{(1, -2, -1, 3), (3, 6, 3, -1), (1, 4, 2, 8)\}$, and $x = (-1, 2, 1, 1)$
- (g) $V = M_{2 \times 2}(R)$, $S = \left\{ \begin{pmatrix} 3 & 5 \\ -1 & 1 \end{pmatrix}, \begin{pmatrix} -1 & 9 \\ 5 & -1 \end{pmatrix}, \begin{pmatrix} 7 & -17 \\ 2 & -6 \end{pmatrix} \right\}$, and $A = \begin{pmatrix} -1 & 27 \\ -4 & 8 \end{pmatrix}$
- (h) $V = M_{2 \times 2}(R)$, $S = \left\{ \begin{pmatrix} 2 & 2 \\ 2 & 1 \end{pmatrix}, \begin{pmatrix} 11 & 4 \\ 2 & 5 \end{pmatrix}, \begin{pmatrix} 4 & -12 \\ 3 & -16 \end{pmatrix} \right\}$, and $A = \begin{pmatrix} 8 & 6 \\ 25 & -13 \end{pmatrix}$
- (i) $V = \text{span}(S)$ with the inner product $\langle f, g \rangle = \int_0^\pi f(t)g(t) dt$, $S = \{\sin t, \cos t, 1, t\}$, and $h(t) = 2t + 1$
- (j) $V = \mathbb{C}^4$, $S = \{(1, i, 2-i, -1), (2+3i, 3i, 1-i, 2i), (-1+7i, 6+10i, 11-4i, 3+4i)\}$, and $x = (-2+7i, 6+9i, 9-3i, 4+4i)$
- (k) $V = \mathbb{C}^4$, $S = \{(-4, 3-2i, i, 1-4i), (-1-5i, 5-4i, -3+5i, 7-2i), (-27-i, -7-6i, -15+25i, -7-6i)\}$, and $x = (-13-7i, -12+3i, -39-11i, -26+5i)$
- (l) $V = M_{2 \times 2}(C)$, $S = \left\{ \begin{pmatrix} 1-i & -2-3i \\ 2+2i & 4+i \end{pmatrix}, \begin{pmatrix} 8i & 4 \\ -3-3i & -4+4i \end{pmatrix}, \begin{pmatrix} -25-38i & -2-13i \\ 12-78i & -7+24i \end{pmatrix} \right\}$, and $A = \begin{pmatrix} -2+8i & -13+i \\ 10-10i & 9-9i \end{pmatrix}$
- (m) $V = M_{2 \times 2}(C)$, $S = \left\{ \begin{pmatrix} -1+i & -i \\ 2-i & 1+3i \end{pmatrix}, \begin{pmatrix} -1-7i & -9-8i \\ 1+10i & -6-2i \end{pmatrix}, \begin{pmatrix} -11-132i & -34-31i \\ 7-126i & -71-5i \end{pmatrix} \right\}$, and $A = \begin{pmatrix} -7+5i & 3+18i \\ 9-6i & -3+7i \end{pmatrix}$

3. In \mathbb{R}^2 , let

$$\beta = \left\{ \left(\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}} \right), \left(\frac{1}{\sqrt{2}}, \frac{-1}{\sqrt{2}} \right) \right\}.$$

Find the Fourier coefficients of $(3, 4)$ relative to β .

4. Let $S = \{(1, 0, i), (1, 2, 1)\}$ in \mathbb{C}^3 . Compute S^\perp .
5. Let $S_0 = \{x_0\}$, where x_0 is a nonzero vector in \mathbb{R}^3 . Describe S_0^\perp geometrically. Now suppose that $S = \{x_1, x_2\}$ is a linearly independent subset of \mathbb{R}^3 . Describe S^\perp geometrically.

6. Let V be an inner product space, and let W be a finite-dimensional subspace of V . If $x \notin W$, prove that there exists $y \in V$ such that $y \in W^\perp$, but $\langle x, y \rangle \neq 0$. *Hint:* Use Theorem 6.6.
7. Let β be a basis for a subspace W of an inner product space V , and let $z \in V$. Prove that $z \in W^\perp$ if and only if $\langle z, v \rangle = 0$ for every $v \in \beta$.
8. Prove that if $\{w_1, w_2, \dots, w_n\}$ is an orthogonal set of nonzero vectors, then the vectors v_1, v_2, \dots, v_n derived from the Gram–Schmidt process satisfy $v_i = w_i$ for $i = 1, 2, \dots, n$. *Hint:* Use mathematical induction.
9. Let $W = \text{span}(\{(i, 0, 1)\})$ in C^3 . Find orthonormal bases for W and W^\perp .
10. Let W be a finite-dimensional subspace of an inner product space V . Prove that $V = W \oplus W^\perp$. Using the definition on page 76, prove that there exists a projection T on W along W^\perp that satisfies $N(T) = W^\perp$. In addition, prove that $\|T(x)\| \leq \|x\|$ for all $x \in V$. *Hint:* Use Theorem 6.6 and Exercise 10 of Section 6.1.
11. Let A be an $n \times n$ matrix with complex entries. Prove that $AA^* = I$ if and only if the rows of A form an orthonormal basis for C^n . Visit goo.gl/iKcC4S for a solution.
12. Prove that for any matrix $A \in M_{m \times n}(F)$, $(R(L_{A^*}))^\perp = N(L_A)$.
13. Let V be an inner product space, S and S_0 be subsets of V , and W be a finite-dimensional subspace of V . Prove the following results.
 - (a) $S_0 \subseteq S$ implies that $S^\perp \subseteq S_0^\perp$.
 - (b) $S \subseteq (S^\perp)^\perp$; so $\text{span}(S) \subseteq (S^\perp)^\perp$.
 - (c) $W = (W^\perp)^\perp$. *Hint:* Use Exercise 6.
 - (d) $V = W \oplus W^\perp$. (See the exercises of Section 1.3.)
14. Let W_1 and W_2 be subspaces of a finite-dimensional inner product space. Prove that $(W_1 + W_2)^\perp = W_1^\perp \cap W_2^\perp$ and $(W_1 \cap W_2)^\perp = W_1^\perp + W_2^\perp$. (See the definition of the sum of subsets of a vector space on page 22.) *Hint for the second equation:* Apply Exercise 13(c) to the first equation.
15. Let V be a finite-dimensional inner product space over F .
 - (a) *Parseval's Identity.* Let $\{v_1, v_2, \dots, v_n\}$ be an orthonormal basis for V . For any $x, y \in V$ prove that

$$\langle x, y \rangle = \sum_{i=1}^n \langle x, v_i \rangle \overline{\langle y, v_i \rangle}.$$

- (b) Use (a) to prove that if β is an orthonormal basis for V with inner product $\langle \cdot, \cdot \rangle$, then for any $x, y \in V$

$$\langle \phi_\beta(x), \phi_\beta(y) \rangle' = \langle [x]_\beta, [y]_\beta \rangle' = \langle x, y \rangle,$$

where $\langle \cdot, \cdot \rangle'$ is the standard inner product on \mathbb{F}^n .

- 16. (a)** *Bessel's Inequality.* Let V be an inner product space, and let $S = \{v_1, v_2, \dots, v_n\}$ be an orthonormal subset of V . Prove that for any $x \in V$ we have

$$\|x\|^2 \geq \sum_{i=1}^n |\langle x, v_i \rangle|^2.$$

Hint: Apply Theorem 6.6 to $x \in V$ and $W = \text{span}(S)$. Then use Exercise 10 of Section 6.1.

- (b)** In the context of (a), prove that Bessel's inequality is an equality if and only if $x \in \text{span}(S)$.
- 17.** Let T be a linear operator on an inner product space V . If $\langle T(x), y \rangle = 0$ for all $x, y \in V$, prove that $T = T_0$. In fact, prove this result if the equality holds for all x and y in some basis for V .

- 18.** Let $V = C([-1, 1])$. Suppose that W_e and W_o denote the subspaces of V consisting of the even and odd functions, respectively. (See Exercise 22 of Section 1.3.) Prove that $W_e^\perp = W_o$, where the inner product on V is defined by

$$\langle f, g \rangle = \int_{-1}^1 f(t)g(t) dt.$$

- 19.** In each of the following parts, find the orthogonal projection of the given vector on the given subspace W of the inner product space V .
- (a)** $V = \mathbb{R}^2$, $u = (2, 6)$, and $W = \{(x, y) : y = 4x\}$.
 - (b)** $V = \mathbb{R}^3$, $u = (2, 1, 3)$, and $W = \{(x, y, z) : x + 3y - 2z = 0\}$.
 - (c)** $V = P(R)$ with the inner product $\langle f(x), g(x) \rangle = \int_0^1 f(t)g(t) dt$, $h(x) = 4 + 3x - 2x^2$, and $W = P_1(R)$.
- 20.** In each part of Exercise 19, find the distance from the given vector to the subspace W .
- 21.** Let $V = C([-1, 1])$ with the inner product $\langle f, g \rangle = \int_{-1}^1 f(t)g(t) dt$, and let W be the subspace $P_2(R)$, viewed as a space of functions. Use the orthonormal basis obtained in Example 5 to compute the “best” (closest) second-degree polynomial approximation of the function $h(t) = e^t$ on the interval $[-1, 1]$.
- 22.** Let $V = C([0, 1])$ with the inner product $\langle f, g \rangle = \int_0^1 f(t)g(t) dt$. Let W be the subspace spanned by the linearly independent set $\{t, \sqrt{t}\}$.
- (a)** Find an orthonormal basis for W .
 - (b)** Let $h(t) = t^2$. Use the orthonormal basis obtained in (a) to obtain the “best” (closest) approximation of h in W .

23. Let V be the vector space defined in Example 5 of Section 1.2, the space of all sequences σ in F (where $F = R$ or $F = C$) such that $\sigma(n) \neq 0$ for only finitely many positive integers n . For $\sigma, \mu \in V$, we define $\langle \sigma, \mu \rangle = \sum_{n=1}^{\infty} \sigma(n)\overline{\mu(n)}$. Since all but a finite number of terms of the series are zero, the series converges.
- Prove that $\langle \cdot, \cdot \rangle$ is an inner product on V , and hence V is an inner product space.
 - For each positive integer n , let e_n be the sequence defined by $e_n(k) = \delta_{nk}$, where δ_{nk} is the Kronecker delta. Prove that $\{e_1, e_2, \dots\}$ is an orthonormal basis for V .
 - Let $\sigma_n = e_1 + e_n$ and $W = \text{span}(\{\sigma_n : n \geq 2\})$.
 - Prove that $e_1 \notin W$, so $W \neq V$.
 - Prove that $W^\perp = \{0\}$, and conclude that $W \neq (W^\perp)^\perp$.

Thus the assumption in Exercise 13(c) that W is finite-dimensional is essential.

6.3 THE ADJOINT OF A LINEAR OPERATOR

In Section 6.1, we defined the conjugate transpose A^* of a matrix A . For a linear operator T on an inner product space V , we now define a related linear operator on V called the *adjoint* of T , whose matrix representation with respect to any orthonormal basis β for V is $[T]_\beta^*$. The analogy between conjugation of complex numbers and adjoints of linear operators will become apparent. We first need a preliminary result.

Let V be an inner product space, and let $y \in V$. The function $g: V \rightarrow F$ defined by $g(x) = \langle x, y \rangle$ is clearly linear. More interesting is the fact that if V is finite-dimensional, every linear transformation from V into F is of this form.

Theorem 6.8. *Let V be a finite-dimensional inner product space over F , and let $g: V \rightarrow F$ be a linear transformation. Then there exists a unique vector $y \in V$ such that $g(x) = \langle x, y \rangle$ for all $x \in V$.*

Proof. Let $\beta = \{v_1, v_2, \dots, v_n\}$ be an orthonormal basis for V , and let

$$y = \sum_{i=1}^n \overline{g(v_i)} v_i.$$

Define $h: V \rightarrow F$ by $h(x) = \langle x, y \rangle$, which is clearly linear. Furthermore, for $1 \leq j \leq n$ we have

$$h(v_j) = \langle v_j, y \rangle = \left\langle v_j, \sum_{i=1}^n \overline{g(v_i)} v_i \right\rangle$$

$$= \sum_{i=1}^n g(v_i) \langle v_j, v_i \rangle = \sum_{i=1}^n g(v_i) \delta_{ji} = g(v_j).$$

Since g and h agree on β , we have that $g = h$ by the corollary to Theorem 2.6 (p. 73).

To show that y is unique, suppose that $g(x) = \langle x, y' \rangle$ for all x . Then $\langle x, y \rangle = \langle x, y' \rangle$ for all x ; so by Theorem 6.1(e) (p. 331), we have $y = y'$. ■

Example 1

Define $g: \mathbb{R}^2 \rightarrow \mathbb{R}$ by $g(a_1, a_2) = 2a_1 + a_2$; clearly g is a linear transformation. Let $\beta = \{e_1, e_2\}$, and let $y = g(e_1)e_1 + g(e_2)e_2 = 2e_1 + e_2 = (2, 1)$, as in the proof of Theorem 6.8. Then $g(a_1, a_2) = \langle (a_1, a_2), (2, 1) \rangle = 2a_1 + a_2$. ♦

Theorem 6.9. *Let V be a finite-dimensional inner product space, and let T be a linear operator on V . Then there exists a unique function $T^*: V \rightarrow V$ such that $\langle T(x), y \rangle = \langle x, T^*(y) \rangle$ for all $x, y \in V$. Furthermore, T^* is linear.*

Proof. Let $y \in V$. Define $g: V \rightarrow F$ by $g(x) = \langle T(x), y \rangle$ for all $x \in V$. We first show that g is linear. Let $x_1, x_2 \in V$ and $c \in F$. Then

$$\begin{aligned} g(cx_1 + x_2) &= \langle T(cx_1 + x_2), y \rangle = \langle cT(x_1) + T(x_2), y \rangle \\ &= c \langle T(x_1), y \rangle + \langle T(x_2), y \rangle = cg(x_1) + g(x_2). \end{aligned}$$

Hence g is linear.

We now apply Theorem 6.8 to obtain a unique vector $y' \in V$ such that $g(x) = \langle x, y' \rangle$; that is, $\langle T(x), y \rangle = \langle x, y' \rangle$ for all $x \in V$. Defining $T^*: V \rightarrow V$ by $T^*(y) = y'$, we have $\langle T(x), y \rangle = \langle x, T^*(y) \rangle$.

To show that T^* is linear, let $y_1, y_2 \in V$ and $c \in F$. Then for any $x \in V$, we have

$$\begin{aligned} \langle x, T^*(cy_1 + y_2) \rangle &= \langle T(x), cy_1 + y_2 \rangle \\ &= \bar{c} \langle T(x), y_1 \rangle + \langle T(x), y_2 \rangle \\ &= \bar{c} \langle x, T^*(y_1) \rangle + \langle x, T^*(y_2) \rangle \\ &= \langle x, cT^*(y_1) + T^*(y_2) \rangle. \end{aligned}$$

Since x is arbitrary, $T^*(cy_1 + y_2) = cT^*(y_1) + T^*(y_2)$ by Theorem 6.1(e) (p. 331).

Finally, we need to show that T^* is unique. Suppose that $U: V \rightarrow V$ is linear and that it satisfies $\langle T(x), y \rangle = \langle x, U(y) \rangle$ for all $x, y \in V$. Then $\langle x, T^*(y) \rangle = \langle x, U(y) \rangle$ for all $x, y \in V$, so $T^* = U$. ■

The linear operator T^* described in Theorem 6.9 is called the **adjoint** of the operator T . The symbol T^* is read “ T star.”

Thus T^* is the unique operator on V satisfying $\langle T(x), y \rangle = \langle x, T^*(y) \rangle$ for all $x, y \in V$. Note that we also have

$$\langle x, T(y) \rangle = \overline{\langle T(y), x \rangle} = \overline{\langle y, T^*(x) \rangle} = \langle T^*(x), y \rangle;$$

so $\langle x, T(y) \rangle = \langle T^*(x), y \rangle$ for all $x, y \in V$. We may view these equations symbolically as adding a * to T when shifting its position inside the inner product symbol.

For an infinite-dimensional inner product space, the adjoint of a linear operator T may be defined to be the function T^* such that $\langle T(x), y \rangle = \langle x, T^*(y) \rangle$ for all $x, y \in V$, provided it exists. Although the uniqueness and linearity of T^* follow as before, the existence of the adjoint is not guaranteed (see Exercise 24). The reader should observe the necessity of the hypothesis of finite-dimensionality in the proof of Theorem 6.8. Many of the theorems we prove about adjoints, nevertheless, do not depend on V being finite-dimensional.

Theorem 6.10 is a useful result for computing adjoints.

Theorem 6.10. *Let V be a finite-dimensional inner product space, and let β be an orthonormal basis for V . If T is a linear operator on V , then*

$$[T^*]_\beta = [T]_\beta^*.$$

Proof. Let $A = [T]_\beta$, $B = [T^*]_\beta$, and $\beta = \{v_1, v_2, \dots, v_n\}$. Then from the corollary to Theorem 6.5 (p. 344), we have

$$B_{ij} = \langle T^*(v_j), v_i \rangle = \overline{\langle v_i, T^*(v_j) \rangle} = \overline{\langle T(v_i), v_j \rangle} = \overline{A}_{ji} = (A^*)_{ij}.$$

Hence $B = A^*$. ■

Corollary. *Let A be an $n \times n$ matrix. Then $L_{A^*} = (L_A)^*$.*

Proof. If β is the standard ordered basis for F^n , then, by Theorem 2.16 (p. 94), we have $[L_A]_\beta = A$. Hence $[(L_A)^*]_\beta = [L_A]_\beta^* = A^* = [L_{A^*}]_\beta$, and so $(L_A)^* = L_{A^*}$. ■

As an illustration of Theorem 6.10, we compute the adjoint of a specific linear operator.

Example 2

Let T be the linear operator on C^2 defined by $T(a_1, a_2) = (2ia_1 + 3a_2, a_1 - a_2)$. If β is the standard ordered basis for C^2 , then

$$[T]_\beta = \begin{pmatrix} 2i & 3 \\ 1 & -1 \end{pmatrix}.$$

So

$$[\mathbf{T}^*]_{\beta} = [\mathbf{T}]_{\beta}^* = \begin{pmatrix} -2i & 1 \\ 3 & -1 \end{pmatrix}.$$

Hence

$$\mathbf{T}^*(a_1, a_2) = (-2ia_1 + a_2, 3a_1 - a_2). \quad \blacklozenge$$

The following theorem suggests an analogy between the conjugates of complex numbers and the adjoints of linear operators.

Theorem 6.11. *Let V be an inner product space, and let T and U be linear operators on V whose adjoints exist. Then*

- (a) $T + U$ has an adjoint, and $(T + U)^* = T^* + U^*$.
- (b) cT has an adjoint, and $(cT)^* = \bar{c}T^*$ for any $c \in F$.
- (c) TU has an adjoint, and $(TU)^* = U^*T^*$.
- (d) T^* has an adjoint, and $T^{**} = T$.
- (e) I has an adjoint, and $I^* = I$.

Proof. We prove (a) and (d); the rest are proved similarly. Let $x, y \in V$.

(a) Because

$$\begin{aligned} \langle (T + U)(x), y \rangle &= \langle T(x) + U(x), y \rangle \\ &= \langle x, T^*(y) \rangle + \langle x, U^*(y) \rangle \\ &= \langle x, T^*(y) + U^*(y) \rangle = \langle x, (T^* + U^*)(y) \rangle, \end{aligned}$$

it follows that $(T + U)^*$ exists and is equal to $T^* + U^*$.

(d) Similarly, since

$$\langle T^*(x), y \rangle, = \langle x, T(y) \rangle,$$

(d) follows. ■

Unless stated otherwise, for the remainder of this chapter we adopt the convention that a reference to the adjoint of a linear operator on an infinite-dimensional inner product space assumes its existence.

Corollary. *Let A and B be $n \times n$ matrices. Then*

- (a) $(A + B)^* = A^* + B^*$.
- (b) $(cA)^* = \bar{c}A^*$ for all $c \in F$.
- (c) $(AB)^* = B^*A^*$.
- (d) $A^{**} = A$.
- (e) $I^* = I$.

Proof. We prove only (c); the remaining parts can be proved similarly.

Since $L_{(AB)}^* = (L_{AB})^* = (L_A L_B)^* = (L_B)^*(L_A)^* = L_{B^*} L_{A^*} = L_{B^* A^*}$, we have $(AB)^* = B^*A^*$. ■

In the preceding proof, we relied on the corollary to Theorem 6.10. An alternative proof, which holds even for nonsquare matrices, can be given by appealing directly to the definition of the conjugate transpose of a matrix (see Exercise 5).

Least Squares Approximation

Consider the following problem: An experimenter collects data by taking measurements y_1, y_2, \dots, y_m at times t_1, t_2, \dots, t_m , respectively. For example, he or she may be measuring unemployment at various times during some period. Suppose that the data $(t_1, y_1), (t_2, y_2), \dots, (t_m, y_m)$ are plotted as points in the plane. (See Figure 6.3.) From this plot, the experimenter feels that there exists an essentially linear relationship between y and t , say $y = ct + d$, and would like to find the constants c and d so that the line $y = ct + d$ represents the best possible fit to the data collected. One such estimate of fit is to calculate the error E that represents the sum of the squares of the vertical distances from the points to the line; that is,

$$E = \sum_{i=1}^m (y_i - ct_i - d)^2.$$

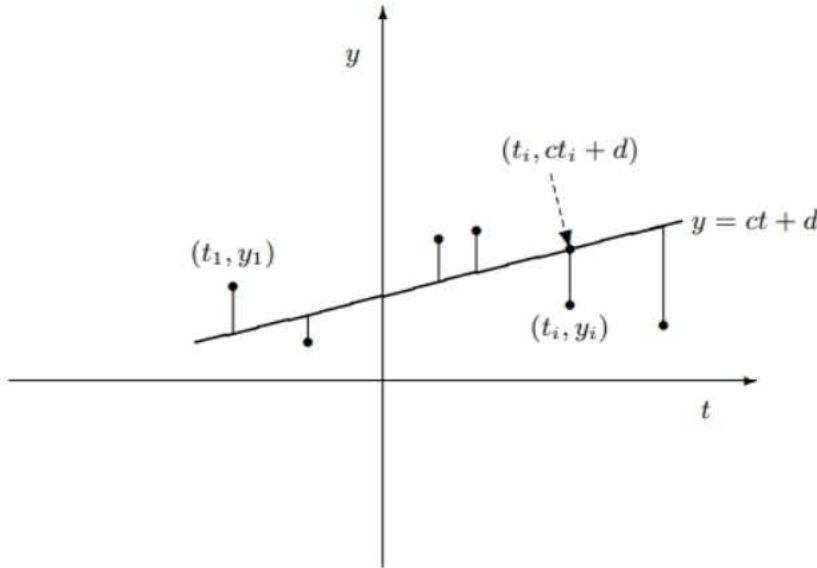


Figure 6.3

Thus the problem is reduced to finding the constants c and d that minimize E . (For this reason the line $y = ct + d$ is called the **least squares line**.) If

we let

$$A = \begin{pmatrix} t_1 & 1 \\ t_2 & 1 \\ \vdots & \vdots \\ t_m & 1 \end{pmatrix}, \quad x = \begin{pmatrix} c \\ d \end{pmatrix}, \quad \text{and} \quad y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix},$$

then it follows that $E = \|y - Ax\|^2$.

We develop a general method for finding an explicit vector $x_0 \in \mathbb{F}^n$ that minimizes E ; that is, given an $m \times n$ matrix A , we find $x_0 \in \mathbb{F}^n$ such that $\|y - Ax_0\| \leq \|y - Ax\|$ for all vectors $x \in \mathbb{F}^n$. This method not only allows us to find the linear function that best fits the data, but also, for any positive integer k , the best fit using a polynomial of degree at most k .

First, we need some notation and two simple lemmas. For $x, y \in \mathbb{F}^n$, let $\langle x, y \rangle_n$ denote the standard inner product of x and y in \mathbb{F}^n . Recall that if x and y are regarded as column vectors, then $\langle x, y \rangle_n = y^*x$.

Lemma 1. Let $A \in \mathbb{M}_{m \times n}(F)$, $x \in \mathbb{F}^n$, and $y \in \mathbb{F}^m$. Then

$$\langle Ax, y \rangle_m = \langle x, A^*y \rangle_n.$$

Proof. By a generalization of the corollary to Theorem 6.11 (see Exercise 5(b)), we have

$$\langle Ax, y \rangle_m = y^*(Ax) = (y^*A)x = (A^*y)^*x = \langle x, A^*y \rangle_n. \quad \blacksquare$$

Lemma 2. Let $A \in \mathbb{M}_{m \times n}(F)$. Then $\text{rank}(A^*A) = \text{rank}(A)$.

Proof. By the dimension theorem, we need only show that, for $x \in \mathbb{F}^n$, we have $A^*Ax = 0$ if and only if $Ax = 0$. Clearly, $Ax = 0$ implies that $A^*Ax = 0$. So assume that $A^*Ax = 0$. Then

$$0 = \langle A^*Ax, x \rangle_n = \langle Ax, A^{**}x \rangle_m = \langle Ax, Ax \rangle_m,$$

so that $Ax = 0$. ■

Corollary. If A is an $m \times n$ matrix such that $\text{rank}(A) = n$, then A^*A is invertible.

Now let A be an $m \times n$ matrix and $y \in \mathbb{F}^m$. Define $W = \{Ax : x \in \mathbb{F}^n\}$; that is, $W = R(L_A)$. By the corollary to Theorem 6.6 (p. 347), there exists a unique vector in W that is closest to y . Call this vector Ax_0 , where $x_0 \in \mathbb{F}^n$. Then $\|Ax_0 - y\| \leq \|Ax - y\|$ for all $x \in \mathbb{F}^n$; so x_0 has the property that $E = \|Ax_0 - y\|^2$ is minimal, as desired.

To develop a practical method for finding such an x_0 , we note from Theorem 6.6 and its corollary that $Ax_0 - y \in W^\perp$; so $\langle Ax, Ax_0 - y \rangle_m = 0$ for

all $x \in \mathbb{F}^n$. Thus, by Lemma 1, we have that $\langle x, A^*(Ax_0 - y) \rangle_n = 0$ for all $x \in \mathbb{F}^n$; that is, $A^*(Ax_0 - y) = 0$. So we need only find a solution x_0 to $A^*Ax = A^*y$. If, in addition, we assume that $\text{rank}(A) = n$, then by Lemma 2 we have $x_0 = (A^*A)^{-1}A^*y$. We summarize this discussion in the following theorem.

Theorem 6.12. *Let $A \in M_{m \times n}(F)$ and $y \in \mathbb{F}^m$. Then there exists $x_0 \in \mathbb{F}^n$ such that $(A^*A)x_0 = A^*y$ and $\|Ax_0 - y\| \leq \|Ax - y\|$ for all $x \in \mathbb{F}^n$. Furthermore, if $\text{rank}(A) = n$, then $x_0 = (A^*A)^{-1}A^*y$.*

To return to our experimenter, let us suppose that the data collected are $(1, 2), (2, 3), (3, 5)$, and $(4, 7)$. Then

$$A = \begin{pmatrix} 1 & 1 \\ 2 & 1 \\ 3 & 1 \\ 4 & 1 \end{pmatrix} \quad \text{and} \quad y = \begin{pmatrix} 2 \\ 3 \\ 5 \\ 7 \end{pmatrix};$$

hence

$$A^*A = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 1 & 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 2 & 1 \\ 3 & 1 \\ 4 & 1 \end{pmatrix} = \begin{pmatrix} 30 & 10 \\ 10 & 4 \end{pmatrix}.$$

Thus

$$(A^*A)^{-1} = \frac{1}{20} \begin{pmatrix} 4 & -10 \\ -10 & 30 \end{pmatrix}.$$

Therefore

$$\begin{pmatrix} c \\ d \end{pmatrix} = x_0 = \frac{1}{20} \begin{pmatrix} 4 & -10 \\ -10 & 30 \end{pmatrix} \begin{pmatrix} 1 & 2 & 3 & 4 \\ 1 & 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} 2 \\ 3 \\ 5 \\ 7 \end{pmatrix} = \begin{pmatrix} 1.7 \\ 0 \end{pmatrix}.$$

It follows that the line $y = 1.7t$ is the least squares line. The error E may be computed directly as $\|Ax_0 - y\|^2 = 0.3$.

Suppose that the experimenter chose the times t_i ($1 \leq i \leq m$) to satisfy

$$\sum_{i=1}^m t_i = 0.$$

Then the two columns of A would be orthogonal, so A^*A would be a diagonal matrix (see Exercise 19). In this case, the computations are greatly simplified.

In practice, the $m \times 2$ matrix A in our least squares application has rank equal to two, and hence A^*A is invertible by the corollary to Lemma 2. For, otherwise, the first column of A is a multiple of the second column, which

consists only of ones. But this would occur only if the experimenter collects all the data at exactly one time.

Finally, the method above may also be applied if, for some k , the experimenter wants to fit a polynomial of degree at most k to the data. For instance, if a polynomial $y = ct^2 + dt + e$ of degree at most 2 is desired, the appropriate model is

$$x = \begin{pmatrix} c \\ d \\ e \end{pmatrix}, \quad y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix}, \quad \text{and} \quad A = \begin{pmatrix} t_1^2 & t_1 & 1 \\ \vdots & \vdots & \vdots \\ t_m^2 & t_m & 1 \end{pmatrix}.$$

Minimal Solutions to Systems of Linear Equations

Even when a system of linear equations $Ax = b$ is consistent, there may be no unique solution. In such cases, it may be desirable to find a solution of minimal norm. A solution s to $Ax = b$ is called a **minimal solution** if $\|s\| \leq \|u\|$ for all other solutions u . The next theorem assures that every consistent system of linear equations has a unique minimal solution and provides a method for computing it.

Theorem 6.13. *Let $A \in M_{m \times n}(F)$ and $b \in F^m$. Suppose that $Ax = b$ is consistent. Then the following statements are true.*

- (a) *There exists exactly one minimal solution s of $Ax = b$, and $s \in R(L_{A^*})$.*
- (b) *The vector s is the only solution to $Ax = b$ that lies in $R(L_{A^*})$; in fact, if u satisfies $(AA^*)u = b$, then $s = A^*u$.*

Proof. (a) For simplicity of notation, we let $W = R(L_{A^*})$ and $W' = N(L_A)$. Let x be any solution to $Ax = b$. By Theorem 6.6 (p. 347), $x = s + y$ for some $s \in W$ and $y \in W^\perp$. But $W^\perp = W'$ by Exercise 12, and therefore $b = Ax = As + Ay = As$. So s is a solution to $Ax = b$ that lies in W . To prove (a), we need only show that s is the unique minimal solution. Let v be any solution to $Ax = b$. By Theorem 3.9 (p. 172), we have that $v = s + u$, where $u \in W'$. Since $s \in W$, which equals W'^\perp by Exercise 12, we have

$$\|v\|^2 = \|s + u\|^2 = \|s\|^2 + \|u\|^2 \geq \|s\|^2$$

by Exercise 10 of Section 6.1. Thus s is a minimal solution. We can also see from the preceding calculation that if $\|v\| = \|s\|$, then $u = 0$; hence $v = s$. Therefore s is the unique minimal solution to $Ax = b$, proving (a).

- (b) Assume that v is also a solution to $Ax = b$ that lies in W . Then

$$v - s \in W \cap W' = W \cap W^\perp = \{0\};$$

so $v = s$.

Finally, suppose that $(AA^*)u = b$, and let $v = A^*u$. Then $v \in W$ and $Av = b$. Therefore $s = v = A^*u$ by the discussion above. ■

Example 3

Consider the system

$$\begin{aligned}x + 2y + z &= 4 \\x - y + 2z &= -11 \\x + 5y &= 19.\end{aligned}$$

Let

$$A = \begin{pmatrix} 1 & 2 & 1 \\ 1 & -1 & 2 \\ 1 & 5 & 0 \end{pmatrix} \quad \text{and} \quad b = \begin{pmatrix} 4 \\ -11 \\ 19 \end{pmatrix}.$$

To find the minimal solution to this system, we must first find some solution u to $AA^*x = b$. Now

$$AA^* = \begin{pmatrix} 6 & 1 & 11 \\ 1 & 6 & -4 \\ 11 & -4 & 26 \end{pmatrix};$$

so we consider the system

$$\begin{aligned}6x + y + 11z &= 4 \\x + 6y - 4z &= -11 \\11x - 4y + 26z &= 19,\end{aligned}$$

for which one solution is

$$u = \begin{pmatrix} 1 \\ -2 \\ 0 \end{pmatrix}.$$

(Any solution will suffice.) Hence

$$s = A^*u = \begin{pmatrix} -1 \\ 4 \\ -3 \end{pmatrix}$$

is the minimal solution to the given system. ◆

EXERCISES

- Label the following statements as true or false. Assume that the underlying inner product spaces are finite-dimensional.
 - Every linear operator has an adjoint.
 - Every linear operator on V has the form $x \rightarrow \langle x, y \rangle$ for some $y \in V$.
 - For every linear operator T on V and every ordered basis β for V , we have $[T^*]_\beta = ([T]_\beta)^*$.
 - The adjoint of a linear operator is unique.

- (e) For any linear operators T and U and scalars a and b ,

$$(aT + bU)^* = aT^* + bU^*.$$

- (f) For any $n \times n$ matrix A , we have $(L_A)^* = L_{A^*}$.

- (g) For any linear operator T , we have $(T^*)^* = T$.

2. For each of the following inner product spaces V (over F) and linear transformations $g: V \rightarrow F$, find a vector y such that $g(x) = \langle x, y \rangle$ for all $x \in V$.
- (a) $V = \mathbb{R}^3$, $g(a_1, a_2, a_3) = a_1 - 2a_2 + 4a_3$
 - (b) $V = \mathbb{C}^2$, $g(z_1, z_2) = z_1 - 2z_2$
 - (c) $V = P_2(R)$ with $\langle f(x), h(x) \rangle = \int_0^1 f(t)h(t) dt$, $g(f) = f(0) + f'(1)$
3. For each of the following inner product spaces V and linear operators T on V , evaluate T^* at the given vector in V .
- (a) $V = \mathbb{R}^2$, $T(a, b) = (2a + b, a - 3b)$, $x = (3, 5)$.
 - (b) $V = \mathbb{C}^2$, $T(z_1, z_2) = (2z_1 + iz_2, (1 - i)z_1)$, $x = (3 - i, 1 + 2i)$.
 - (c) $V = P_1(R)$ with $\langle f(x), g(x) \rangle = \int_{-1}^1 f(t)g(t) dt$, $T(f) = f' + 3f$,
 $f(t) = 4 - 2t$
4. Complete the proof of Theorem 6.11.
5. (a) Complete the proof of the corollary to Theorem 6.11 by using Theorem 6.11, as in the proof of (c).
(b) State a result for nonsquare matrices that is analogous to the corollary to Theorem 6.11, and prove it using a matrix argument.
6. Let T be a linear operator on an inner product space V . Let $U_1 = T + T^*$ and $U_2 = TT^*$. Prove that $U_1 = U_1^*$ and $U_2 = U_2^*$.
7. Give an example of a linear operator T on an inner product space V such that $N(T) \neq N(T^*)$.
8. Let V be a finite-dimensional inner product space, and let T be a linear operator on V . Prove that if T is invertible, then T^* is invertible and $(T^*)^{-1} = (T^{-1})^*$.
9. Prove that if $V = W \oplus W^\perp$ and T is the projection on W along W^\perp , then $T = T^*$. Hint: Recall that $N(T) = W^\perp$. (For definitions, see the exercises of Sections 1.3 and 2.1.)
10. Let T be a linear operator on an inner product space V . Prove that $\|T(x)\| = \|x\|$ for all $x \in V$ if and only if $\langle T(x), T(y) \rangle = \langle x, y \rangle$ for all $x, y \in V$. Hint: Use Exercise 20 of Section 6.1.

11. For a linear operator T on an inner product space V , prove that $T^*T = T_0$ implies $T = T_0$. Is the same result true if we assume that $TT^* = T_0$?
12. Let V be an inner product space, and let T be a linear operator on V . Prove the following results.
 - (a) $R(T^*)^\perp = N(T)$.
 - (b) If V is finite-dimensional, then $R(T^*) = N(T)^\perp$. Hint: Use Exercise 13(c) of Section 6.2.
13. Let T be a linear operator on a finite-dimensional inner product space V . Prove the following results.
 - (a) $N(T^*T) = N(T)$. Deduce that $\text{rank}(T^*T) = \text{rank}(T)$.
 - (b) $\text{rank}(T) = \text{rank}(T^*)$. Deduce from (a) that $\text{rank}(TT^*) = \text{rank}(T)$.
 - (c) For any $n \times n$ matrix A , $\text{rank}(A^*A) = \text{rank}(AA^*) = \text{rank}(A)$.
14. Let V be an inner product space, and let $y, z \in V$. Define $T: V \rightarrow V$ by $T(x) = \langle x, y \rangle z$ for all $x \in V$. First prove that T is linear. Then show that T^* exists, and find an explicit expression for it.

The following definition is used in Exercises 15–17 and is an extension of the definition of the *adjoint* of a linear operator.

Definition. Let $T: V \rightarrow W$ be a linear transformation, where V and W are finite-dimensional inner product spaces with inner products $\langle \cdot, \cdot \rangle_1$ and $\langle \cdot, \cdot \rangle_2$, respectively. A function $T^*: W \rightarrow V$ is called an **adjoint** of T if $\langle T(x), y \rangle_2 = \langle x, T^*(y) \rangle_1$ for all $x \in V$ and $y \in W$.

15. Let $T: V \rightarrow W$ be a linear transformation, where V and W are finite-dimensional inner product spaces with inner products $\langle \cdot, \cdot \rangle_1$ and $\langle \cdot, \cdot \rangle_2$, respectively. Prove the following results.
 - (a) There is a unique adjoint T^* of T , and T^* is linear.
 - (b) If β and γ are orthonormal bases for V and W , respectively, then $[T^*]_\gamma^\beta = ([T]_\beta^\gamma)^*$.
 - (c) $\text{rank}(T^*) = \text{rank}(T)$.
 - (d) $\langle T^*(x), y \rangle_1 = \langle x, T(y) \rangle_2$ for all $x \in W$ and $y \in V$.
 - (e) For all $x \in V$, $T^*T(x) = 0$ if and only if $T(x) = 0$.
16. State and prove a result that extends the first four parts of Theorem 6.11 using the preceding definition.
17. Let $T: V \rightarrow W$ be a linear transformation, where V and W are finite-dimensional inner product spaces. Prove that $(R(T^*))^\perp = N(T)$, using the preceding definition.

18.[†] Let A be an $n \times n$ matrix. Prove that $\det(A^*) = \overline{\det(A)}$. Visit goo.gl/csqoFY for a solution.

19. Suppose that A is an $m \times n$ matrix in which no two columns are identical. Prove that A^*A is a diagonal matrix if and only if every pair of columns of A is orthogonal.
20. For each of the sets of data that follows, use the least squares approximation to find the best fits with both (i) a linear function and (ii) a quadratic function. Compute the error E in both cases.
- $\{(-3, 9), (-2, 6), (0, 2), (1, 1)\}$
 - $\{(1, 2), (3, 4), (5, 7), (7, 9), (9, 12)\}$
 - $\{(-2, 4), (-1, 3), (0, 1), (1, -1), (2, -3)\}$
21. In physics, *Hooke's law* states that (within certain limits) there is a linear relationship between the length x of a spring and the force y applied to (or exerted by) the spring. That is, $y = cx + d$, where c is called the **spring constant**. Use the following data to estimate the spring constant (the length is given in inches and the force is given in pounds).

Length		Force
x		y
3.5		1.0
4.0		2.2
4.5		2.8
5.0		4.3

22. Find the minimal solution to each of the following systems of linear equations.

$$(a) \quad x + 2y - z = 12$$

$$(b) \quad \begin{aligned} x + 2y - z &= 1 \\ 2x + 3y + z &= 2 \\ 4x + 7y - z &= 4 \end{aligned}$$

$$(c) \quad \begin{aligned} x + y - z &= 0 \\ 2x - y + z &= 3 \\ x - y + z &= 2 \end{aligned}$$

$$(d) \quad \begin{aligned} x + y + z - w &= 1 \\ 2x - y &+ w = 1 \end{aligned}$$

23. Consider the problem of finding the least squares line $y = ct + d$ corresponding to the m observations $(t_1, y_1), (t_2, y_2), \dots, (t_m, y_m)$.
- (a) Show that the equation $(A^*A)x_0 = A^*y$ of Theorem 6.12 takes the form of the *normal equations*:

$$\left(\sum_{i=1}^m t_i^2 \right) c + \left(\sum_{i=1}^m t_i \right) d = \sum_{i=1}^m t_i y_i$$

and

$$\left(\sum_{i=1}^m t_i \right) c + md = \sum_{i=1}^m y_i.$$

These equations may also be obtained from the error E by setting the partial derivatives of E with respect to both c and d equal to zero.

- (b) Use the second normal equation of (a) to show that the least squares line must pass through the *center of mass*, (\bar{t}, \bar{y}) , where

$$\bar{t} = \frac{1}{m} \sum_{i=1}^m t_i \quad \text{and} \quad \bar{y} = \frac{1}{m} \sum_{i=1}^m y_i.$$

24. Let V and $\{e_1, e_2, \dots\}$ be defined as in Exercise 23 of Section 6.2. Define $T: V \rightarrow V$ by

$$T(\sigma)(k) = \sum_{i=k}^{\infty} \sigma(i) \quad \text{for every positive integer } k.$$

Notice that the infinite series in the definition of T converges because $\sigma(i) \neq 0$ for only finitely many i .

- (a) Prove that T is a linear operator on V .
- (b) Prove that for any positive integer n , $T(e_n) = \sum_{i=1}^n e_i$.
- (c) Prove that T has no adjoint. *Hint:* By way of contradiction, suppose that T^* exists. Prove that for any positive integer n , $T^*(e_n)(k) \neq 0$ for infinitely many k .

6.4 NORMAL AND SELF-ADJOINT OPERATORS

We have seen the importance of diagonalizable operators in Chapter 5. For an operator on a vector space V to be diagonalizable, it is necessary and sufficient for V to contain a basis of eigenvectors for this operator. As V is an inner product space in this chapter, it is reasonable to seek conditions that guarantee that V has an orthonormal basis of eigenvectors. A very important result that helps achieve our goal is Schur's theorem (Theorem 6.14). The formulation that follows is in terms of linear operators. The next section contains the more familiar matrix form. We begin with a lemma.

Lemma. *Let T be a linear operator on a finite-dimensional inner product space V . If T has an eigenvector, then so does T^* .*

Proof. Suppose that v is an eigenvector of T with corresponding eigenvalue λ . Then for any $x \in V$,

$$0 = \langle \theta, x \rangle = \langle (T - \lambda I)(v), x \rangle = \langle v, (T - \lambda I)^*(x) \rangle = \langle v, (T^* - \bar{\lambda} I)(x) \rangle,$$

and hence v is orthogonal to the range of $T^* - \bar{\lambda}I$. So $T^* - \bar{\lambda}I$ is not onto and hence is not one-to-one. Thus $T^* - \bar{\lambda}I$ has a nonzero null space, and any nonzero vector in this null space is an eigenvector of T^* with corresponding eigenvalue $\bar{\lambda}$. \blacksquare

Recall (see the exercises of Section 2.1 and see Section 5.4) that a subspace W of V is said to be **T -invariant** if $T(W)$ is contained in W . If W is T -invariant, we may define the restriction $T_W: W \rightarrow W$ by $T_W(x) = T(x)$ for all $x \in W$. It is clear that T_W is a linear operator on W . Recall from Section 5.2 that a polynomial is said to **split** if it factors into linear polynomials.

Theorem 6.14 (Schur). *Let T be a linear operator on a finite-dimensional inner product space V . Suppose that the characteristic polynomial of T splits. Then there exists an orthonormal basis γ for V such that the matrix $[T]_\gamma$ is upper triangular.*

Proof. By Exercise 12(a) of Section 5.2, there exists an ordered basis $\beta = \{w_1, w_2, \dots, w_n\}$ for V such that $[T]_\beta$ is upper triangular. Now apply the Gram-Schmidt process to β to obtain an orthogonal basis $\beta' = \{v_1, v_2, \dots, v_n\}$ for V . For each k , $1 \leq k \leq n$, let

$$S_k = \{w_1, w_2, \dots, w_k\} \quad \text{and} \quad S'_k = \{v_1, v_2, \dots, v_k\}.$$

As in the proof of Theorem 6.4, $\text{span}(S_k) = \text{span}(S'_k)$ for all k . By Exercise 12 of Section 2.2, $T(w_k) \in \text{span}(S_k)$ for all k . Hence $T(v_k) \in \text{span}(S'_k)$ for all k , and so $[T]_{\beta'}$ is upper triangular by the same exercise. Finally, let $z_i = \frac{1}{\|v_i\|}v_i$ for all $1 \leq i \leq n$ and $\gamma = \{z_1, z_2, \dots, z_n\}$. Then γ is an orthonormal basis for V , and $[T]_\gamma$ is upper triangular. \blacksquare

We now return to our original goal of finding an orthonormal basis of eigenvectors of a linear operator T on a finite-dimensional inner product space V . Note that if such an orthonormal basis β exists, then $[T]_\beta$ is a diagonal matrix, and hence $[T^*]_\beta = [T]_\beta^*$ is also a diagonal matrix. Because diagonal matrices commute, we conclude that T and T^* commute. Thus if V possesses an orthonormal basis of eigenvectors of T , then $TT^* = T^*T$.

Definitions. *Let V be an inner product space, and let T be a linear operator on V . We say that T is **normal** if $TT^* = T^*T$. An $n \times n$ real or complex matrix A is **normal** if $AA^* = A^*A$.*

It follows immediately from Theorem 6.10 (p. 356) that T is normal if and only if $[T]_\beta$ is normal, where β is an orthonormal basis.

Example 1

Let $T: \mathbb{R}^2 \rightarrow \mathbb{R}^2$ be rotation by θ , where $0 < \theta < \pi$. The matrix representation of T in the standard ordered basis is given by

$$A = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix}.$$

Note that $AA^* = I = A^*A$; so A , and hence T , is normal. \spadesuit

Example 2

Suppose that A is a real *skew-symmetric* matrix; that is, $A^t = -A$. Then A is normal because both AA^t and A^tA are equal to $-A^2$. \clubsuit

Clearly, the operator T in Example 1 does not even possess one eigenvector. So in the case of a real inner product space, we see that normality is not sufficient to guarantee an orthonormal basis of eigenvectors. All is not lost, however. We show that normality suffices if V is a *complex* inner product space.

Before we prove the promised result for normal operators, we need some general properties of normal operators.

Theorem 6.15. *Let V be an inner product space, and let T be a normal operator on V . Then the following statements are true.*

- (a) $\|T(x)\| = \|T^*(x)\|$ for all $x \in V$.
- (b) $T - cl$ is normal for every $c \in F$.
- (c) If x is an eigenvector of T corresponding to eigenvalue λ , then x is also an eigenvector of T^* corresponding to eigenvalue $\bar{\lambda}$. That is, if $T(x) = \lambda x$, then $T^*(x) = \bar{\lambda}x$.
- (d) If λ_1 and λ_2 are distinct eigenvalues of T with corresponding eigenvectors x_1 and x_2 , then x_1 and x_2 are orthogonal.

Proof. (a) For any $x \in V$, we have

$$\begin{aligned} \|T(x)\|^2 &= \langle T(x), T(x) \rangle = \langle T^*T(x), x \rangle = \langle TT^*(x), x \rangle \\ &= \langle T^*(x), T^*(x) \rangle = \|T^*(x)\|^2. \end{aligned}$$

The proof of (b) is left as an exercise.

(c) Suppose that $T(x) = \lambda x$ for some $x \in V$. Let $U = T - \lambda I$. Then $U(x) = 0$, and U is normal by (b). Thus (a) implies that

$$0 = \|U(x)\| = \|U^*(x)\| = \|(T^* - \bar{\lambda}I)(x)\| = \|T^*(x) - \bar{\lambda}x\|.$$

Hence $T^*(x) = \bar{\lambda}x$. So x is an eigenvector of T^* .

(d) Let λ_1 and λ_2 be distinct eigenvalues of T with corresponding eigenvectors x_1 and x_2 . Then, using (c), we have

$$\lambda_1 \langle x_1, x_2 \rangle = \langle \lambda_1 x_1, x_2 \rangle = \langle T(x_1), x_2 \rangle = \langle x_1, T^*(x_2) \rangle$$

$$= \langle x_1, \bar{\lambda}_2 x_2 \rangle = \lambda_2 \langle x_1, x_2 \rangle.$$

Since $\lambda_1 \neq \lambda_2$, we conclude that $\langle x_1, x_2 \rangle = 0$. ■

Theorem 6.16. *Let T be a linear operator on a finite-dimensional complex inner product space V . Then T is normal if and only if there exists an orthonormal basis for V consisting of eigenvectors of T .*

Proof. Suppose that T is normal. By the fundamental theorem of algebra (Theorem D.4), the characteristic polynomial of T splits. So we may apply Schur's theorem to obtain an orthonormal basis $\beta = \{v_1, v_2, \dots, v_n\}$ for V such that $[T]_\beta = A$ is upper triangular. We know that v_1 is an eigenvector of T because A is upper triangular. Assume that v_1, v_2, \dots, v_{k-1} are eigenvectors of T . We claim that v_k is also an eigenvector of T . It then follows by mathematical induction on k that all of the v_i 's are eigenvectors of T . Consider any $j < k$, and let λ_j denote the eigenvalue of T corresponding to v_j . By Theorem 6.15, $T^*(v_j) = \bar{\lambda}_j v_j$. Since A is upper triangular,

$$T(v_k) = A_{1k}v_1 + A_{2k}v_2 + \cdots + A_{jk}v_j + \cdots + A_{kk}v_k.$$

Furthermore, by the corollary to Theorem 6.5 (p. 345),

$$A_{jk} = \langle T(v_k), v_j \rangle = \langle v_k, T^*(v_j) \rangle = \langle v_k, \bar{\lambda}_j v_j \rangle = \lambda_j \langle v_k, v_j \rangle = 0.$$

It follows that $T(v_k) = A_{kk}v_k$, and hence v_k is an eigenvector of T . So by induction, all the vectors in β are eigenvectors of T .

The converse was already proved on page 367. ■

Interestingly, as the next example shows, Theorem 6.16 does not extend to infinite-dimensional complex inner product spaces.

Example 3

Consider the inner product space H with the orthonormal set S from Example 9 in Section 6.1. Let $V = \text{span}(S)$, and let T and U be the linear operators on V defined by $T(f) = f_1 f$ and $U(f) = f_{-1} f$. Then

$$T(f_n) = f_{n+1} \quad \text{and} \quad U(f_n) = f_{n-1}$$

for all integers n . Thus

$$\langle T(f_m), f_n \rangle = \langle f_{m+1}, f_n \rangle = \delta_{(m+1),n} = \delta_{m,(n-1)} = \langle f_m, f_{n-1} \rangle = \langle f_m, U(f_n) \rangle.$$

It follows that $U = T^*$. Furthermore, $TT^* = I = T^*T$; so T is normal.

We show that T has no eigenvectors. Suppose that f is an eigenvector of T , say, $T(f) = \lambda f$ for some λ . Since V equals the span of S , we may write

$$f = \sum_{i=n}^m a_i f_i, \quad \text{where } a_m \neq 0.$$

Hence

$$\sum_{i=n}^m a_i f_{i+1} = T(f) = \lambda f = \sum_{i=n}^m \lambda a_i f_i.$$

Since $a_m \neq 0$, we can write f_{m+1} as a linear combination of f_n, f_{n+1}, \dots, f_m . But this is a contradiction because S is linearly independent. ♦

Example 1 illustrates that normality is not sufficient to guarantee the existence of an orthonormal basis of eigenvectors for real inner product spaces. For real inner product spaces, we must replace normality by the stronger condition that $T = T^*$ in order to guarantee such a basis.

Definitions. Let T be a linear operator on an inner product space V . We say that T is **self-adjoint** (or **Hermitian**) if $T = T^*$. An $n \times n$ real or complex matrix A is **self-adjoint** (or **Hermitian**) if $A = A^*$.

It follows immediately that if β is an orthonormal basis, then T is self-adjoint if and only if $[T]_\beta$ is self-adjoint. For real matrices, this condition reduces to the requirement that A be symmetric.

Before we state our main result for self-adjoint operators, we need some preliminary work.

By definition, a linear operator on a real inner product space has only real eigenvalues. The lemma that follows shows that the same can be said for self-adjoint operators on a complex inner product space. Similarly, the characteristic polynomial of every linear operator on a complex inner product space splits, and the same is true for self-adjoint operators on a real inner product space.

Lemma. Let T be a self-adjoint operator on a finite-dimensional inner product space V . Then

- (a) Every eigenvalue of T is real.
- (b) Suppose that V is a real inner product space. Then the characteristic polynomial of T splits.

Proof. (a) Suppose that $T(x) = \lambda x$ for $x \neq 0$. Because a self-adjoint operator is also normal, we can apply Theorem 6.15(c) to obtain

$$\lambda x = T(x) = T^*(x) = \bar{\lambda}x.$$

So $\lambda = \bar{\lambda}$; that is, λ is real.

(b) Let $n = \dim(V)$, β be an orthonormal basis for V , and $A = [T]_\beta$. Then A is self-adjoint. Let T_A be the linear operator on C^n defined by $T_A(x) = Ax$ for all $x \in C^n$. Note that T_A is self-adjoint because $[T_A]_\gamma = A$, where γ is the standard ordered (orthonormal) basis for C^n . So, by (a), the eigenvalues of T_A are real. By the fundamental theorem of algebra, the characteristic polynomial of T_A splits into factors of the form $t - \lambda$. Since each

λ is real, the characteristic polynomial splits over R . But T_A has the same characteristic polynomial as A , which has the same characteristic polynomial as T . Therefore the characteristic polynomial of T splits. ■

We are now able to establish one of the major results of this chapter.

Theorem 6.17. *Let T be a linear operator on a finite-dimensional real inner product space V . Then T is self-adjoint if and only if there exists an orthonormal basis β for V consisting of eigenvectors of T .*

Proof. Suppose that T is self-adjoint. By the lemma, we may apply Schur's theorem to obtain an orthonormal basis β for V such that the matrix $A = [T]_\beta$ is upper triangular. But

$$A^* = [T]^*_\beta = [T^*]_\beta = [T]_\beta = A.$$

So A and A^* are both upper triangular, and therefore A is a diagonal matrix. Thus β must consist of eigenvectors of T . ■

The converse is left as an exercise. ■

We restate this theorem in matrix form in the next section (as Theorem 6.20 on p. 381).

Example 4

As we noted earlier, real symmetric matrices are self-adjoint, and self-adjoint matrices are normal. The following matrix A is complex and symmetric:

$$A = \begin{pmatrix} i & i \\ i & 1 \end{pmatrix} \quad \text{and} \quad A^* = \begin{pmatrix} -i & -i \\ -i & 1 \end{pmatrix}.$$

But A is not normal, because $(AA^*)_{12} = 1+i$ and $(A^*A)_{12} = 1-i$. Therefore complex symmetric matrices need not be normal. ♦

EXERCISES

1. Label the following statements as true or false. Assume that the underlying inner product spaces are finite-dimensional.
 - (a) Every self-adjoint operator is normal.
 - (b) Operators and their adjoints have the same eigenvectors.
 - (c) If T is an operator on an inner product space V , then T is normal if and only if $[T]_\beta$ is normal, where β is any ordered basis for V .
 - (d) A real or complex matrix A is normal if and only if L_A is normal.
 - (e) The eigenvalues of a self-adjoint operator must all be real.
 - (f) The identity and zero operators are self-adjoint.

- (g) Every normal operator is diagonalizable.
 (h) Every self-adjoint operator is diagonalizable.
2. For each linear operator T on an inner product space V , determine whether T is normal, self-adjoint, or neither. If possible, produce an orthonormal basis of eigenvectors of T for V and list the corresponding eigenvalues.
- (a) $V = \mathbb{R}^2$ and T is defined by $T(a, b) = (2a - 2b, -2a + 5b)$.
 (b) $V = \mathbb{R}^3$ and T is defined by $T(a, b, c) = (-a + b, 5b, 4a - 2b + 5c)$.
 (c) $V = \mathbb{C}^2$ and T is defined by $T(a, b) = (2a + ib, a + 2b)$.
 (d) $V = P_2(R)$ and T is defined by $T(f) = f'$, where
- $$\langle f(x), g(x) \rangle = \int_0^1 f(t)g(t) dt.$$
- (e) $V = M_{2 \times 2}(R)$ and T is defined by $T(A) = A^t$.
 (f) $V = M_{2 \times 2}(R)$ and T is defined by $T \begin{pmatrix} a & b \\ c & d \end{pmatrix} = \begin{pmatrix} c & d \\ a & b \end{pmatrix}$.
3. Give an example of a linear operator T on \mathbb{R}^2 and an ordered basis for \mathbb{R}^2 that provides a counterexample to the statement in Exercise 1(c).
4. Let T and U be self-adjoint operators on an inner product space V . Prove that TU is self-adjoint if and only if $TU = UT$.
5. Prove (b) of Theorem 6.15.
6. Let V be a complex inner product space, and let T be a linear operator on V . Define
- $$T_1 = \frac{1}{2}(T + T^*) \quad \text{and} \quad T_2 = \frac{1}{2i}(T - T^*).$$
- (a) Prove that T_1 and T_2 are self-adjoint and that $T = T_1 + iT_2$.
 (b) Suppose also that $T = U_1 + iU_2$, where U_1 and U_2 are self-adjoint. Prove that $U_1 = T_1$ and $U_2 = T_2$.
 (c) Prove that T is normal if and only if $T_1T_2 = T_2T_1$.
7. Let T be a linear operator on an inner product space V , and let W be a T -invariant subspace of V . Prove the following results.
- (a) If T is self-adjoint, then T_W is self-adjoint.
 (b) W^\perp is T^* -invariant.
 (c) If W is both T - and T^* -invariant, then $(T_W)^* = (T^*)_W$.
 (d) If W is both T - and T^* -invariant and T is normal, then T_W is normal.

8. Let T be a normal operator on a finite-dimensional complex inner product space V , and let W be a subspace of V . Prove that if W is T -invariant, then W is also T^* -invariant. *Hint:* Use Exercise 24 of Section 5.4.
9. Let T be a normal operator on a finite-dimensional inner product space V . Prove that $N(T) = N(T^*)$ and $R(T) = R(T^*)$. *Hint:* Use Theorem 6.15 and Exercise 12 of Section 6.3.
10. Let T be a self-adjoint operator on a finite-dimensional inner product space V . Prove that for all $x \in V$

$$\|T(x) \pm ix\|^2 = \|T(x)\|^2 + \|x\|^2.$$

Deduce that $T - iI$ is invertible and that the adjoint of $(T - iI)^{-1}$ is $(T + iI)^{-1}$.

11. Assume that T is a linear operator on a complex (not necessarily finite-dimensional) inner product space V with an adjoint T^* . Prove the following results.
 - (a) If T is self-adjoint, then $\langle T(x), x \rangle$ is real for all $x \in V$.
 - (b) If T satisfies $\langle T(x), x \rangle = 0$ for all $x \in V$, then $T = T_0$. *Hint:* Replace x by $x + y$ and then by $x + iy$, and expand the resulting inner products.
 - (c) If $\langle T(x), x \rangle$ is real for all $x \in V$, then T is self-adjoint.
12. Let T be a normal operator on a finite-dimensional real inner product space V whose characteristic polynomial splits. Prove that V has an orthonormal basis of eigenvectors of T . Hence prove that T is self-adjoint.
13. An $n \times n$ real matrix A is said to be a **Gramian** matrix if there exists a real (square) matrix B such that $A = B^t B$. Prove that A is a Gramian matrix if and only if A is symmetric and all of its eigenvalues are non-negative. *Hint:* Apply Theorem 6.17 to $T = L_A$ to obtain an orthonormal basis $\{v_1, v_2, \dots, v_n\}$ of eigenvectors with the associated eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_n$. Define the linear operator U by $U(v_i) = \sqrt{\lambda_i} v_i$.
14. *Simultaneous Diagonalization.* Let V be a finite-dimensional real inner product space, and let U and T be self-adjoint linear operators on V such that $UT = TU$. Prove that there exists an orthonormal basis for V consisting of vectors that are eigenvectors of both U and T . (The complex version of this result appears as Exercise 10 of Section 6.6.) *Hint:* For any eigenspace $W = E_\lambda$ of T , we have that W is both T - and U -invariant. By Exercise 7, we have that W^\perp is both T - and U -invariant. Apply Theorem 6.17 and Theorem 6.6 (p. 347).

15. Let A and B be symmetric $n \times n$ matrices such that $AB = BA$. Use Exercise 14 to prove that there exists an orthogonal matrix P such that $P^t AP$ and $P^t BP$ are both diagonal matrices.
16. Prove the *Cayley–Hamilton theorem* for a complex $n \times n$ matrix A . That is, if $f(t)$ is the characteristic polynomial of A , prove that $f(A) = O$. *Hint:* Use Schur's theorem to show that A may be assumed to be upper triangular, in which case

$$f(t) = \prod_{i=1}^n (A_{ii} - t).$$

Now if $T = L_A$, we have $(A_{jj}I - T)(e_j) \in \text{span}(\{e_1, e_2, \dots, e_{j-1}\})$ for $j \geq 2$, where $\{e_1, e_2, \dots, e_n\}$ is the standard ordered basis for \mathbb{C}^n . (The general case is proved in Section 5.4.)

The following definitions are used in Exercises 17 through 23.

Definitions. A linear operator T on a finite-dimensional inner product space is called **positive definite [positive semidefinite]** if T is self-adjoint and $\langle T(x), x \rangle > 0$ [$\langle T(x), x \rangle \geq 0$] for all $x \neq 0$.

An $n \times n$ matrix A with entries from R or C is called **positive definite [positive semidefinite]** if L_A is positive definite [positive semidefinite].

17. Let T and U be self-adjoint linear operators on an n -dimensional inner product space V , and let $A = [T]_\beta$, where β is an orthonormal basis for V . Prove the following results.
- (a) T is positive definite [semidefinite] if and only if all of its eigenvalues are positive [nonnegative].
 - (b) T is positive definite if and only if

$$\sum_{i,j} A_{ij} a_j \bar{a}_i > 0 \text{ for all nonzero } n\text{-tuples } (a_1, a_2, \dots, a_n).$$

- (c) T is positive semidefinite if and only if $A = B^*B$ for some square matrix B .
- (d) If T and U are positive semidefinite operators such that $T^2 = U^2$, then $T = U$.
- (e) If T and U are positive definite operators such that $TU = UT$, then TU is positive definite.
- (f) T is positive definite [semidefinite] if and only if A is positive definite [semidefinite].

Because of (f), results analogous to items (a) through (d) hold for matrices as well as operators.

18. Let $T: V \rightarrow W$ be a linear transformation, where V and W are finite-dimensional inner product spaces. Prove the following results.

- (a) T^*T and TT^* are positive semidefinite. (See Exercise 15 of Section 6.3.)
(b) $\text{rank}(T^*T) = \text{rank}(TT^*) = \text{rank}(T)$.

19. Let T and U be positive definite operators on an inner product space V . Prove the following results.

- (a) $T + U$ is positive definite.
(b) If $c > 0$, then cT is positive definite.
(c) T^{-1} is positive definite.

Visit [goo.gl/cQch7i](#) for a solution.

20. Let V be an inner product space with inner product $\langle \cdot, \cdot \rangle$, and let T be a positive definite linear operator on V . Prove that $\langle x, y \rangle' = \langle T(x), y \rangle$ defines another inner product on V .

21. Let V be a finite-dimensional inner product space, and let T and U be self-adjoint operators on V such that T is positive definite. Prove that both TU and UT are diagonalizable linear operators that have only real eigenvalues. *Hint:* Show that UT is self-adjoint with respect to the inner product $\langle x, y \rangle' = \langle T(x), y \rangle$. To show that TU is self-adjoint, repeat the argument with T^{-1} in place of T .

22. This exercise provides a converse to Exercise 20. Let V be a finite-dimensional inner product space with inner product $\langle \cdot, \cdot \rangle$, and let $\langle \cdot, \cdot \rangle'$ be any other inner product on V .

- (a) Prove that there exists a unique linear operator T on V such that $\langle x, y \rangle' = \langle T(x), y \rangle$ for all x and y in V . *Hint:* Let $\beta = \{v_1, v_2, \dots, v_n\}$ be an orthonormal basis for V with respect to $\langle \cdot, \cdot \rangle$, and define a matrix A by $A_{ij} = \langle v_j, v_i \rangle'$ for all i and j . Let T be the unique linear operator on V such that $[T]_\beta = A$.
(b) Prove that the operator T of (a) is positive definite with respect to both inner products.

23. Let U be a diagonalizable linear operator on a finite-dimensional inner product space V such that all of the eigenvalues of U are real. Prove that there exist positive definite linear operators T_1 and T'_1 and self-adjoint linear operators T_2 and T'_2 such that $U = T_2 T_1 = T'_1 T'_2$. *Hint:* Let $\langle \cdot, \cdot \rangle$ be the inner product associated with V , β a basis of eigenvectors for U , $\langle \cdot, \cdot \rangle'$ the inner product on V with respect to which β is orthonormal (see Exercise 22(a) of Section 6.1), and T_1 the positive definite operator according to Exercise 22. Show that U is self-adjoint with respect to $\langle \cdot, \cdot \rangle'$ and $U = T_1^{-1} U^* T_1$ (the adjoint is with respect to $\langle \cdot, \cdot \rangle$). Let $T_2 = T_1^{-1} U^*$.

6.5 UNITARY AND ORTHOGONAL OPERATORS AND THEIR MATRICES

In this section, we continue our analogy between complex numbers and linear operators. Recall that the adjoint of a linear operator acts similarly to the conjugate of a complex number (see, for example, Theorem 6.11 p. 357). A complex number z has length 1 if $z\bar{z} = 1$. In this section, we study those linear operators T on an inner product space V such that $TT^* = T^*T = I$. We will see that these are precisely the linear operators that “preserve length” in the sense that $\|T(x)\| = \|x\|$ for all $x \in V$. As another characterization, we prove that, on a finite-dimensional complex inner product space, these are the normal operators whose eigenvalues all have absolute value 1.

In past chapters, we were interested in studying those functions that preserve the structure of the underlying space. In particular, linear operators preserve the operations of vector addition and scalar multiplication, and isomorphisms preserve all the vector space structure. It is now natural to consider those linear operators T on an inner product space that preserve length. We will see that this condition guarantees, in fact, that T preserves the inner product.

Definitions. Let T be a linear operator on a finite-dimensional inner product space V (over F). If $\|T(x)\| = \|x\|$ for all $x \in V$, we call T a **unitary operator** if $F = C$ and an **orthogonal operator** if $F = R$.

It should be noted that in the infinite-dimensional case, an operator that preserves the norm is one-to-one, but not necessarily onto. If it is also onto, then we call it a **unitary or orthogonal operator**.

Clearly, any rotation or reflection in R^2 preserves length and hence is an orthogonal operator. We study these operators in much more detail in Section 6.11.

Example 1

Recall the inner product space H defined on page 330. Let $h \in H$ satisfy $|h(x)| = 1$ for all x . Define the linear operator T on H by $T(f) = hf$. Then

$$\|T(f)\|^2 = \|hf\|^2 = \frac{1}{2\pi} \int_0^{2\pi} h(t)f(t)\overline{h(t)f(t)} dt = \|f\|^2$$

since $|h(t)|^2 = 1$ for all t . So T is a unitary operator. ◆

Theorem 6.18. Let T be a linear operator on a finite-dimensional inner product space V . Then the following statements are equivalent.

- (a) $T^*T = I$.
- (b) $TT^* = I$.
- (c) $\langle T(x), T(y) \rangle = \langle x, y \rangle$ for all $x, y \in V$.

- (d) If β is an orthonormal basis for V , then $T(\beta)$ is an orthonormal basis for V .
- (e) There exists an orthonormal basis β for V such that $T(\beta)$ is an orthonormal basis for V .
- (f) $\|T(x)\| = \|x\|$ for all $x \in V$.

Thus all the conditions above are equivalent to the definition of a unitary or orthogonal operator. From (a) and (b), it follows that unitary or orthogonal operators are normal.

Before proving the theorem, we first prove a lemma. Compare this lemma to Exercise 11(b) of Section 6.4.

Lemma. Let U be a self-adjoint operator on an inner product space V , and suppose that $\langle x, U(x) \rangle = 0$ for all $x \in V$. Then $U = T_0$.

Proof. For any $x \in V$,

$$\begin{aligned} 0 &= \langle x + U(x), U(x + U(x)) \rangle \\ &= \langle x + U(x), U(x) + U^2(x) \rangle \\ &= \langle x, U(x) \rangle + \langle x, U^2(x) \rangle + \langle U(x), U(x) \rangle + \langle U(x), U^2(x) \rangle \\ &= 0 + \langle x, U^2(x) \rangle + \langle U(x), U(x) \rangle + 0 \\ &= \langle x, U^*U(x) \rangle + \|U(x)\|^2 \\ &= 2\|U(x)\|^2. \end{aligned}$$

So for any $x \in V$, $\|U(x)\| = 0$. It follows that $U = T_0$. ■

Proof of Theorem 6.18. Part (a) implies (b) by Theorem 6.10 and Exercise 10(c) of Section 2.4.

To prove that (b) implies (c), let $x, y \in V$. Then $\langle x, y \rangle = \langle T^*T(x), y \rangle = \langle T(x), T(y) \rangle$.

Next, we prove that (c) implies (d). Let $\beta = \{v_1, v_2, \dots, v_n\}$ be an orthonormal basis for V ; so $T(\beta) = \{T(v_1), T(v_2), \dots, T(v_n)\}$. It follows that $\langle T(v_i), T(v_j) \rangle = \langle v_i, v_j \rangle = \delta_{ij}$. Therefore $T(\beta)$ is an orthonormal basis for V .

That (d) implies (e) is obvious.

Next we prove that (e) implies (f). Let $x \in V$, and let $\beta = \{v_1, v_2, \dots, v_n\}$. Now

$$x = \sum_{i=1}^n a_i v_i$$

for some scalars a_i , and so

$$\|x\|^2 = \left\langle \sum_{i=1}^n a_i v_i, \sum_{j=1}^n a_j v_j \right\rangle = \sum_{i=1}^n \sum_{j=1}^n a_i \bar{a}_j \langle v_i, v_j \rangle$$

$$= \sum_{i=1}^n \sum_{j=1}^n a_i \bar{a}_j \delta_{ij} = \sum_{i=1}^n |a_i|^2$$

since β is orthonormal.

Applying the same manipulations to

$$\mathbf{T}(x) = \sum_{i=1}^n a_i \mathbf{T}(v_i)$$

and using the fact that $\mathbf{T}(\beta)$ is also orthonormal, we obtain

$$\|\mathbf{T}(x)\|^2 = \sum_{i=1}^n |a_i|^2.$$

Hence $\|\mathbf{T}(x)\| = \|x\|$.

Finally, we prove that (f) implies (a). For any $x \in V$, we have

$$\langle x, x \rangle = \|x\|^2 = \|\mathbf{T}(x)\|^2 = \langle \mathbf{T}(x), \mathbf{T}(x) \rangle = \langle x, \mathbf{T}^* \mathbf{T}(x) \rangle.$$

So $\langle x, (\mathbf{I} - \mathbf{T}^* \mathbf{T})(x) \rangle = 0$ for all $x \in V$. Let $U = \mathbf{I} - \mathbf{T}^* \mathbf{T}$; then U is self-adjoint, and $\langle x, U(x) \rangle = 0$ for all $x \in V$. Hence, by the lemma, we have $\mathbf{T}_0 = U = \mathbf{I} - \mathbf{T}^* \mathbf{T}$, and therefore $\mathbf{T}^* \mathbf{T} = \mathbf{I}$. ■

In the case that \mathbf{T} satisfies (c), we say that \mathbf{T} *preserves inner products*. In the case that \mathbf{T} satisfies (f), we say that \mathbf{T} *preserves norms*.

It follows immediately from the definition that every eigenvalue of a unitary or orthogonal operator has absolute value 1. In fact, even more is true.

Corollary 1. *Let \mathbf{T} be a linear operator on a finite-dimensional real inner product space V . Then V has an orthonormal basis of eigenvectors of \mathbf{T} with corresponding eigenvalues of absolute value 1 if and only if \mathbf{T} is both self-adjoint and orthogonal.*

Proof. Suppose that V has an orthonormal basis $\{v_1, v_2, \dots, v_n\}$ such that $\mathbf{T}(v_i) = \lambda_i v_i$ and $|\lambda_i| = 1$ for all i . By Theorem 6.17 (p. 371), \mathbf{T} is self-adjoint. Thus $(\mathbf{T}\mathbf{T}^*)(v_i) = \mathbf{T}(\lambda_i v_i) = \lambda_i \lambda_i v_i = \lambda_i^2 v_i = v_i$ for each i . So $\mathbf{T}\mathbf{T}^* = \mathbf{I}$, and again by Exercise 10 of Section 2.4, \mathbf{T} is orthogonal by Theorem 6.18(a).

If \mathbf{T} is self-adjoint, then, by Theorem 6.17, we have that V possesses an orthonormal basis $\{v_1, v_2, \dots, v_n\}$ such that $\mathbf{T}(v_i) = \lambda_i v_i$ for all i . If \mathbf{T} is also orthogonal, we have

$$|\lambda_i| \cdot \|v_i\| = \|\lambda_i v_i\| = \|\mathbf{T}(v_i)\| = \|v_i\|;$$

so $|\lambda_i| = 1$ for every i . ■

Corollary 2. Let T be a linear operator on a finite-dimensional complex inner product space V . Then V has an orthonormal basis of eigenvectors of T with corresponding eigenvalues of absolute value 1 if and only if T is unitary.

Proof. The proof is similar to the proof of Corollary 1. ■

Example 2

Let $T: \mathbb{R}^2 \rightarrow \mathbb{R}^2$ be a rotation by θ , where $0 < \theta < \pi$. It is clear geometrically that T “preserves length,” that is, that $\|T(x)\| = \|x\|$ for all $x \in \mathbb{R}^2$. The fact that rotations by a fixed angle preserve perpendicularity not only can be seen geometrically but now follows from (b) of Theorem 6.18. Perhaps the fact that such a transformation preserves the inner product is not so obvious; however, we obtain this fact from (b) also. Finally, an inspection of the matrix representation of T with respect to the standard ordered basis, which is

$$\begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix},$$

reveals that T is not self-adjoint for the given restriction on θ . As we mentioned earlier, this fact also follows from the geometric observation that T has no eigenvectors and from Theorem 6.15 (p. 368). It is seen easily from the preceding matrix that T^* is the rotation by $-\theta$. ♦

Definition. Let L be a one-dimensional subspace of \mathbb{R}^2 . We may view L as a line in the plane through the origin. A linear operator T on \mathbb{R}^2 is called a **reflection of \mathbb{R}^2 about L** if $T(x) = x$ for all $x \in L$ and $T(x) = -x$ for all $x \in L^\perp$.

As an example of a reflection, consider the operator defined in Example 3 of Section 2.5.

Example 3

Let T be a reflection of \mathbb{R}^2 about a line L through the origin. We show that T is an orthogonal operator. Select vectors $v_1 \in L$ and $v_2 \in L^\perp$ such that $\|v_1\| = \|v_2\| = 1$. Then $T(v_1) = v_1$ and $T(v_2) = -v_2$. Thus v_1 and v_2 are eigenvectors of T with corresponding eigenvalues 1 and -1 , respectively. Furthermore, $\{v_1, v_2\}$ is an orthonormal basis for \mathbb{R}^2 . It follows that T is an orthogonal operator by Corollary 1 to Theorem 6.18. ♦

We now examine the matrices that represent unitary and orthogonal transformations.

Definitions. A square matrix A is called an **orthogonal matrix** if $A^t A = AA^t = I$ and **unitary** if $A^* A = AA^* = I$.

Since for a real matrix A we have $A^* = A^t$, a real unitary matrix is also orthogonal. In this case, we call A **orthogonal** rather than unitary.

Note that the condition $AA^* = I$ is equivalent to the statement that the rows of A form an orthonormal basis for \mathbb{F}^n because

$$\delta_{ij} = I_{ij} = (AA^*)_{ij} = \sum_{k=1}^n A_{ik}(A^*)_{kj} = \sum_{k=1}^n A_{ik}\bar{A}_{jk},$$

and the last term represents the inner product of the i th and j th rows of A .

A similar remark can be made about the columns of A and the condition $A^*A = I$.

It also follows from the definition above and from Theorem 6.10 (p. 356) that a linear operator T on an inner product space V is unitary [orthogonal] if and only if $[T]_\beta$ is unitary [orthogonal] for some orthonormal basis β for V .

Example 4

From Example 2, the matrix

$$\begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix}$$

is clearly orthogonal. One can easily see that the rows of the matrix form an orthonormal basis for \mathbb{R}^2 . Similarly, the columns of the matrix form an orthonormal basis for \mathbb{R}^2 . ◆

Example 5

Let T be a reflection of \mathbb{R}^2 about a line L through the origin, let β be the standard ordered basis for \mathbb{R}^2 , and let $A = [T]_\beta$. Then $T = L_A$. Since T is an orthogonal operator and β is an orthonormal basis, A is an orthogonal matrix. We describe A .

Suppose that α is the angle from the positive x -axis to L . Let $v_1 = (\cos \alpha, \sin \alpha)$ and $v_2 = (-\sin \alpha, \cos \alpha)$. Then $\|v_1\| = \|v_2\| = 1$, $v_1 \in L$, and $v_2 \in L^\perp$. Hence $\gamma = \{v_1, v_2\}$ is an orthonormal basis for \mathbb{R}^2 . Because $T(v_1) = v_1$ and $T(v_2) = -v_2$, we have

$$[T]_\gamma = [L_A]_\gamma = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}.$$

Let

$$Q = \begin{pmatrix} \cos \alpha & -\sin \alpha \\ \sin \alpha & \cos \alpha \end{pmatrix}.$$

By the corollary to Theorem 2.23 (p. 115),

$$A = Q[L_A]_\gamma Q^{-1}$$

$$\begin{aligned}
 &= \begin{pmatrix} \cos \alpha & -\sin \alpha \\ \sin \alpha & \cos \alpha \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix} \begin{pmatrix} \cos \alpha & \sin \alpha \\ -\sin \alpha & \cos \alpha \end{pmatrix} \\
 &= \begin{pmatrix} \cos^2 \alpha - \sin^2 \alpha & 2 \sin \alpha \cos \alpha \\ 2 \sin \alpha \cos \alpha & -(\cos^2 \alpha - \sin^2 \alpha) \end{pmatrix} \\
 &= \begin{pmatrix} \cos 2\alpha & \sin 2\alpha \\ \sin 2\alpha & -\cos 2\alpha \end{pmatrix}. \quad \blacklozenge
 \end{aligned}$$

We know that, for a complex normal [real symmetric] matrix A , there exists an orthonormal basis β for \mathbb{F}^n consisting of eigenvectors of A . Hence A is similar to a diagonal matrix D . By the corollary to Theorem 2.23 (p. 115), the matrix Q whose columns are the vectors in β is such that $D = Q^{-1}AQ$. But since the columns of Q are an orthonormal basis for \mathbb{F}^n , it follows that Q is unitary [orthogonal]. In this case, we say that A is **unitarily equivalent** [**orthogonally equivalent**] to D . It is easily seen (see Exercise 18) that this relation is an equivalence relation on $M_{n \times n}(C)$ [$M_{n \times n}(R)$]. More generally, *A and B are unitarily equivalent [orthogonally equivalent] if and only if there exists a unitary [orthogonal] matrix P such that $A = P^*BP$.*

The preceding paragraph has proved half of each of the next two theorems.

Theorem 6.19. *Let A be a complex $n \times n$ matrix. Then A is normal if and only if A is unitarily equivalent to a diagonal matrix.*

Proof. By the preceding remarks, we need only prove that if A is unitarily equivalent to a diagonal matrix, then A is normal.

Suppose that $A = P^*DP$, where P is a unitary matrix and D is a diagonal matrix. Then

$$AA^* = (P^*DP)(P^*DP)^* = (P^*DP)(P^*D^*P) = P^*DID^*P = P^*DD^*P.$$

Similarly, $A^*A = P^*D^*DP$. Since D is a diagonal matrix, however, we have $DD^* = D^*D$. Thus $AA^* = A^*A$. ■

Theorem 6.20. *Let A be a real $n \times n$ matrix. Then A is symmetric if and only if A is orthogonally equivalent to a real diagonal matrix.*

Proof. The proof is similar to the proof of Theorem 6.19 and is left as an exercise. ■

Theorem 6.20 is used extensively in many areas of mathematics and statistics. See, for example, goo.gl/cbqApK.

Example 6

Let

$$A = \begin{pmatrix} 4 & 2 & 2 \\ 2 & 4 & 2 \\ 2 & 2 & 4 \end{pmatrix}.$$

Since A is symmetric, Theorem 6.20 tells us that A is orthogonally equivalent to a diagonal matrix. We find an orthogonal matrix P and a diagonal matrix D such that $P^t AP = D$.

To find P , we obtain an orthonormal basis of eigenvectors. It is easy to show that the eigenvalues of A are 2 and 8. The set $\{(-1, 1, 0), (-1, 0, 1)\}$ is a basis for the eigenspace corresponding to 2. Because this set is not orthogonal, we apply the Gram–Schmidt process to obtain the orthogonal set $\{(-1, 1, 0), -\frac{1}{2}(1, 1, -2)\}$. The set $\{(1, 1, 1)\}$ is a basis for the eigenspace corresponding to 8. Notice that $(1, 1, 1)$ is orthogonal to the preceding two vectors, as predicted by Theorem 6.15(d) (p. 368). Taking the union of these two bases and normalizing the vectors, we obtain the following orthonormal basis for \mathbb{R}^3 consisting of eigenvectors of A :

$$\left\{ \frac{1}{\sqrt{2}}(-1, 1, 0), \frac{1}{\sqrt{6}}(1, 1, -2), \frac{1}{\sqrt{3}}(1, 1, 1) \right\}.$$

Thus one possible choice for P is

$$P = \begin{pmatrix} \frac{-1}{\sqrt{2}} & \frac{1}{\sqrt{6}} & \frac{1}{\sqrt{3}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{6}} & \frac{1}{\sqrt{3}} \\ 0 & \frac{-2}{\sqrt{6}} & \frac{1}{\sqrt{3}} \end{pmatrix}, \quad \text{and} \quad D = \begin{pmatrix} 2 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 8 \end{pmatrix}. \quad \blacklozenge$$

Because of Schur's theorem (Theorem 6.14 p. 367), the next result is immediate. As it is the matrix form of Schur's theorem, we also refer to it as Schur's theorem.

Theorem 6.21 (Schur). *Let $A \in M_{n \times n}(F)$ be a matrix whose characteristic polynomial splits over F .*

- (a) *If $F = C$, then A is unitarily equivalent to a complex upper triangular matrix.*
- (b) *If $F = R$, then A is orthogonally equivalent to a real upper triangular matrix.*

Rigid Motions*

The purpose of this application is to characterize the so-called *rigid motions* of a finite-dimensional real inner product space. One may think intuitively of such a motion as a transformation that does not affect the shape of

a figure under its action, hence the term *rigid*. The key requirement for such a transformation is that it preserves distances.

Definition. Let V be a real inner product space. A function $f: V \rightarrow V$ is called a **rigid motion** if

$$\|f(x) - f(y)\| = \|x - y\|$$

for all $x, y \in V$.

For example, any orthogonal operator on a finite-dimensional real inner product space is a rigid motion.

Another class of rigid motions are the *translations*. A function $g: V \rightarrow V$, where V is a real inner product space, is called a **translation** if there exists a vector $v_0 \in V$ such that $g(x) = x + v_0$ for all $x \in V$. We say that g is the *translation by v_0* . It is a simple exercise to show that translations, as well as composites of rigid motions on a real inner product space, are also rigid motions. (See Exercise 22.) Thus an orthogonal operator on a finite-dimensional real inner product space V followed by a translation on V is a rigid motion on V . Remarkably, every rigid motion on V may be characterized in this way.

Theorem 6.22. Let $f: V \rightarrow V$ be a rigid motion on a finite-dimensional real inner product space V . Then there exists a unique orthogonal operator T on V and a unique translation g on V such that $f = g \circ T$.

Any orthogonal operator is a special case of this composite, in which the translation is by 0 . Any translation is also a special case, in which the orthogonal operator is the identity operator.

Proof. Let $T: V \rightarrow V$ be defined by

$$T(x) = f(x) - f(0)$$

for all $x \in V$. Note that $f = g \circ T$, where g is the translation by $f(0)$. Moreover, T is the composite of f and the translation by $-f(0)$; hence T is a rigid motion. We begin by showing that T is a linear operator. For any $x \in V$, we have

$$\|T(x)\|^2 = \|f(x) - f(0)\|^2 = \|x - 0\|^2 = \|x\|^2,$$

and consequently $\|T(x)\| = \|x\|$ for any $x \in V$. Thus for any $x, y \in V$,

$$\begin{aligned} \|T(x) - T(y)\|^2 &= \|T(x)\|^2 - 2 \langle T(x), T(y) \rangle + \|T(y)\|^2 \\ &= \|x\|^2 - 2 \langle T(x), T(y) \rangle + \|y\|^2 \end{aligned}$$

and

$$\|x - y\|^2 = \|x\|^2 - 2 \langle x, y \rangle + \|y\|^2.$$

But $\|\mathbf{T}(x) - \mathbf{T}(y)\|^2 = \|x - y\|^2$; so $\langle \mathbf{T}(x), \mathbf{T}(y) \rangle = \langle x, y \rangle$ for all $x, y \in V$.

We are now in a position to show that \mathbf{T} is a linear transformation. Let $x, y \in V$, and let $a \in R$. Then

$$\begin{aligned}\|\mathbf{T}(x + ay) - \mathbf{T}(x) - a\mathbf{T}(y)\|^2 &= \|[\mathbf{T}(x + ay) - \mathbf{T}(x)] - a\mathbf{T}(y)\|^2 \\&= \|\mathbf{T}(x + ay) - \mathbf{T}(x)\|^2 + a^2\|\mathbf{T}(y)\|^2 - 2a\langle \mathbf{T}(x + ay) - \mathbf{T}(x), \mathbf{T}(y) \rangle \\&= \|(x + ay) - x\|^2 + a^2\|y\|^2 - 2a[\langle \mathbf{T}(x + ay), \mathbf{T}(y) \rangle - \langle \mathbf{T}(x), \mathbf{T}(y) \rangle] \\&= a^2\|y\|^2 + a^2\|y\|^2 - 2a[\langle x + ay, y \rangle - \langle x, y \rangle] \\&= 2a^2\|y\|^2 - 2a[\langle x, y \rangle + a\|y\|^2 - \langle x, y \rangle] \\&= 0.\end{aligned}$$

Thus $\mathbf{T}(x + ay) = \mathbf{T}(x) + a\mathbf{T}(y)$, and hence \mathbf{T} is linear. Since we have already shown that \mathbf{T} preserves inner products, \mathbf{T} is an orthogonal operator.

To prove uniqueness, suppose that u_0 and v_0 are in V and \mathbf{T} and \mathbf{U} are orthogonal operators on V such that

$$f(x) = \mathbf{T}(x) + u_0 = \mathbf{U}(x) + v_0$$

for all $x \in V$. Substituting $x = 0$ in the preceding equation yields $u_0 = v_0$, and hence the translation is unique. This equation, therefore, reduces to $\mathbf{T}(x) = \mathbf{U}(x)$ for all $x \in V$, and hence $\mathbf{T} = \mathbf{U}$. ■

Orthogonal Operators on R^2

Because of Theorem 6.22, an understanding of rigid motions requires a characterization of orthogonal operators. The next result characterizes orthogonal operators on R^2 . We postpone the case of orthogonal operators on more general spaces to Section 6.11.

Theorem 6.23. *Let \mathbf{T} be an orthogonal operator on R^2 , and let $A = [\mathbf{T}]_\beta$, where β is the standard ordered basis for R^2 . Then exactly one of the following conditions is satisfied:*

- (a) \mathbf{T} is a rotation, and $\det(A) = 1$.
- (b) \mathbf{T} is a reflection about a line through the origin, and $\det(A) = -1$.

Proof. Because \mathbf{T} is an orthogonal operator, $\mathbf{T}(\beta) = \{\mathbf{T}(e_1), \mathbf{T}(e_2)\}$ is an orthonormal basis for R^2 by Theorem 6.18(c). Since $\mathbf{T}(e_1)$ is a unit vector, there is a unique angle θ , $0 \leq \theta < 2\pi$, such that $\mathbf{T}(e_1) = (\cos \theta, \sin \theta)$. Since $\mathbf{T}(e_2)$ is a unit vector and is orthogonal to $\mathbf{T}(e_1)$, there are only two possible choices for $\mathbf{T}(e_2)$. Either

$$\mathbf{T}(e_2) = (-\sin \theta, \cos \theta) \quad \text{or} \quad \mathbf{T}(e_2) = (\sin \theta, -\cos \theta).$$

First, suppose that $\mathbf{T}(e_2) = (-\sin \theta, \cos \theta)$. Then $A = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix}$.

It follows from Example 1 of Section 6.4 that T is a rotation by the angle θ . Also

$$\det(A) = \cos^2 \theta + \sin^2 \theta = 1.$$

Now suppose that $T(e_2) = (\sin \theta, -\cos \theta)$. Then $A = \begin{pmatrix} \cos \theta & \sin \theta \\ \sin \theta & -\cos \theta \end{pmatrix}$.

Comparing this matrix to the matrix A of Example 5, we see that T is the reflection of \mathbb{R}^2 about a line L such that $\alpha = \theta/2$ is the angle from the positive x -axis to L . Furthermore,

$$\det(A) = -\cos^2 \theta - \sin^2 \theta = -1. \quad \blacksquare$$

Combining Theorems 6.22 and 6.23, we obtain the following characterization of rigid motions on \mathbb{R}^2 .

Corollary. *Any rigid motion on \mathbb{R}^2 is either a rotation followed by a translation or a reflection about a line through the origin followed by a translation.*

Example 7

Let

$$A = \begin{pmatrix} \frac{1}{\sqrt{5}} & \frac{2}{\sqrt{5}} \\ \frac{2}{\sqrt{5}} & \frac{-1}{\sqrt{5}} \end{pmatrix}.$$

We show that L_A is the reflection of \mathbb{R}^2 about a line L through the origin, and then describe L .

Clearly $AA^* = A^*A = I$, and therefore A is an orthogonal matrix. Hence L_A is an orthogonal operator. Furthermore,

$$\det(A) = -\frac{1}{5} - \frac{4}{5} = -1,$$

and thus L_A is a reflection of \mathbb{R}^2 about a line L through the origin by Theorem 6.23. Since L is the one-dimensional eigenspace corresponding to the eigenvalue 1 of L_A , it suffices to find an eigenvector of L_A corresponding to 1. One such vector is $v = (2, \sqrt{5} - 1)$. Thus L is the span of $\{v\}$. Alternatively, L is the line through the origin with slope $(\sqrt{5} - 1)/2$, and hence is the line with the equation

$$y = \frac{\sqrt{5} - 1}{2} x. \quad \blacklozenge$$

Conic Sections

As an application of Theorem 6.20, we consider the quadratic equation

$$ax^2 + 2bxy + cy^2 + dx + ey + f = 0. \quad (2)$$

For special choices of the coefficients in (2), we obtain the various conic sections. For example, if $a = c = 1$, $b = d = e = 0$, and $f = -1$, we obtain the circle $x^2 + y^2 = 1$ with center at the origin. The remaining conic sections, namely, the ellipse, parabola, and hyperbola, are obtained by other choices of the coefficients. If $b = 0$, then it is easy to graph the equation by the method of completing the square because the xy -term is absent. For example, the equation $x^2 + 2x + y^2 + 4y + 2 = 0$ may be rewritten as $(x+1)^2 + (y+2)^2 = 3$, which describes a circle with radius $\sqrt{3}$ and center at $(-1, -2)$ in the xy -coordinate system. If we consider the transformation of coordinates $(x, y) \rightarrow (x', y')$, where $x' = x + 1$ and $y' = y + 2$, then our equation simplifies to $(x')^2 + (y')^2 = 3$. This change of variable allows us to eliminate the x - and y -terms.

We now concentrate solely on the elimination of the xy -term. To accomplish this, we consider the expression

$$ax^2 + 2bxy + cy^2, \quad (3)$$

which is called the **associated quadratic form** of (2). Quadratic forms are studied in more generality in Section 6.8.

If we let

$$A = \begin{pmatrix} a & b \\ b & c \end{pmatrix} \quad \text{and} \quad X = \begin{pmatrix} x \\ y \end{pmatrix},$$

then (3) may be written as $X^t AX = \langle AX, X \rangle$. For example, the quadratic form $3x^2 + 4xy + 6y^2$ may be written as

$$X^t \begin{pmatrix} 3 & 2 \\ 2 & 6 \end{pmatrix} X.$$

The fact that A is symmetric is crucial in our discussion. For, by Theorem 6.20, we may choose an orthogonal matrix P and a diagonal matrix D with real diagonal entries λ_1 and λ_2 such that $P^t AP = D$. Now define

$$X' = \begin{pmatrix} x' \\ y' \end{pmatrix}$$

by $X' = P^t X$ or, equivalently, by $PX' = PP^t X = X$. Then

$$X^t AX = (PX')^t A(PX') = X'^t (P^t AP) X' = X'^t DX' = \lambda_1(x')^2 + \lambda_2(y')^2.$$

Thus the transformation $(x, y) \rightarrow (x', y')$ allows us to eliminate the xy -term in (3), and hence in (2).

Furthermore, since P is orthogonal, we have by Theorem 6.23 (with $T = L_P$) that $\det(P) = \pm 1$. If $\det(P) = -1$, we may interchange the columns of P to obtain a matrix Q . Because the columns of P form an orthonormal basis of eigenvectors of A , the same is true of the columns of Q . Therefore

$$Q^t AQ = \begin{pmatrix} \lambda_2 & 0 \\ 0 & \lambda_1 \end{pmatrix}.$$

Notice that $\det(Q) = -\det(P) = 1$. So, if $\det(P) = -1$, we can take Q for our new P ; consequently, we may always choose P so that $\det(P) = 1$. By Lemma 4 to Theorem 6.22 (with $T = L_P$), it follows that matrix P represents a rotation.

In summary, the xy -term in (2) may be eliminated by a rotation of the x -axis and y -axis to new axes x' and y' given by $X = PX'$, where P is an orthogonal matrix and $\det(P) = 1$. Furthermore, the coefficients of $(x')^2$ and $(y')^2$ are the eigenvalues of

$$A = \begin{pmatrix} a & b \\ b & c \end{pmatrix}.$$

This result is a restatement of a result known as the *principal axis theorem* for \mathbb{R}^2 . The arguments above, of course, are easily extended to quadratic equations in n variables. For example, in the case $n = 3$, by special choices of the coefficients, we obtain the quadratic surfaces—the elliptic cone, the ellipsoid, the hyperbolic paraboloid, etc.

As an illustration of the preceding transformation, consider the quadratic equation

$$2x^2 - 4xy + 5y^2 - 36 = 0,$$

for which the associated quadratic form is $2x^2 - 4xy + 5y^2$. In the notation we have been using,

$$A = \begin{pmatrix} 2 & -2 \\ -2 & 5 \end{pmatrix},$$

so that the eigenvalues of A are 1 and 6 with associated eigenvectors

$$\begin{pmatrix} 2 \\ 1 \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} -1 \\ 2 \end{pmatrix}.$$

As expected (from Theorem 6.15(d) p. 368), these vectors are orthogonal. The corresponding orthonormal basis of eigenvectors

$$\beta = \left\{ \begin{pmatrix} \frac{2}{\sqrt{5}} \\ \frac{1}{\sqrt{5}} \end{pmatrix}, \begin{pmatrix} \frac{-1}{\sqrt{5}} \\ \frac{2}{\sqrt{5}} \end{pmatrix} \right\}$$

determines new axes x' and y' as in Figure 6.4. Hence if

$$P = \begin{pmatrix} \frac{2}{\sqrt{5}} & \frac{-1}{\sqrt{5}} \\ \frac{1}{\sqrt{5}} & \frac{2}{\sqrt{5}} \end{pmatrix} = \frac{1}{\sqrt{5}} \begin{pmatrix} 2 & -1 \\ 1 & 2 \end{pmatrix},$$

then

$$P^t A P = \begin{pmatrix} 1 & 0 \\ 0 & 6 \end{pmatrix}.$$

Under the transformation $X = PX'$ or

$$\begin{aligned} x &= \frac{2}{\sqrt{5}}x' - \frac{1}{\sqrt{5}}y' \\ y &= \frac{1}{\sqrt{5}}x' + \frac{2}{\sqrt{5}}y', \end{aligned}$$

we have the new quadratic form $(x')^2 + 6(y')^2$. Thus the original equation $2x^2 - 4xy + 5y^2 = 36$ may be written in the form $(x')^2 + 6(y')^2 = 36$ relative to a new coordinate system with the x' - and y' -axes in the directions of the first and second vectors of β , respectively. It is clear that this equation represents

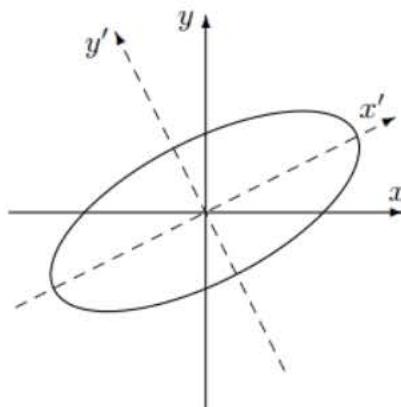


Figure 6.4

an ellipse. (See Figure 6.4.) Note that the preceding matrix P has the form

$$\begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix},$$

where $\theta = \cos^{-1} \frac{2}{\sqrt{5}} \approx 26.6^\circ$. So P is the matrix representation of a rotation of \mathbb{R}^2 through the angle θ . Thus the change of variable $X = PX'$ can be accomplished by this rotation of the x - and y -axes. There is another possibility for P , however. If the eigenvector of A corresponding to the eigenvalue 6 is taken to be $(1, -2)$ instead of $(-1, 2)$, and the eigenvalues are interchanged, then we obtain the matrix

$$\begin{pmatrix} \frac{1}{\sqrt{5}} & \frac{2}{\sqrt{5}} \\ -2 & \frac{1}{\sqrt{5}} \end{pmatrix},$$

which is the matrix representation of a rotation through the angle $\theta = \sin^{-1} \left(-\frac{2}{\sqrt{5}} \right) \approx -63.4^\circ$. This possibility produces the same ellipse as the one in Figure 6.4, but interchanges the names of the x' - and y' -axes.

EXERCISES

1. Label the following statements as true or false. Assume that the underlying inner product spaces are finite-dimensional.
 - (a) Every unitary operator is normal.
 - (b) Every orthogonal operator is diagonalizable.
 - (c) A matrix is unitary if and only if it is invertible.
 - (d) If two matrices are unitarily equivalent, then they are also similar.
 - (e) The sum of unitary matrices is unitary.
 - (f) The adjoint of a unitary operator is unitary.
 - (g) If T is an orthogonal operator on V , then $[T]_\beta$ is an orthogonal matrix for any ordered basis β for V .
 - (h) If all the eigenvalues of a linear operator are 1, then the operator must be unitary or orthogonal.
 - (i) A linear operator may preserve norms without preserving inner products.
2. For each of the following matrices A , find an orthogonal or unitary matrix P and a diagonal matrix D such that $P^*AP = D$.

(a) $\begin{pmatrix} 1 & 2 \\ 2 & 1 \end{pmatrix}$	(b) $\begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix}$	(c) $\begin{pmatrix} 2 & 3-3i \\ 3+3i & 5 \end{pmatrix}$
(d) $\begin{pmatrix} 0 & 2 & 2 \\ 2 & 0 & 2 \\ 2 & 2 & 0 \end{pmatrix}$	(e) $\begin{pmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{pmatrix}$	
3. Prove that the composite of unitary [orthogonal] operators is unitary [orthogonal].
4. For $z \in C$, define $T_z: C \rightarrow C$ by $T_z(u) = zu$. Characterize those z for which T_z is normal, self-adjoint, or unitary.
5. Which of the following pairs of matrices are unitarily equivalent?

(a) $\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$ and $\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$	(b) $\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$ and $\begin{pmatrix} 0 & \frac{1}{2} \\ \frac{1}{2} & 0 \end{pmatrix}$
---	---

$$(a) \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \quad (b) \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} 0 & \frac{1}{2} \\ \frac{1}{2} & 0 \end{pmatrix}$$

$$\begin{array}{ll}
 \text{(c)} & \begin{pmatrix} 0 & 1 & 0 \\ -1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} 2 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & 0 \end{pmatrix} \\
 \text{(d)} & \begin{pmatrix} 0 & 1 & 0 \\ -1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} 1 & 0 & 0 \\ 0 & i & 0 \\ 0 & 0 & -i \end{pmatrix} \\
 \text{(e)} & \begin{pmatrix} 1 & 1 & 0 \\ 0 & 2 & 2 \\ 0 & 0 & 3 \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{pmatrix}
 \end{array}$$

6. Let V be the inner product space of complex-valued continuous functions on $[0, 1]$ with the inner product

$$\langle f, g \rangle = \int_0^1 f(t) \overline{g(t)} dt.$$

Let $h \in V$, and define $T: V \rightarrow V$ by $T(f) = hf$. Prove that T is a unitary operator if and only if $|h(t)| = 1$ for $0 \leq t \leq 1$.

red Hint for the “only if” part: Suppose that T is unitary. Set $f(t) = 1 - |h(t)|^2$ and $g(t) = 1$. Show that

$$\int_0^1 (1 - |h(t)|^2)^2 dt = 0,$$

and use the fact that if the integral of a nonnegative continuous function is zero, then the function is identically zero.

7. Prove that if T is a unitary operator on a finite-dimensional inner product space V , then T has a unitary *square root*; that is, there exists a unitary operator U such that $T = U^2$. Visit goo.gl/jADTaS for a solution.
8. Let T be a self-adjoint linear operator on a finite-dimensional inner product space. Prove that $(T + iI)(T - iI)^{-1}$ is unitary using Exercise 10 of Section 6.4.
9. Let U be a linear operator on a finite-dimensional inner product space V . If $\|U(x)\| = \|x\|$ for all x in some orthonormal basis for V , must U be unitary? Justify your answer with a proof or a counterexample.
10. Let A be an $n \times n$ real symmetric or complex normal matrix. Prove that

$$\operatorname{tr}(A) = \sum_{i=1}^n \lambda_i \quad \text{and} \quad \operatorname{tr}(A^* A) = \sum_{i=1}^n |\lambda_i|^2,$$

where the λ_i 's are the (not necessarily distinct) eigenvalues of A .

11. Find an orthogonal matrix whose first row is $(\frac{1}{3}, \frac{2}{3}, \frac{2}{3})$.
12. Let A be an $n \times n$ real symmetric or complex normal matrix. Prove that

$$\det(A) = \prod_{i=1}^n \lambda_i,$$

where the λ_i 's are the (not necessarily distinct) eigenvalues of A .

13. Suppose that A and B are diagonalizable matrices. Prove or disprove that A is similar to B if and only if A and B are unitarily equivalent.
14. Prove that if A and B are unitarily equivalent matrices, then A is positive definite [semidefinite] if and only if B is positive definite [semidefinite]. (See the definitions in the exercises in Section 6.4.)
15. Let U be a unitary operator on an inner product space V , and let W be a finite-dimensional U -invariant subspace of V . Prove that

- (a) $U(W) = W$;
 (b) W^\perp is U -invariant.

Contrast (b) with Exercise 16.

16. Find an example of a unitary operator U on an inner product space and a U -invariant subspace W such that W^\perp is not U -invariant.
17. Prove that a matrix that is both unitary and upper triangular must be a diagonal matrix.
18. Show that “is unitarily equivalent to” is an equivalence relation on $M_{n \times n}(C)$.
19. Let W be a finite-dimensional subspace of an inner product space V . By Theorem 6.7 (p. 349) and the exercises of Section 1.3, $V = W \oplus W^\perp$. Define $U: V \rightarrow V$ by $U(v_1 + v_2) = v_1 - v_2$, where $v_1 \in W$ and $v_2 \in W^\perp$. Prove that U is a self-adjoint unitary operator.
20. Let V be a finite-dimensional inner product space. A linear operator U on V is called a **partial isometry** if there exists a subspace W of V such that $\|U(x)\| = \|x\|$ for all $x \in W$ and $U(x) = 0$ for all $x \in W^\perp$. Observe that W need *not* be U -invariant. Suppose that U is such an operator and $\{v_1, v_2, \dots, v_k\}$ is an orthonormal basis for W . Prove the following results.
- (a) $\langle U(x), U(y) \rangle = \langle x, y \rangle$ for all $x, y \in W$. Hint: Use Exercise 20 of Section 6.1.
- (b) $\{U(v_1), U(v_2), \dots, U(v_k)\}$ is an orthonormal basis for $R(U)$.

- (c) There exists an orthonormal basis γ for V such that the first k columns of $[U]_\gamma$ form an orthonormal set and the remaining columns are zero.
- (d) Let $\{w_1, w_2, \dots, w_j\}$ be an orthonormal basis for $R(U)^\perp$ and $\beta = \{U(v_1), U(v_2), \dots, U(v_k), w_1, \dots, w_j\}$. Then β is an orthonormal basis for V .
- (e) Let T be the linear operator on V that satisfies $T(U(v_i)) = v_i$ ($1 \leq i \leq k$) and $T(w_i) = 0$ ($1 \leq i \leq j$). Then T is well defined, and $T = U^*$. Hint: Show that $\langle U(x), y \rangle = \langle x, T(y) \rangle$ for all $x, y \in \beta$. There are four cases.
- (f) U^* is a partial isometry.

This exercise is continued in Exercise 9 of Section 6.6.

- 21.** Let A and B be $n \times n$ matrices that are unitarily equivalent.
- (a) Prove that $\text{tr}(A^*A) = \text{tr}(B^*B)$.
 - (b) Use (a) to prove that
- $$\sum_{i,j=1}^n |A_{ij}|^2 = \sum_{i,j=1}^n |B_{ij}|^2.$$
- (c) Use (b) to show that the matrices
- $$\begin{pmatrix} 1 & 2 \\ 2 & i \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} i & 4 \\ 1 & 1 \end{pmatrix}$$
- are *not* unitarily equivalent.
- 22.** Let V be a real inner product space.
- (a) Prove that any translation on V is a rigid motion.
 - (b) Prove that the composite of any two rigid motions on V is a rigid motion on V .
- 23.** Prove the following variant of Theorem 6.22: If $f: V \rightarrow V$ is a rigid motion on a finite-dimensional real inner product space V , then there exists a unique orthogonal operator T on V and a unique translation g on V such that $f = T \circ g$. (Note that the conclusion of Theorem 6.22 has $f = g \circ T$).
- 24.** Let T and U be orthogonal operators on \mathbb{R}^2 . Use Theorem 6.23 to prove the following results.
- (a) If T and U are both reflections about lines through the origin, then UT is a rotation.
 - (b) If T is a rotation and U is a reflection about a line through the origin, then both UT and TU are reflections about lines through the origin.

25. Suppose that T and U are reflections of \mathbb{R}^2 about the respective lines L and L' through the origin and that ϕ and ψ are the angles from the positive x -axis to L and L' , respectively. By Exercise 24, UT is a rotation. Find its angle of rotation.
26. Suppose that T and U are orthogonal operators on \mathbb{R}^2 such that T is the rotation by the angle ϕ and U is the reflection about the line L through the origin. Let ψ be the angle from the positive x -axis to L . By Exercise 24, both UT and TU are reflections about lines L_1 and L_2 , respectively, through the origin.
- Find the angle θ from the positive x -axis to L_1 .
 - Find the angle θ from the positive x -axis to L_2 .
27. Find new coordinates x', y' so that the following quadratic forms can be written as $\lambda_1(x')^2 + \lambda_2(y')^2$.
- $x^2 + 4xy + y^2$
 - $2x^2 + 2xy + 2y^2$
 - $x^2 - 12xy - 4y^2$
 - $3x^2 + 2xy + 3y^2$
 - $x^2 - 2xy + y^2$
28. Consider the expression $X^t AX$, where $X^t = (x, y, z)$ and A is as defined in Exercise 2(e). Find a change of coordinates x', y', z' so that the preceding expression is of the form $\lambda_1(x')^2 + \lambda_2(y')^2 + \lambda_3(z')^2$.
29. *QR-Factorization.* Let w_1, w_2, \dots, w_n be linearly independent vectors in \mathbb{F}^n , and let v_1, v_2, \dots, v_n be the orthogonal vectors obtained from w_1, w_2, \dots, w_n by the Gram–Schmidt process. Let u_1, u_2, \dots, u_n be the orthonormal basis obtained by normalizing the v_i 's.
- Solving (1) in Section 6.2 for w_k in terms of u_k , show that

$$w_k = \|v_k\|u_k + \sum_{j=1}^{k-1} \langle w_k, u_j \rangle u_j \quad (1 \leq k \leq n).$$

- (b) Let A and Q denote the $n \times n$ matrices in which the k th columns are w_k and u_k , respectively. Define $R \in \mathbb{M}_{n \times n}(F)$ by

$$R_{jk} = \begin{cases} \|v_j\| & \text{if } j = k \\ \langle w_k, u_j \rangle & \text{if } j < k \\ 0 & \text{if } j > k. \end{cases}$$

Prove $A = QR$.

- (c) Compute Q and R as in (b) for the 3×3 matrix whose columns are the vectors $(1, 1, 0)$, $(2, 0, 1)$, and $(2, 2, 1)$.

- (d) Since Q is unitary [orthogonal] and R is upper triangular in (b), we have shown that every invertible matrix is the product of a unitary [orthogonal] matrix and an upper triangular matrix. Suppose that $A \in M_{n \times n}(F)$ is invertible and $A = Q_1 R_1 = Q_2 R_2$, where $Q_1, Q_2 \in M_{n \times n}(F)$ are unitary and $R_1, R_2 \in M_{n \times n}(F)$ are upper triangular. Prove that $D = R_2 R_1^{-1}$ is a unitary diagonal matrix.

Hint: Use Exercise 17.

- (e) The QR factorization described in (b) provides an orthogonalization method for solving a linear system $Ax = b$ when A is invertible. Decompose A to QR , by the Gram–Schmidt process or other means, where Q is unitary and R is upper triangular. Then $QRx = b$, and hence $Rx = Q^*b$. This last system can be easily solved since R is upper triangular.¹

Use the orthogonalization method and (c) to solve the system

$$\begin{array}{rl} x_1 + 2x_2 + 2x_3 &= 1 \\ x_1 &+ 2x_3 = 11 \\ x_2 + x_3 &= -1. \end{array}$$

30. Suppose that β and γ are ordered bases for an n -dimensional real [complex] inner product space V . Prove that if Q is an orthogonal [unitary] $n \times n$ matrix that changes γ -coordinates into β -coordinates, then β is orthonormal if and only if γ is orthonormal.

The following definition is used in Exercises 31 and 32.

Definition. Let V be a finite-dimensional complex [real] inner product space, and let u be a unit vector in V . Define the **Householder** operator $H_u : V \rightarrow V$ by $H_u(x) = x - 2\langle x, u \rangle u$ for all $x \in V$.

31. Let H_u be a Householder operator on a finite-dimensional inner product space V . Prove the following results.

- (a) H_u is linear.
- (b) $H_u(x) = x$ if and only if x is orthogonal to u .
- (c) $H_u(u) = -u$.
- (d) $H_u^* = H_u$ and $H_u^2 = I$, and hence H_u is a unitary [orthogonal] operator on V .

(Note: If V is a real inner product space, then in the language of Section 6.11, H_u is a reflection.)

¹At one time, because of its great stability, this method for solving large systems of linear equations with a computer was being advocated as a better method than Gaussian elimination even though it requires about three times as much work. (Later, however, J. H. Wilkinson showed that if Gaussian elimination is done “properly,” then it is nearly as stable as the orthogonalization method.)

32. Let V be a finite-dimensional inner product space over F . Let x and y be linearly independent vectors in V such that $\|x\| = \|y\|$.
- If $F = C$, prove that there exists a unit vector u in V and a complex number θ with $|\theta| = 1$ such that $H_u(x) = \theta y$. Hint: Choose θ so that $\langle x, \theta y \rangle$ is real, and set $u = \frac{1}{\|x - \theta y\|}(x - \theta y)$.
 - If $F = R$, prove that there exists a unit vector u in V such that $H_u(x) = y$.

6.6 ORTHOGONAL PROJECTIONS AND THE SPECTRAL THEOREM

In this section, we rely heavily on Theorems 6.16 (p. 369) and 6.17 (p. 371) to develop an elegant representation of a normal (if $F = C$) or a self-adjoint (if $F = R$) operator T on a finite-dimensional inner product space. We prove that T can be written in the form $\lambda_1 T_1 + \lambda_2 T_2 + \cdots + \lambda_k T_k$, where $\lambda_1, \lambda_2, \dots, \lambda_k$ are the distinct eigenvalues of T and T_1, T_2, \dots, T_k are *orthogonal projections*. We must first develop some results about these special projections.

We assume that the reader is familiar with the results about direct sums developed at the end of Section 5.2. The special case where V is a direct sum of two subspaces is considered in the exercises of Section 1.3.

Recall from the exercises of Section 2.1 that if $V = W_1 \oplus W_2$, then a linear operator T on V is the **projection on W_1 along W_2** if, whenever $x = x_1 + x_2$, with $x_1 \in W_1$ and $x_2 \in W_2$, we have $T(x) = x_1$. By Exercise 27 of Section 2.1, we have

$$R(T) = W_1 = \{x \in V : T(x) = x\} \quad \text{and} \quad N(T) = W_2.$$

So $V = R(T) \oplus N(T)$. Thus there is no ambiguity if we refer to T as a “projection on W_1 ” or simply as a “projection.” In fact, it can be shown (see Exercise 17 of Section 2.3) that T is a projection if and only if $T = T^2$. Because $V = W_1 \oplus W_2 = W_1 \oplus W_3$ does *not* imply that $W_2 = W_3$, we see that W_1 does not uniquely determine T . For an *orthogonal projection* T , however, T is uniquely determined by its range.

Definition. Let V be an inner product space, and let $T: V \rightarrow V$ be a projection. We say that T is an **orthogonal projection** if $R(T)^\perp = N(T)$ and $N(T)^\perp = R(T)$.

Note that by Exercise 13(c) of Section 6.2, if V is finite-dimensional, we need only assume that one of the equalities in this definition holds. For example, if $R(T)^\perp = N(T)$, then $R(T) = R(T)^{\perp\perp} = N(T)^\perp$.

An orthogonal projection is *not* the same as an orthogonal operator. In Figure 6.5, T is an orthogonal projection, but T is clearly not an orthogonal operator because $\|T(v)\| \neq \|v\|$.

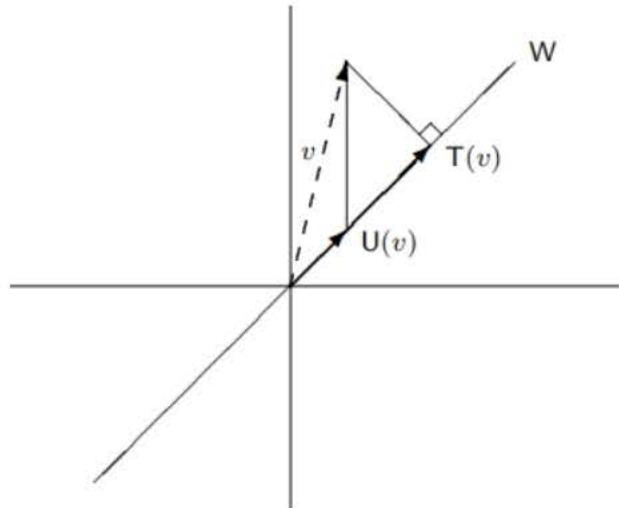


Figure 6.5

Now assume that W is a finite-dimensional subspace of an inner product space V . In the notation of Theorem 6.6 (p. 347), we can define a function $T: V \rightarrow V$ by $T(y) = u$. It is easy to show that T is an orthogonal projection on W . We can say even more—there exists exactly one orthogonal projection on W . For if T and U are orthogonal projections on W , then $R(T) = W = R(U)$. Hence $N(T) = R(T)^\perp = R(U)^\perp = N(U)$, and since every projection is uniquely determined by its range and null space, we have $T = U$. We call T the **orthogonal projection** of V on W .

To understand the geometric difference between an arbitrary projection on W and the orthogonal projection on W , let $V = \mathbb{R}^2$ and $W = \text{span}\{(1, 1)\}$. Define U and T as in Figure 6.5, where $T(v)$ is the foot of a perpendicular from v on the line $y = x$ and $U(a_1, a_2) = (a_1, a_1)$. Then T is the orthogonal projection of V on W , and U is a different projection on W . Note that $v - T(v) \in W^\perp$, whereas $v - U(v) \notin W^\perp$. From Figure 6.5, we see that $T(v)$ is the “best approximation in W to v ”; that is, if $w \in W$, then $\|w - v\| \geq \|T(v) - v\|$. In fact, this approximation property characterizes T . These results follow immediately from the corollary to Theorem 6.6 (p. 348).

As an application to Fourier analysis, recall the inner product space H and the orthonormal set S in Example 9 of Section 6.1. Define a **trigonometric polynomial of degree n** to be a function $g \in H$ of the form

$$g(t) = \sum_{j=-n}^n a_j f_j(t) = \sum_{j=-n}^n a_j e^{ijt},$$

where a_n or a_{-n} is nonzero.

Let $f \in H$. We show that the best approximation to f by a trigonometric

polynomial of degree less than or equal to n is the trigonometric polynomial whose coefficients are the Fourier coefficients of f relative to the orthonormal set S . For this result, let $W = \text{span}(\{f_j : |j| \leq n\})$, and let T be the orthogonal projection of H on W . The corollary to Theorem 6.6 (p. 348) tells us that the best approximation to f by a function in W is

$$T(f) = \sum_{j=-n}^n \langle f, f_j \rangle f_j.$$

For an application of this material to electronic music, visit [goo.gl/EN5Fai](#).

An algebraic characterization of orthogonal projections follows in the next theorem.

Theorem 6.24. *Let V be an inner product space, and let T be a linear operator on V . Then T is an orthogonal projection if and only if T has an adjoint T^* and $T^2 = T = T^*$.*

Proof. Suppose that T is an orthogonal projection. Since $T^2 = T$ because T is a projection, we need only show that T^* exists and $T = T^*$. Now $V = R(T) \oplus N(T)$ and $R(T)^\perp = N(T)$. Let $x, y \in V$. Then we can write $x = x_1 + x_2$ and $y = y_1 + y_2$, where $x_1, y_1 \in R(T)$ and $x_2, y_2 \in N(T)$. Hence

$$\langle x, T(y) \rangle = \langle x_1 + x_2, y_1 \rangle = \langle x_1, y_1 \rangle + \langle x_2, y_1 \rangle = \langle x_1, y_1 \rangle$$

and

$$\langle T(x), y \rangle = \langle x_1, y_1 + y_2 \rangle = \langle x_1, y_1 \rangle + \langle x_1, y_2 \rangle = \langle x_1, y_1 \rangle.$$

So $\langle x, T(y) \rangle = \langle T(x), y \rangle$ for all $x, y \in V$; thus T^* exists and $T = T^*$.

Now suppose that $T^2 = T = T^*$. Then T is a projection by Exercise 17 of Section 2.3, and hence we must show that $R(T) = N(T)^\perp$ and $R(T)^\perp = N(T)$. Let $x \in R(T)$ and $y \in N(T)$. Then $x = T(x) = T^*(x)$, and so

$$\langle x, y \rangle = \langle T^*(x), y \rangle = \langle x, T(y) \rangle = \langle x, 0 \rangle = 0.$$

Therefore $x \in N(T)^\perp$, from which it follows that $R(T) \subseteq N(T)^\perp$.

Let $y \in N(T)^\perp$. We must show that $y \in R(T)$, that is, $T(y) = y$. Now

$$\begin{aligned} \|y - T(y)\|^2 &= \langle y - T(y), y - T(y) \rangle \\ &= \langle y, y - T(y) \rangle - \langle T(y), y - T(y) \rangle. \end{aligned}$$

Since $y - T(y) \in N(T)$, the first term must equal zero. But also

$$\langle T(y), y - T(y) \rangle = \langle y, T^*(y - T(y)) \rangle = \langle y, T(y - T(y)) \rangle = \langle y, 0 \rangle = 0.$$

Thus $y - T(y) = 0$; that is, $y = T(y) \in R(T)$. Hence $R(T) = N(T)^\perp$.

Using the preceding results, we have $R(T)^\perp = N(T)^{\perp\perp} \supseteq N(T)$ by Exercise 13(b) of Section 6.2. Now suppose that $x \in R(T)^\perp$. For any $y \in V$, we have $\langle T(x), y \rangle = \langle x, T^*(y) \rangle = \langle x, T(y) \rangle = 0$. So $T(x) = 0$, and thus $x \in N(T)$. Hence $R(T)^\perp = N(T)$. ■

Let V be a finite-dimensional inner product space, W be a subspace of V , and T be the orthogonal projection of V on W . We may choose an orthonormal basis $\beta = \{v_1, v_2, \dots, v_n\}$ for V such that $\{v_1, v_2, \dots, v_k\}$ is a basis for W . Then $[T]_\beta$ is a diagonal matrix with ones as the first k diagonal entries and zeros elsewhere. In fact, $[T]_\beta$ has the form

$$\begin{pmatrix} I_k & O_1 \\ O_2 & O_3 \end{pmatrix}.$$

If U is any projection on W , we may choose a basis γ for V such that $[U]_\gamma$ has the form above; however γ is not necessarily orthonormal.

We are now ready for the principal theorem of this section.

Theorem 6.25 (The Spectral Theorem). Suppose that T is a linear operator on a finite-dimensional inner product space V over F with the distinct eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_k$. Assume that T is normal if $F = C$ and that T is self-adjoint if $F = R$. For each i ($1 \leq i \leq k$), let W_i be the eigenspace of T corresponding to the eigenvalue λ_i , and let T_i be the orthogonal projection of V on W_i . Then the following statements are true.

- (a) $V = W_1 \oplus W_2 \oplus \cdots \oplus W_k$.
- (b) If W'_i denotes the direct sum of the subspaces W_j for $j \neq i$, then $W_i^\perp = W'_i$.
- (c) $T_i T_j = \delta_{ij} T_i$ for $1 \leq i, j \leq k$.
- (d) $I = T_1 + T_2 + \cdots + T_k$.
- (e) $T = \lambda_1 T_1 + \lambda_2 T_2 + \cdots + \lambda_k T_k$.

Proof. (a) By Theorems 6.16 (p. 369) and 6.17 (p. 371), T is diagonalizable; so

$$V = W_1 \oplus W_2 \oplus \cdots \oplus W_k$$

by Theorem 5.10 (p. 277).

(b) If $x \in W_i$ and $y \in W_j$ for some $i \neq j$, then $\langle x, y \rangle = 0$ by Theorem 6.15(d) (p. 368). It follows easily from this result that $W'_i \subseteq W_i^\perp$. From (a), we have

$$\dim(W'_i) = \sum_{j \neq i} \dim(W_j) = \dim(V) - \dim(W_i).$$

On the other hand, we have $\dim(W_i^\perp) = \dim(V) - \dim(W_i)$ by Theorem 6.7(c) (p. 349). Hence $W'_i = W_i^\perp$, proving (b).

(c) The proof of (c) is left as an exercise.

(d) Since T_i is the orthogonal projection of V on W_i , it follows from (b) that $N(T_i) = R(T_i)^\perp = W_i^\perp = W'_i$. Hence, for $x \in V$, we have $x = x_1 + x_2 + \cdots + x_k$, where $T_i(x) = x_i \in W_i$, proving (d).

(e) For $x \in V$, write $x = x_1 + x_2 + \cdots + x_k$, where $x_i \in W_i$. Then

$$T(x) = T(x_1) + T(x_2) + \cdots + T(x_k)$$

$$\begin{aligned}
 &= \lambda_1 x_1 + \lambda_2 x_2 + \cdots + \lambda_k x_k \\
 &= \lambda_1 T_1(x) + \lambda_2 T_2(x) + \cdots + \lambda_k T_k(x) \\
 &= (\lambda_1 T_1 + \lambda_2 T_2 + \cdots + \lambda_k T_k)(x).
 \end{aligned}$$

■

The set $\{\lambda_1, \lambda_2, \dots, \lambda_k\}$ of eigenvalues of T is called the **spectrum** of T , the sum $I = T_1 + T_2 + \cdots + T_k$ in (d) is called the **resolution of the identity operator** induced by T , and the sum $T = \lambda_1 T_1 + \lambda_2 T_2 + \cdots + \lambda_k T_k$ in (e) is called the **spectral decomposition** of T . The spectral decomposition of T is unique up to the order of its eigenvalues.

With the preceding notation, let β be the union of orthonormal bases of the W_i 's and let $m_i = \dim(W_i)$. (Thus m_i is the multiplicity of λ_i .) Then $[T]_\beta$ has the form

$$\begin{pmatrix} \lambda_1 I_{m_1} & O & \cdots & O \\ O & \lambda_2 I_{m_2} & \cdots & O \\ \vdots & \vdots & & \vdots \\ O & O & \cdots & \lambda_k I_{m_k} \end{pmatrix};$$

that is, $[T]_\beta$ is a diagonal matrix in which the diagonal entries are the eigenvalues λ_i of T , and each λ_i is repeated m_i times. If $\lambda_1 T_1 + \lambda_2 T_2 + \cdots + \lambda_k T_k$ is the spectral decomposition of T , then it follows (from Exercise 7) that $g(T) = g(\lambda_1)T_1 + g(\lambda_2)T_2 + \cdots + g(\lambda_k)T_k$ for any polynomial g . This fact is used below.

We now list several interesting corollaries of the spectral theorem; many more results are found in the exercises. For what follows, we assume that T is a linear operator on a finite-dimensional inner product space V over F .

Corollary 1. If $F = C$, then T is normal if and only if $T^* = g(T)$ for some polynomial g .

Proof. Suppose first that T is normal. Let $T = \lambda_1 T_1 + \lambda_2 T_2 + \cdots + \lambda_k T_k$ be the spectral decomposition of T . Taking the adjoint of both sides of the preceding equation, we have $T^* = \bar{\lambda}_1 T_1 + \bar{\lambda}_2 T_2 + \cdots + \bar{\lambda}_k T_k$ since each T_i is self-adjoint. Using the Lagrange interpolation formula (see page 53), we may choose a polynomial g such that $g(\lambda_i) = \bar{\lambda}_i$ for $1 \leq i \leq k$. Then

$$g(T) = g(\lambda_1)T_1 + g(\lambda_2)T_2 + \cdots + g(\lambda_k)T_k = \bar{\lambda}_1 T_1 + \bar{\lambda}_2 T_2 + \cdots + \bar{\lambda}_k T_k = T^*.$$

Conversely, if $T^* = g(T)$ for some polynomial g , then T commutes with T^* since T commutes with every polynomial in T . So T is normal. ■

Corollary 2. If $F = C$, then T is unitary if and only if T is normal and $|\lambda| = 1$ for every eigenvalue λ of T .

Proof. If T is unitary, then T is normal and every eigenvalue of T has absolute value 1 by Corollary 2 to Theorem 6.18 (p. 379).

Let $T = \lambda_1 T_1 + \lambda_2 T_2 + \cdots + \lambda_k T_k$ be the spectral decomposition of T . If $|\lambda| = 1$ for every eigenvalue λ of T , then by (c) of the spectral theorem,

$$\begin{aligned} TT^* &= (\lambda_1 T_1 + \lambda_2 T_2 + \cdots + \lambda_k T_k)(\bar{\lambda}_1 T_1 + \bar{\lambda}_2 T_2 + \cdots + \bar{\lambda}_k T_k) \\ &= |\lambda_1|^2 T_1 + |\lambda_2|^2 T_2 + \cdots + |\lambda_k|^2 T_k \\ &= T_1 + T_2 + \cdots + T_k \\ &= I. \end{aligned}$$

Hence T is unitary. ■

Corollary 3. *If $F = C$, then T is self-adjoint if and only if T is normal and every eigenvalue of T is real.*

Proof. Suppose that T is normal and that its eigenvalues are real. Let $T = \lambda_1 T_1 + \lambda_2 T_2 + \cdots + \lambda_k T_k$ be the spectral decomposition of T . Then

$$T^* = \bar{\lambda}_1 T_1 + \bar{\lambda}_2 T_2 + \cdots + \bar{\lambda}_k T_k = \lambda_1 T_1 + \lambda_2 T_2 + \cdots + \lambda_k T_k = T.$$

Hence T is self-adjoint.

Now suppose that T is self-adjoint and hence normal. That its eigenvalues are real has been proved in the lemma to Theorem 6.17 (p. 371). ■

Corollary 4. *Let T be as in the spectral theorem with spectral decomposition $T = \lambda_1 T_1 + \lambda_2 T_2 + \cdots + \lambda_k T_k$. Then each T_j is a polynomial in T .*

Proof. Choose a polynomial g_j ($1 \leq j \leq k$) such that $g_j(\lambda_i) = \delta_{ij}$. Then

$$\begin{aligned} g_j(T) &= g_j(\lambda_1)T_1 + g_j(\lambda_2)T_2 + \cdots + g_j(\lambda_k)T_k \\ &= \delta_{1j} T_1 + \delta_{2j} T_2 + \cdots + \delta_{kj} T_k = T_j. \end{aligned}$$
■

EXERCISES

- Label the following statements as true or false. Assume that the underlying inner product spaces are finite-dimensional.
 - All projections are self-adjoint.
 - An orthogonal projection is uniquely determined by its range.
 - Every self-adjoint operator is a linear combination of orthogonal projections.
 - If T is a projection on W , then $T(x)$ is the vector in W that is closest to x .

- (e) Every orthogonal projection is a unitary operator.
2. Let $V = \mathbb{R}^2$, $W = \text{span}(\{(1, 2)\})$, and β be the standard ordered basis for V . Compute $[T]_\beta$, where T is the orthogonal projection of V on W . Do the same for $V = \mathbb{R}^3$ and $W = \text{span}(\{(1, 0, 1)\})$.
3. For each of the matrices A in Exercise 2 of Section 6.5:
- (1) Verify that L_A possesses a spectral decomposition.
 - (2) For each eigenvalue of L_A , explicitly define the orthogonal projection on the corresponding eigenspace.
 - (3) Verify your results using the spectral theorem.
4. Let W be a finite-dimensional subspace of an inner product space V . Show that if T is the orthogonal projection of V on W , then $I - T$ is the orthogonal projection of V on W^\perp .
5. Let T be a linear operator on a finite-dimensional inner product space V .
- (a) If T is an orthogonal projection, prove that $\|T(x)\| \leq \|x\|$ for all $x \in V$. Give an example of a projection for which this inequality does not hold. What can be concluded about a projection for which the inequality is actually an equality for all $x \in V$?
 - (b) Suppose that T is a projection such that $\|T(x)\| \leq \|x\|$ for all $x \in V$. Prove that T is an orthogonal projection.
6. Let T be a normal operator on a finite-dimensional inner product space. Prove that if T is a projection, then T is also an orthogonal projection.
7. Let T be a normal operator on a finite-dimensional complex inner product space V . Use the spectral decomposition $\lambda_1 T_1 + \lambda_2 T_2 + \cdots + \lambda_k T_k$ of T to prove the following results.
- (a) If g is a polynomial, then
$$g(T) = \sum_{i=1}^k g(\lambda_i) T_i.$$
 - (b) If $T^n = T_0$ for some n , then $T = T_0$.
 - (c) Let U be a linear operator on V . Then U commutes with T if and only if U commutes with each T_i .
 - (d) There exists a normal operator U on V such that $U^2 = T$.
 - (e) T is invertible if and only if $\lambda_i \neq 0$ for $1 \leq i \leq k$.
 - (f) T is a projection if and only if every eigenvalue of T is 1 or 0.
 - (g) $T = -T^*$ if and only if every λ_i is an imaginary number.

8. Use Corollary 1 of the spectral theorem to show that if T is a normal operator on a complex finite-dimensional inner product space and U is a linear operator that commutes with T , then U commutes with T^* .
9. Referring to Exercise 20 of Section 6.5, prove the following facts about a partial isometry U .
 - (a) U^*U is an orthogonal projection on W .
 - (b) $UU^*U = U$.
10. *Simultaneous diagonalization.* Let U and T be normal operators on a finite-dimensional complex inner product space V such that $TU = UT$. Prove that there exists an orthonormal basis for V consisting of vectors that are eigenvectors of both T and U . *Hint:* Use the hint of Exercise 14 of Section 6.4 along with Exercise 8.
11. Prove (c) of the spectral theorem. Visit goo.gl/utQ9Pb for a solution.

6.7* THE SINGULAR VALUE DECOMPOSITION AND THE PSEUDOINVERSE

In Section 6.4, we characterized normal operators on complex spaces and self-adjoint operators on real spaces in terms of orthonormal bases of eigenvectors and their corresponding eigenvalues (Theorems 6.16, p. 369, and 6.17, p. 371). In this section, we establish a comparable theorem whose scope is the entire class of linear transformations on both complex and real finite-dimensional inner product spaces—the *singular value theorem for linear transformations* (Theorem 6.26). There are similarities and differences among these theorems. All rely on the use of orthonormal bases and numerical invariants. However, because of its general scope, the singular value theorem is concerned with two (usually distinct) inner product spaces and with two (usually distinct) orthonormal bases. If the two spaces and the two bases are identical, then the transformation would, in fact, be a normal or self-adjoint operator. Another difference is that the numerical invariants in the singular value theorem, the *singular values*, are nonnegative, in contrast to their counterparts, the eigenvalues, for which there is no such restriction. This property is necessary to guarantee the uniqueness of singular values.

The singular value theorem encompasses both real and complex spaces. For brevity, in this section we use the terms *unitary operator* and *unitary matrix* to include orthogonal operators and orthogonal matrices in the context of real spaces. Thus any operator T for which $\langle T(x), T(y) \rangle = \langle x, y \rangle$, or any matrix A for which $\langle Ax, Ay \rangle = \langle x, y \rangle$, for all x and y is called *unitary* for the purposes of this section.

In Exercise 15 of Section 6.3, the definition of the adjoint of an operator is extended to any linear transformation $T: V \rightarrow W$, where V and W are

finite-dimensional inner product spaces. By this exercise, the adjoint T^* of T is a linear transformation from W to V and $[T^*]_{\gamma}^{\beta} = ([T]_{\beta}^{\gamma})^*$, where β and γ are orthonormal bases for V and W , respectively. Furthermore, the linear operator T^*T on V is positive semidefinite and $\text{rank}(T^*T) = \text{rank}(T)$ by Exercise 18 of Section 6.4.

With these facts in mind, we begin with the principal result.

Theorem 6.26 (Singular Value Theorem for Linear Transformations). *Let V and W be finite-dimensional inner product spaces, and let $T: V \rightarrow W$ be a linear transformation of rank r . Then there exist orthonormal bases $\{v_1, v_2, \dots, v_n\}$ for V and $\{u_1, u_2, \dots, u_m\}$ for W and positive scalars $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r$ such that*

$$T(v_i) = \begin{cases} \sigma_i u_i & \text{if } 1 \leq i \leq r \\ 0 & \text{if } i > r. \end{cases} \quad (4)$$

Conversely, suppose that the preceding conditions are satisfied. Then for $1 \leq i \leq n$, v_i is an eigenvector of T^*T with corresponding eigenvalue σ_i^2 if $1 \leq i \leq r$ and 0 if $i > r$. Therefore the scalars $\sigma_1, \sigma_2, \dots, \sigma_r$ are uniquely determined by T .

Proof. We first establish the existence of the bases and scalars. By Exercises 18 of Section 6.4 and 15(d) of Section 6.3, T^*T is a positive semidefinite linear operator of rank r on V ; hence there is an orthonormal basis $\{v_1, v_2, \dots, v_n\}$ for V consisting of eigenvectors of T^*T with corresponding eigenvalues λ_i , where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0$, and $\lambda_i = 0$ for $i > r$. For $1 \leq i \leq r$, define $\sigma_i = \sqrt{\lambda_i}$ and $u_i = \frac{1}{\sigma_i} T(v_i)$. We show that $\{u_1, u_2, \dots, u_r\}$ is an orthonormal subset of W . Suppose $1 \leq i, j \leq r$. Then

$$\begin{aligned} \langle u_i, u_j \rangle &= \left\langle \frac{1}{\sigma_i} T(v_i), \frac{1}{\sigma_j} T(v_j) \right\rangle \\ &= \frac{1}{\sigma_i \sigma_j} \langle T^*T(v_i), v_j \rangle \\ &= \frac{1}{\sigma_i \sigma_j} \langle \lambda_i v_i, v_j \rangle \\ &= \frac{\sigma_i^2}{\sigma_i \sigma_j} \langle v_i, v_j \rangle \\ &= \delta_{ij}, \end{aligned}$$

and hence $\{u_1, u_2, \dots, u_r\}$ is orthonormal. By Theorem 6.7(a) (p. 349), this set extends to an orthonormal basis $\{u_1, u_2, \dots, u_r, \dots, u_m\}$ for W . Clearly

$T(v_i) = \sigma_i u_i$ if $1 \leq i \leq r$. If $i > r$, then $T^*T(v_i) = 0$, and so $T(v_i) = 0$ by Exercise 15(d) of Section 6.3.

To establish uniqueness, suppose that $\{v_1, v_2, \dots, v_n\}$, $\{u_1, u_2, \dots, u_m\}$, and $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$ satisfy the properties stated in the first part of the theorem. Then for $1 \leq i \leq m$ and $1 \leq j \leq n$,

$$\begin{aligned}\langle T^*(u_i), v_j \rangle &= \langle u_i, T(v_j) \rangle \\ &= \begin{cases} \sigma_i & \text{if } i = j \leq r \\ 0 & \text{otherwise,} \end{cases}\end{aligned}$$

and hence for any $1 \leq i \leq m$,

$$T^*(u_i) = \sum_{j=1}^n \langle T^*(u_i), v_j \rangle v_j = \begin{cases} \sigma_i v_i & \text{if } i = j \leq r \\ 0 & \text{otherwise.} \end{cases} \quad (5)$$

So for $i \leq r$,

$$T^*T(v_i) = T^*(\sigma_i u_i) = \sigma_i T^*(u_i) = \sigma_i^2 v_i$$

and $T^*T(v_i) = T^*(0) = 0$ for $i > r$. Therefore each v_i is an eigenvector of T^*T with corresponding eigenvalue σ_i^2 if $i \leq r$ and 0 if $i > r$. ■

Definition. The unique scalars $\sigma_1, \sigma_2, \dots, \sigma_r$ in Theorem 6.26 are called the **singular values** of T . If r is less than both m and n , then the term singular value is extended to include $\sigma_{r+1} = \dots = \sigma_k = 0$, where k is the minimum of m and n .

Although the singular values of a linear transformation T are uniquely determined by T , the orthonormal bases given in the statement of Theorem 6.26 are not uniquely determined because there is more than one orthonormal basis of eigenvectors of T^*T .

In view of (5), the singular values of a linear transformation $T: V \rightarrow W$ and its adjoint T^* are identical. Furthermore, the orthonormal bases for V and W given in Theorem 6.26 are simply reversed for T^* .

Example 1

Let $P_2(R)$ and $P_1(R)$ be the polynomial spaces with inner products defined by

$$\langle f(x), g(x) \rangle = \int_{-1}^1 f(t)g(t) dt.$$

Let $T: P_2(R) \rightarrow P_1(R)$ be the linear transformation defined by $T(f(x)) = f'(x)$. Find orthonormal bases $\beta = \{v_1, v_2, v_3\}$ for $P_2(R)$ and $\gamma = \{u_1, u_2\}$ for $P_1(R)$ such that $T(v_i) = \sigma_i u_i$ for $i = 1, 2$ and $T(v_3) = 0$, where $\sigma_1 \geq \sigma_2 > 0$ are the nonzero singular values of T .

To facilitate the computations, we translate this problem into the corresponding problem for a matrix representation of T . Caution is advised here because not any matrix representation will do. Since the adjoint is defined in terms of inner products, we must use a matrix representation constructed from orthonormal bases for $P_2(R)$ and $P_1(R)$ to guarantee that the adjoint of the matrix representation of T is the same as the matrix representation of the adjoint of T . (See Exercise 15 of Section 6.3.) For this purpose, we use the results of Exercise 21(a) of Section 6.2 to obtain orthonormal bases

$$\alpha = \left\{ \frac{1}{\sqrt{2}}, \sqrt{\frac{3}{2}}x, \sqrt{\frac{5}{8}}(3x^2 - 1) \right\} \quad \text{and} \quad \alpha' = \left\{ \frac{1}{\sqrt{2}}, \sqrt{\frac{3}{2}}x \right\}$$

for $P_2(R)$ and $P_1(R)$, respectively.

Let

$$A = [T]_{\alpha'}^{\alpha} = \begin{pmatrix} 0 & \sqrt{3} & 0 \\ 0 & 0 & \sqrt{15} \end{pmatrix}.$$

Then

$$A^* A = \begin{pmatrix} 0 & 0 \\ \sqrt{3} & 0 \\ 0 & \sqrt{15} \end{pmatrix} \begin{pmatrix} 0 & \sqrt{3} & 0 \\ 0 & 0 & \sqrt{15} \end{pmatrix} = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 15 \end{pmatrix},$$

which has eigenvalues (listed in descending order of size) $\lambda_1 = 15$, $\lambda_2 = 3$, and $\lambda_3 = 0$. These eigenvalues correspond, respectively, to the orthonormal eigenvectors $e_3 = (0, 0, 1)$, $e_2 = (0, 1, 0)$, and $e_1 = (1, 0, 0)$ in \mathbb{R}^3 . Translating everything into the context of T , $P_2(R)$, and $P_1(R)$, let

$$v_1 = \sqrt{\frac{5}{8}}(3x^2 - 1), \quad v_2 = \sqrt{\frac{3}{2}}x, \quad \text{and} \quad v_3 = \frac{1}{\sqrt{2}}.$$

Then $\beta = \{v_1, v_2, v_3\}$ is an orthonormal basis for $P_2(R)$ consisting of eigenvectors of $T^* T$ with corresponding eigenvalues λ_1 , λ_2 , and λ_3 . Now set $\sigma_1 = \sqrt{\lambda_1} = \sqrt{15}$ and $\sigma_2 = \sqrt{\lambda_2} = \sqrt{3}$, the nonzero singular values of T , and take

$$u_1 = \frac{1}{\sigma_1} T(v_1) = \sqrt{\frac{3}{2}}x \quad \text{and} \quad u_2 = \frac{1}{\sigma_2} T(v_2) = \frac{1}{\sqrt{2}},$$

to obtain the required basis $\gamma = \{u_1, u_2\}$ for $\mathbb{P}_1(R)$. \blacklozenge

We can use singular values to describe how a figure is distorted by a linear transformation. This is illustrated in the next example.

Example 2

Let T be an invertible linear operator on \mathbb{R}^2 and $S = \{x \in \mathbb{R}^2 : \|x\| = 1\}$, the unit circle in \mathbb{R}^2 . We apply Theorem 6.26 to describe $S' = T(S)$.

Since T is invertible, it has rank equal to 2 and hence has singular values $\sigma_1 \geq \sigma_2 > 0$. Let $\{v_1, v_2\}$ and $\beta = \{u_1, u_2\}$ be orthonormal bases for \mathbb{R}^2 so that $T(v_1) = \sigma_1 u_1$ and $T(v_2) = \sigma_2 u_2$, as in Theorem 6.26. Then β determines a coordinate system, which we will call the $x'y'$ -coordinate system for \mathbb{R}^2 , where the x' -axis contains u_1 and the y' -axis contains u_2 . For any vector $u \in \mathbb{R}^2$, if $u = x'_1 u_1 + x'_2 u_2$, then $[u]_\beta = \begin{pmatrix} x'_1 \\ x'_2 \end{pmatrix}$ is the coordinate vector of u relative to β .

We characterize S' in terms of an equation relating x'_1 and x'_2 .

For any vector $v = x_1 v_1 + x_2 v_2 \in \mathbb{R}^2$, the equation $u = T(v)$ means that

$$u = T(x_1 v_1 + x_2 v_2) = x_1 T(v_1) + x_2 T(v_2) = x_1 \sigma_1 u_1 + x_2 \sigma_2 u_2.$$

Thus for $u = x'_1 u_1 + x'_2 u_2$, we have $x'_1 = x_1 \sigma_1$ and $x'_2 = x_2 \sigma_2$. Furthermore, $u \in S'$ if and only if $v \in S$ if and only if

$$\frac{(x'_1)^2}{\sigma_1^2} + \frac{(x'_2)^2}{\sigma_2^2} = x_1^2 + x_2^2 = 1.$$

If $\sigma_1 = \sigma_2$, this is the equation of a circle of radius σ_1 , and if $\sigma_1 > \sigma_2$, this is the equation of an ellipse with major axis and minor axis oriented along the x' -axis and the y' -axis, respectively. (See Figure 6.6.) \blacklozenge

The singular value theorem for linear transformations is useful in its matrix form because we can perform numerical computations on matrices. We begin with the definition of the singular values of a matrix.

Definition. Let A be an $m \times n$ matrix. We define the **singular values** of A to be the singular values of the linear transformation L_A .

Theorem 6.27 (Singular Value Decomposition Theorem for Matrices). Let A be an $m \times n$ matrix of rank r with the positive singular values $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r$, and let Σ be the $m \times n$ matrix defined by

$$\Sigma_{ij} = \begin{cases} \sigma_i & \text{if } i = j \leq r \\ 0 & \text{otherwise.} \end{cases}$$

Then there exists an $m \times m$ unitary matrix U and an $n \times n$ unitary matrix V such that

$$A = U \Sigma V^*.$$

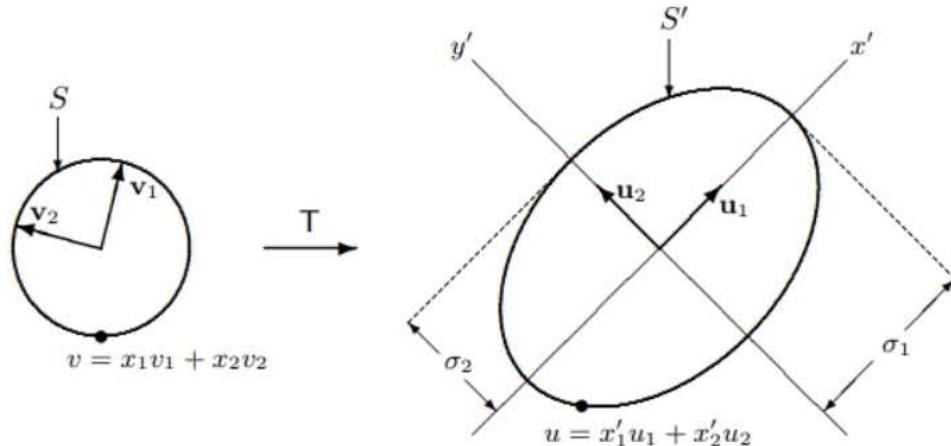


Figure 6.6

Proof. Let $T = L_A: \mathbb{F}^n \rightarrow \mathbb{F}^m$. By Theorem 6.26, there exist orthonormal bases $\beta = \{v_1, v_2, \dots, v_n\}$ for \mathbb{F}^n and $\gamma = \{u_1, u_2, \dots, u_m\}$ for \mathbb{F}^m such that $T(v_i) = \sigma_i u_i$ for $1 \leq i \leq r$ and $T(v_i) = 0$ for $i > r$. Let U be the $m \times m$ matrix whose j th column is u_j for all j , and let V be the $n \times n$ matrix whose j th column is v_j for all j . Note that both U and V are unitary matrices.

By Theorem 2.13(a) (p. 91), the j th column of AV is $Av_j = \sigma_j u_j$. Observe that the j th column of Σ is $\sigma_j e_j$, where e_j is the j th standard vector of \mathbb{F}^m . So by Theorem 2.13(a) and (b), the j th column of $U\Sigma$ is given by

$$U(\sigma_j e_j) = \sigma_j U e_j = \sigma_j u_j.$$

It follows that AV and $U\Sigma$ are $m \times n$ matrices whose corresponding columns are equal, and hence $AV = U\Sigma$. Therefore $A = AVV^* = U\Sigma V^*$. ■

Definition. Let A be an $m \times n$ matrix of rank r with positive singular values $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r$. A factorization $A = U\Sigma V^*$ where U and V are unitary matrices and Σ is the $m \times n$ matrix defined as in Theorem 6.27 is called a **singular value decomposition** of A .

In the proof of Theorem 6.27, the columns of V are the vectors in β , and the columns of U are the vectors in γ . Furthermore, the nonzero singular values of A are the same as those of L_A ; hence they are the square roots of the nonzero eigenvalues of A^*A or of AA^* . (See Exercise 9.)

Example 3

We find a singular value decomposition for $A = \begin{pmatrix} 1 & 1 & -1 \\ 1 & 1 & -1 \end{pmatrix}$.

First observe that for

$$v_1 = \frac{1}{\sqrt{3}} \begin{pmatrix} 1 \\ 1 \\ -1 \end{pmatrix}, \quad v_2 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ -1 \\ 0 \end{pmatrix}, \quad \text{and} \quad v_3 = \frac{1}{\sqrt{6}} \begin{pmatrix} 1 \\ 1 \\ 2 \end{pmatrix},$$

the set $\beta = \{v_1, v_2, v_3\}$ is an orthonormal basis for \mathbb{R}^3 consisting of eigenvectors of A^*A with corresponding eigenvalues $\lambda_1 = 6$, and $\lambda_2 = \lambda_3 = 0$. Consequently, $\sigma_1 = \sqrt{6}$ is the only nonzero singular value of A . Hence, as in the proof of Theorem 6.27, we let V be the matrix whose columns are the vectors in β . Then

$$\Sigma = \begin{pmatrix} \sqrt{6} & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} \quad \text{and} \quad V = \begin{pmatrix} \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{6}} \\ \frac{1}{\sqrt{3}} & \frac{-1}{\sqrt{2}} & \frac{1}{\sqrt{6}} \\ \frac{-1}{\sqrt{3}} & 0 & \frac{2}{\sqrt{6}} \end{pmatrix}.$$

Also, as in Theorem 6.27, we take

$$u_1 = \frac{1}{\sigma_i} L_A(v_1) = \frac{1}{\sigma_i} Av_1 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ 1 \end{pmatrix}.$$

Next choose $u_2 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ -1 \end{pmatrix}$, a unit vector orthogonal to u_1 , to obtain the orthonormal basis $\gamma = \{u_1, u_2\}$ for \mathbb{R}^2 , and set

$$U = \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{-1}{\sqrt{2}} \end{pmatrix}.$$

Then $A = U\Sigma V^*$ is the desired singular value decomposition. ♦

The Polar Decomposition of a Square Matrix

A singular value decomposition of a matrix can be used to factor a square matrix in a manner analogous to the factoring of a complex number as the product of a complex number of length 1 and a nonnegative number. In the case of matrices, the complex number of length 1 is replaced by a unitary matrix, and the nonnegative number is replaced by a positive semidefinite matrix.

Theorem 6.28 (Polar Decomposition). *For any square matrix A , there exists a unitary matrix W and a positive semidefinite matrix P such that*

$$A = WP.$$

Furthermore, if A is invertible, then the representation is unique.

Proof. By Theorem 6.27, there exist unitary matrices U and V and a diagonal matrix Σ with nonnegative diagonal entries such that $A = U\Sigma V^*$. So

$$A = U\Sigma V^* = UV^*V\Sigma V^* = WP,$$

where $W = UV^*$ and $P = V\Sigma V^*$. Since W is the product of unitary matrices, W is unitary, and since Σ is positive semidefinite and P is unitarily equivalent to Σ , P is positive semidefinite by Exercise 14 of Section 6.5.

Now suppose that A is invertible and factors as the products

$$A = WP = ZQ,$$

where W and Z are unitary and P and Q are positive semidefinite. Since A is invertible, it follows that P and Q are positive definite and invertible, and therefore $Z^*W = QP^{-1}$. Thus QP^{-1} is unitary, and so

$$I = (QP^{-1})^*(QP^{-1}) = P^{-1}Q^2P^{-1}.$$

Hence $P^2 = Q^2$. Since both P and Q are positive definite, it follows that $P = Q$ by Exercise 17 of Section 6.4. Therefore $W = Z$, and consequently the factorization is unique. ■

The factorization of a square matrix A as WP where W is unitary and P is positive semidefinite, is called a **polar decomposition** of A .

Example 4

To find the polar decomposition of $A = \begin{pmatrix} 11 & -5 \\ -2 & 10 \end{pmatrix}$, we begin by finding a singular value decomposition $U\Sigma V^*$ of A . The object is to find an orthonormal basis β for \mathbb{R}^2 consisting of eigenvectors of A^*A . It can be shown that

$$v_1 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ -1 \end{pmatrix} \quad \text{and} \quad v_2 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

are orthonormal eigenvectors of A^*A with corresponding eigenvalues $\lambda_1 = 200$ and $\lambda_2 = 50$. So $\beta = \{v_1, v_2\}$ is an appropriate basis. Thus $\sigma_1 = \sqrt{200} = 10\sqrt{2}$ and $\sigma_2 = \sqrt{50} = 5\sqrt{2}$ are the singular values of A . So we have

$$V = \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{-1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix} \quad \text{and} \quad \Sigma = \begin{pmatrix} 10\sqrt{2} & 0 \\ 0 & 5\sqrt{2} \end{pmatrix}.$$

Next, we find the columns u_1 and u_2 of U :

$$u_1 = \frac{1}{\sigma_1} Av_1 = \frac{1}{5} \begin{pmatrix} 4 \\ -3 \end{pmatrix} \quad \text{and} \quad u_2 = \frac{1}{\sigma_2} Av_2 = \frac{1}{5} \begin{pmatrix} 3 \\ 4 \end{pmatrix}.$$

Thus

$$U = \begin{pmatrix} \frac{4}{5} & \frac{3}{5} \\ -\frac{3}{5} & \frac{4}{5} \end{pmatrix}.$$

Therefore, in the notation of Theorem 6.28, we have

$$W = UV^* = \begin{pmatrix} \frac{4}{5} & \frac{3}{5} \\ -\frac{3}{5} & \frac{4}{5} \end{pmatrix} \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{-1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix} = \frac{1}{5\sqrt{2}} \begin{pmatrix} 7 & -1 \\ 1 & 7 \end{pmatrix},$$

and

$$P = V\Sigma V^* = \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{-1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix} \begin{pmatrix} 10\sqrt{2} & 0 \\ 0 & 5\sqrt{2} \end{pmatrix} \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{-1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix} = \frac{5}{\sqrt{2}} \begin{pmatrix} 3 & -1 \\ -1 & 3 \end{pmatrix}.$$

◆

The Pseudoinverse

Let V and W be finite-dimensional inner product spaces over the same field, and let $T: V \rightarrow W$ be a linear transformation. It is desirable to have a linear transformation from W to V that captures some of the essence of an inverse of T even if T is not invertible. A simple approach to this problem is to focus on the “part” of T that is invertible, namely, the restriction of T to $N(T)^\perp$. Let $L: N(T)^\perp \rightarrow R(T)$ be the linear transformation defined by $L(x) = T(x)$ for all $x \in N(T)^\perp$. Then L is invertible, and we can use the inverse of L to construct a linear transformation from W to V that salvages some of the benefits of an inverse of T .

Definition. Let V and W be finite-dimensional inner product spaces over the same field, and let $T: V \rightarrow W$ be a linear transformation. Let $L: N(T)^\perp \rightarrow R(T)$ be the linear transformation defined by $L(x) = T(x)$ for all $x \in N(T)^\perp$. The **pseudoinverse** (or Moore-Penrose generalized inverse) of T , denoted by T^\dagger , is defined as the unique linear transformation from W to V such that

$$T^\dagger(y) = \begin{cases} L^{-1}(y) & \text{for } y \in R(T) \\ 0 & \text{for } y \in R(T)^\perp. \end{cases}$$

Note that $L^{-1}T(x) = x$ for all $x \in N(T)^\perp$.

The pseudoinverse of a linear transformation T on a finite-dimensional inner product space exists even if T is not invertible. Furthermore, if T is invertible, then $T^\dagger = T^{-1}$ because $N(T)^\perp = V$, and L (as defined above) coincides with T .

As an extreme example, consider the zero transformation $T_0: V \rightarrow W$ between two finite-dimensional inner product spaces V and W . Then $R(T_0) = \{0\}$, and therefore T^\dagger is the zero transformation from W to V .

We can use the singular value theorem to describe the pseudoinverse of a linear transformation. Suppose that V and W are finite-dimensional vector spaces and $T: V \rightarrow W$ is a linear transformation of rank r . Let $\{v_1, v_2, \dots, v_n\}$ and $\{u_1, u_2, \dots, u_m\}$ be orthonormal bases for V and W , respectively, and let $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r$ be the nonzero singular values of T satisfying (4) in Theorem 6.26. Then $\{v_1, v_2, \dots, v_r\}$ is a basis for $N(T)^\perp$, $\{v_{r+1}, v_{r+2}, \dots, v_n\}$ is a basis for $N(T)$, $\{u_1, u_2, \dots, u_r\}$ is a basis for $R(T)$, and $\{u_{r+1}, u_{r+2}, \dots, u_m\}$ is a basis for $R(T)^\perp$. Let L be the restriction of T to $N(T)^\perp$, as in the definition of pseudoinverse. Then $L^{-1}(u_i) = \frac{1}{\sigma_i}v_i$ for $1 \leq i \leq r$. Therefore

$$T^\dagger(u_i) = \begin{cases} \frac{1}{\sigma_i}v_i & \text{if } 1 \leq i \leq r \\ 0 & \text{if } r < i \leq m. \end{cases} \quad (6)$$

Example 5

Let $T: P_2(R) \rightarrow P_1(R)$ be the linear transformation defined by $T(f(x)) = f'(x)$, as in Example 1. Let $\beta = \{v_1, v_2, v_3\}$ and $\gamma = \{u_1, u_2\}$ be the orthonormal bases for $P_2(R)$ and $P_1(R)$ in Example 1. Then $\sigma_1 = \sqrt{15}$ and $\sigma_2 = \sqrt{3}$ are the nonzero singular values of T . It follows that

$$T^\dagger\left(\sqrt{\frac{3}{2}}x\right) = T^\dagger(u_1) = \frac{1}{\sigma_1}v_1 = \frac{1}{\sqrt{15}}\sqrt{\frac{5}{8}}(3x^2 - 1),$$

and hence

$$T^\dagger(x) = \frac{1}{6}(3x^2 - 1).$$

Similarly, $T^\dagger(1) = x$. Thus, for any polynomial $a + bx \in P_1(R)$,

$$T^\dagger(a + bx) = aT^\dagger(1) + bT^\dagger(x) = ax + \frac{b}{6}(3x^2 - 1). \quad \blacklozenge$$

The Pseudoinverse of a Matrix

Let A be an $m \times n$ matrix. Then there exists a unique $n \times m$ matrix B such that $(L_A)^\dagger: F^m \rightarrow F^n$ is equal to the left-multiplication transformation L_B . We call B the **pseudoinverse** of A and denote it by $B = A^\dagger$. Thus

$$(L_A)^\dagger = L_{A^\dagger}.$$

Let A be an $m \times n$ matrix of rank r . The pseudoinverse of A can be computed with the aid of a singular value decomposition $A = U\Sigma V^*$. Let β and γ be the ordered bases whose vectors are the columns of V and U , respectively, and let $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r$ be the nonzero singular values of

A. Then β and γ are orthonormal bases for F^n and F^m , respectively, and (4) and (6) are satisfied for $T = L_A$. Reversing the roles of β and γ in the proof of Theorem 6.27, we obtain the following result.

Theorem 6.29. *Let A be an $m \times n$ matrix of rank r with a singular value decomposition $A = U\Sigma V^*$ and nonzero singular values $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r$. Let Σ^\dagger be the $n \times m$ matrix defined by*

$$\Sigma_{ij}^\dagger = \begin{cases} \frac{1}{\sigma_i} & \text{if } i = j \leq r \\ 0 & \text{otherwise.} \end{cases}$$

Then $A^\dagger = V\Sigma^\dagger U^$, and this is a singular value decomposition of A^\dagger .*

Notice that Σ^\dagger as defined in Theorem 6.29 is actually the pseudoinverse of Σ .

Example 6

We find A^\dagger for the matrix $A = \begin{pmatrix} 1 & 1 & -1 \\ 1 & 1 & -1 \end{pmatrix}$.

Since A is the matrix of Example 3, we can use the singular value decomposition obtained in that example:

$$A = U\Sigma V^* = \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{-1}{\sqrt{2}} \end{pmatrix} \begin{pmatrix} \sqrt{6} & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{6}} \\ \frac{1}{\sqrt{3}} & \frac{-1}{\sqrt{2}} & \frac{1}{\sqrt{6}} \\ \frac{-1}{\sqrt{3}} & 0 & \frac{2}{\sqrt{6}} \end{pmatrix}^*.$$

By Theorem 6.29, we have

$$A^\dagger = V\Sigma^\dagger U^* = \begin{pmatrix} \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{6}} \\ \frac{1}{\sqrt{3}} & \frac{-1}{\sqrt{2}} & \frac{1}{\sqrt{6}} \\ \frac{-1}{\sqrt{3}} & 0 & \frac{2}{\sqrt{6}} \end{pmatrix} \begin{pmatrix} \frac{1}{\sqrt{6}} & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{-1}{\sqrt{2}} \end{pmatrix} = \frac{1}{6} \begin{pmatrix} 1 & 1 \\ 1 & 1 \\ -1 & -1 \end{pmatrix}.$$



Notice that the linear transformation T of Example 5 is L_A , where A is the matrix of Example 6, and that $T^\dagger = L_{A^\dagger}$.

The Pseudoinverse and Systems of Linear Equations

Let A be an $m \times n$ matrix with entries in F . Then for any $b \in F^m$, the matrix equation $Ax = b$ is a system of linear equations, and so it either has no solutions, a unique solution, or infinitely many solutions. We know that the system has a unique solution for every $b \in F^m$ if and only if A is invertible,

in which case the solution is given by $A^{-1}b$. Furthermore, if A is invertible, then $A^{-1} = A^\dagger$, and so the solution can be written as $x = A^\dagger b$. If, on the other hand, A is not invertible or the system $Ax = b$ is inconsistent, then $A^\dagger b$ still exists. We therefore pose the following question: In general, how is the vector $A^\dagger b$ related to the system of linear equations $Ax = b$?

In order to answer this question, we need the following lemma.

Lemma. Let V and W be finite-dimensional inner product spaces, and let $T: V \rightarrow W$ be linear. Then

- (a) $T^\dagger T$ is the orthogonal projection of V on $N(T)^\perp$.
- (b) TT^\dagger is the orthogonal projection of W on $R(T)$.

Proof. As in the earlier discussion, we define $L: N(T)^\perp \rightarrow W$ by $L(x) = T(x)$ for all $x \in N(T)^\perp$. If $x \in N(T)^\perp$, then $T^\dagger T(x) = L^{-1}L(x) = x$, and if $x \in N(T)$, then $T^\dagger T(x) = T^\dagger(0) = 0$. Consequently $T^\dagger T$ is the orthogonal projection of V on $N(T)^\perp$. This proves (a).

The proof of (b) is similar and is left as an exercise. ■

Theorem 6.30. Consider the system of linear equations $Ax = b$, where A is an $m \times n$ matrix and $b \in F^m$. If $z = A^\dagger b$, then z has the following properties.

- (a) If $Ax = b$ is consistent, then z is the unique solution to the system having minimum norm. That is, z is a solution to the system, and if y is any solution to the system, then $\|z\| \leq \|y\|$ with equality if and only if $z = y$.
- (b) If $Ax = b$ is inconsistent, then z is the unique best approximation to a solution having minimum norm. That is, $\|Az - b\| \leq \|Ay - b\|$ for any $y \in F^n$, with equality if and only if $Az = Ay$. Furthermore, if $Az = Ay$, then $\|z\| \leq \|y\|$ with equality if and only if $z = y$.

Proof. For convenience, let $T = L_A$.

(a) Suppose that $Ax = b$ is consistent, and let $z = A^\dagger b$. Observe that $b \in R(T)$, and therefore $Az = AA^\dagger b = TT^\dagger(b) = b$ by part (b) of the lemma. Thus z is a solution to the system. Now suppose that y is any solution to the system. Then

$$T^\dagger T(y) = A^\dagger Ay = A^\dagger b = z,$$

and hence z is the orthogonal projection of y on $N(T)^\perp$ by part (a) of the lemma. Therefore, by the corollary to Theorem 6.6 (p. 348), we have that $\|z\| \leq \|y\|$ with equality if and only if $z = y$.

(b) Suppose that $Ax = b$ is inconsistent. By the lemma, $Az = AA^\dagger b = TT^\dagger(b) = b$ is the orthogonal projection of b on $R(T)$; therefore, by the corollary to Theorem 6.6 (p. 348), Az is the vector in $R(T)$ nearest b . That is, if Ay is any other vector in $R(T)$, then $\|Az - b\| \leq \|Ay - b\|$ with equality if and only if $Az = Ay$.

Finally, suppose that y is any vector in \mathbb{F}^n such that $Az = Ay = c$. Then

$$A^\dagger c = A^\dagger Az = A^\dagger AA^\dagger b = A^\dagger b = z$$

by Exercise 23; hence we may apply part (a) of this theorem to the system $Ax = c$ to conclude that $\|z\| \leq \|y\|$ with equality if and only if $z = y$. ■

Note that the vector $z = A^\dagger b$ in Theorem 6.30 is the vector x_0 described in Theorem 6.12 that arises in the least squares application on pages 358–361.

Example 7

Consider the linear systems

$$\begin{array}{l} x_1 + x_2 - x_3 = 1 \\ x_1 + x_2 - x_3 = 1 \end{array} \quad \text{and} \quad \begin{array}{l} x_1 + x_2 - x_3 = 1 \\ x_1 + x_2 - x_3 = 2. \end{array}$$

The first system has infinitely many solutions. Let $A = \begin{pmatrix} 1 & 1 & -1 \\ 1 & 1 & -1 \end{pmatrix}$, the coefficient matrix of the system, and let $b = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$. By Example 6,

$$A^\dagger = \frac{1}{6} \begin{pmatrix} 1 & 1 \\ 1 & 1 \\ -1 & -1 \end{pmatrix},$$

and therefore

$$z = A^\dagger b = \frac{1}{6} \begin{pmatrix} 1 & 1 \\ 1 & 1 \\ -1 & -1 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \frac{1}{3} \begin{pmatrix} 1 \\ 1 \\ -1 \end{pmatrix}$$

is the solution of minimal norm by Theorem 6.30(a).

The second system is obviously inconsistent. Let $b = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$. Thus, although

$$z = A^\dagger b = \frac{1}{6} \begin{pmatrix} 1 & 1 \\ 1 & 1 \\ -1 & -1 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \end{pmatrix} = \frac{1}{2} \begin{pmatrix} 1 \\ 1 \\ -1 \end{pmatrix}$$

is not a solution to the second system, it is the “best approximation” to a solution having minimum norm, as described in Theorem 6.30(b). ♦

EXERCISES

1. Label the following statements as true or false.

- (a) The singular values of any linear operator on a finite-dimensional vector space are also eigenvalues of the operator.
- (b) The singular values of any matrix A are the eigenvalues of A^*A .
- (c) For any matrix A and any scalar c , if σ is a singular value of A , then $|c|\sigma$ is a singular value of cA .
- (d) The singular values of any linear operator are nonnegative.
- (e) If λ is an eigenvalue of a self-adjoint matrix A , then λ is a singular value of A .
- (f) For any $m \times n$ matrix A and any $b \in \mathbb{F}^n$, the vector $A^\dagger b$ is a solution to $Ax = b$.
- (g) The pseudoinverse of any linear operator exists even if the operator is not invertible.
2. Let $T: V \rightarrow W$ be a linear transformation of rank r , where V and W are finite-dimensional inner product spaces. In each of the following, find orthonormal bases $\{v_1, v_2, \dots, v_n\}$ for V and $\{u_1, u_2, \dots, u_m\}$ for W , and the nonzero singular values $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r$ of T such that $T(v_i) = \sigma_i u_i$ for $1 \leq i \leq r$.
- (a) $T: \mathbb{R}^2 \rightarrow \mathbb{R}^3$ defined by $T(x_1, x_2) = (x_1, x_1 + x_2, x_1 - x_2)$
- (b) $T: P_2(R) \rightarrow P_1(R)$, where $T(f(x)) = f''(x)$, and the inner products are defined as in Example 1
- (c) Let $V = W = \text{span}(\{1, \sin x, \cos x\})$ with the inner product defined by $\langle f, g \rangle = \int_0^{2\pi} f(t)g(t) dt$, and T is defined by $T(f) = f' + 2f$
- (d) $T: \mathbb{C}^2 \rightarrow \mathbb{C}^2$ defined by $T(z_1, z_2) = ((1-i)z_2, (1+i)z_1 + z_2)$
3. Find a singular value decomposition for each of the following matrices.
- (a) $\begin{pmatrix} 1 & 1 \\ 1 & 1 \\ -1 & -1 \end{pmatrix}$ (b) $\begin{pmatrix} 1 & 0 & 1 \\ 1 & 0 & -1 \end{pmatrix}$ (c) $\begin{pmatrix} 1 & 1 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{pmatrix}$
- (d) $\begin{pmatrix} 1 & 1 & 1 \\ 1 & -1 & 0 \\ 1 & 0 & -1 \end{pmatrix}$ (e) $\begin{pmatrix} 1+i & 1 \\ 1-i & -i \end{pmatrix}$ (f) $\begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 0 & -2 & 1 \\ 1 & -1 & 1 & 1 \end{pmatrix}$
4. Find a polar decomposition for each of the following matrices.
- (a) $\begin{pmatrix} 1 & 1 \\ 2 & -2 \end{pmatrix}$ (b) $\begin{pmatrix} 20 & 4 & 0 \\ 0 & 0 & 1 \\ 4 & 20 & 0 \end{pmatrix}$
5. Find an explicit formula for each of the following expressions.
- (a) $T^\dagger(x_1, x_2, x_3)$, where T is the linear transformation of Exercise 2(a)
- (b) $T^\dagger(a + bx + cx^2)$, where T is the linear transformation of Exercise 2(b)

- (c) $T^\dagger(a + b \sin x + c \cos x)$, where T is the linear transformation of Exercise 2(c)
 (d) $T^\dagger(z_1, z_2)$, where T is the linear transformation of Exercise 2(d)
6. Use the results of Exercise 3 to find the pseudoinverse of each of the following matrices.
- (a) $\begin{pmatrix} 1 & 1 \\ 1 & 1 \\ -1 & -1 \end{pmatrix}$ (b) $\begin{pmatrix} 1 & 0 & 1 \\ 1 & 0 & -1 \end{pmatrix}$ (c) $\begin{pmatrix} 1 & 1 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{pmatrix}$
 (d) $\begin{pmatrix} 1 & 1 & 1 \\ 1 & -1 & 0 \\ 1 & 0 & -1 \end{pmatrix}$ (e) $\begin{pmatrix} 1+i & 1 \\ 1-i & -i \end{pmatrix}$ (f) $\begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 0 & -2 & 1 \\ 1 & -1 & 1 & 1 \end{pmatrix}$
7. For each of the given linear transformations $T: V \rightarrow W$,
- Describe the subspace Z_1 of V such that $T^\dagger T$ is the orthogonal projection of V on Z_1 .
 - Describe the subspace Z_2 of W such that TT^\dagger is the orthogonal projection of W on Z_2 .
- (a) T is the linear transformation of Exercise 2(a)
 (b) T is the linear transformation of Exercise 2(b)
 (c) T is the linear transformation of Exercise 2(c)
 (d) T is the linear transformation of Exercise 2(d)
8. For each of the given systems of linear equations,
- If the system is consistent, find the unique solution having minimum norm.
 - If the system is inconsistent, find the “best approximation to a solution” having minimum norm, as described in Theorem 6.30(b).
- (Use your answers to parts (a) and (f) of Exercise 6.)

$$(a) \begin{array}{l} x_1 + x_2 = 1 \\ x_1 + x_2 = 2 \\ -x_1 - x_2 = 0 \end{array} \quad (b) \begin{array}{l} x_1 + x_2 + x_3 + x_4 = 2 \\ x_1 - 2x_3 + x_4 = -1 \\ x_1 - x_2 + x_3 + x_4 = 2 \end{array}$$

9. Let V and W be finite-dimensional inner product spaces over F , and suppose that $\{v_1, v_2, \dots, v_n\}$ and $\{u_1, u_2, \dots, u_m\}$ are orthonormal bases for V and W , respectively. Let $T: V \rightarrow W$ be a linear transformation of rank r , and suppose that $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$ are such that

$$T(v_i) = \begin{cases} \sigma_i u_i & \text{if } 1 \leq i \leq r \\ 0 & \text{if } r < i. \end{cases}$$

- (a) Prove that $\{u_1, u_2, \dots, u_m\}$ is a set of eigenvectors of TT^* with corresponding eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_m$, where

$$\lambda_i = \begin{cases} \sigma_i^2 & \text{if } 1 \leq i \leq r \\ 0 & \text{if } r < i. \end{cases}$$

- (b) Let A be an $m \times n$ matrix with real or complex entries. Prove that the nonzero singular values of A are the positive square roots of the nonzero eigenvalues of AA^* , including repetitions.
- (c) Prove that TT^* and T^*T have the same nonzero eigenvalues, including repetitions.
- (d) State and prove a result for matrices analogous to (c).
10. Use Exercise 8 of Section 2.5 to obtain another proof of Theorem 6.27, the singular value decomposition theorem for matrices.
11. This exercise relates the singular values of a well-behaved linear operator or matrix to its eigenvalues.
- (a) Let T be a normal linear operator on an n -dimensional inner product space with eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_n$. Prove that the singular values of T are $|\lambda_1|, |\lambda_2|, \dots, |\lambda_n|$.
- (b) State and prove a result for matrices analogous to (a).
12. Let A be a normal matrix with an orthonormal basis of eigenvectors $\beta = \{v_1, v_2, \dots, v_n\}$ and corresponding eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_n$. Let V be the $n \times n$ matrix whose columns are the vectors in β . Prove that for each i there is a scalar θ_i of absolute value 1 such that if U is the $n \times n$ matrix with $\theta_i v_i$ as column i and Σ is the diagonal matrix such that $\Sigma_{ii} = |\lambda_i|$ for each i , then $U\Sigma V^*$ is a singular value decomposition of A .
13. Prove that if A is a positive semidefinite matrix, then the singular values of A are the same as the eigenvalues of A .
14. Prove that if A is a positive definite matrix and $A = U\Sigma V^*$ is a singular value decomposition of A , then $U = V$.
15. Let A be a square matrix with a polar decomposition $A = WP$.
- (a) Prove that A is normal if and only if $WP^2 = P^2W$.
- (b) Use (a) to prove that A is normal if and only if $WP = PW$.
16. Let A be a square matrix. Prove an alternate form of the polar decomposition for A : There exists a unitary matrix W and a positive semidefinite matrix P such that $A = PW$.

17. Let T and U be linear operators on \mathbb{R}^2 defined for all $(x_1, x_2) \in \mathbb{R}^2$ by

$$T(x_1, x_2) = (x_1, 0) \text{ and } U(x_1, x_2) = (x_1 + x_2, 0).$$

- (a) Prove that $(UT)^\dagger \neq T^\dagger U^\dagger$.
- (b) Exhibit matrices A and B such that AB is defined, but $(AB)^\dagger \neq B^\dagger A^\dagger$.

18. Let A be an $m \times n$ matrix. Prove the following results.

- (a) For any $m \times m$ unitary matrix G , $(GA)^\dagger = A^\dagger G^*$.
- (b) For any $n \times n$ unitary matrix H , $(AH)^\dagger = H^* A^\dagger$.

19. Let A be a matrix with real or complex entries. Prove the following results.

- (a) The nonzero singular values of A are the same as the nonzero singular values of A^* , which are the same as the nonzero singular values of A^t .
- (b) $(A^\dagger)^* = (A^*)^\dagger$.
- (c) $(A^\dagger)^t = (A^t)^\dagger$.

20. Let A be a square matrix such that $A^2 = O$. Prove that $(A^\dagger)^2 = O$.

21. Let V and W be finite-dimensional inner product spaces, and let $T: V \rightarrow W$ be linear. Prove the following results.

- (a) $TT^\dagger T = T$.
- (b) $T^\dagger TT^\dagger = T^\dagger$.
- (c) Both $T^\dagger T$ and TT^\dagger are self-adjoint.

Visit [goo.gl/Dz3WQE](#) for a solution. The preceding three statements are called the **Penrose conditions**, and they characterize the pseudoinverse of a linear transformation as shown in Exercise 22.

22. Let V and W be finite-dimensional inner product spaces. Let $T: V \rightarrow W$ and $U: W \rightarrow V$ be linear transformations such that $TUT = T$, $UTU = U$, and both UT and TU are self-adjoint. Prove that $U = T^\dagger$.

23. State and prove a result for matrices that is analogous to the result of Exercise 21.

24. State and prove a result for matrices that is analogous to the result of Exercise 22.

25. Let V and W be finite-dimensional inner product spaces, and let $T: V \rightarrow W$ be linear. Prove the following results.

- (a) If T is one-to-one, then T^*T is invertible and $T^\dagger = (T^*T)^{-1}T^*$.
- (b) If T is onto, then TT^* is invertible and $T^\dagger = T^*(TT^*)^{-1}$.

26. Let V and W be finite-dimensional inner product spaces with orthonormal bases β and γ , respectively, and let $T: V \rightarrow W$ be linear. Prove that $([T]_{\beta}^{\gamma})^{\dagger} = [T^{\dagger}]_{\gamma}^{\beta}$.
27. Let V and W be finite-dimensional inner product spaces, and let $T: V \rightarrow W$ be a linear transformation. Prove part (b) of the lemma to Theorem 6.30: TT^{\dagger} is the orthogonal projection of W on $R(T)$.

6.8* BILINEAR AND QUADRATIC FORMS

There is a certain class of scalar-valued functions of two variables defined on a vector space that arises in the study of such diverse subjects as geometry and multivariable calculus. This is the class of *bilinear forms*. We study the basic properties of this class with a special emphasis on symmetric bilinear forms, and we consider some of its applications to quadratic surfaces and multivariable calculus. In this section, F denotes any field that does not have characteristic two, as defined on page 549.

Bilinear Forms

Definition. Let V be a vector space over a field F . A function H from the set $V \times V$ of ordered pairs of vectors to F is called a **bilinear form** on V if H is linear in each variable when the other variable is held fixed; that is, H is a bilinear form on V if

- $H(ax_1 + x_2, y) = aH(x_1, y) + H(x_2, y)$ for all $x_1, x_2, y \in V$ and $a \in F$
- $H(x, ay_1 + y_2) = aH(x, y_1) + H(x, y_2)$ for all $x, y_1, y_2 \in V$ and $a \in F$.

We denote the set of all bilinear forms on V by $\mathcal{B}(V)$. Observe that an inner product on a vector space is a bilinear form if the underlying field is real, but not if the underlying field is complex.

Example 1

Define a function $H: \mathbb{R}^2 \times \mathbb{R}^2 \rightarrow \mathbb{R}$ by

$$H\left(\begin{pmatrix} a_1 \\ a_2 \end{pmatrix}, \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}\right) = 2a_1b_1 + 3a_1b_2 + 4a_2b_1 - a_2b_2 \quad \text{for } \begin{pmatrix} a_1 \\ a_2 \end{pmatrix}, \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} \in \mathbb{R}^2.$$

We could verify directly that H is a bilinear form on \mathbb{R}^2 . However, it is more enlightening and less tedious to observe that if

$$A = \begin{pmatrix} 2 & 3 \\ 4 & -1 \end{pmatrix}, \quad x = \begin{pmatrix} a_1 \\ a_2 \end{pmatrix}, \quad \text{and} \quad y = \begin{pmatrix} b_1 \\ b_2 \end{pmatrix},$$

then

$$H(x, y) = x^t A y.$$

The bilinearity of H now follows directly from the distributive property of matrix multiplication over matrix addition. ♦

The preceding bilinear form is a special case of the next example.

Example 2

Let $V = F^n$, where the vectors are considered as column vectors. For any $A \in M_{n \times n}(F)$, define $H: V \times V \rightarrow F$ by

$$H(x, y) = x^t A y \quad \text{for } x, y \in V.$$

Notice that since x and y are $n \times 1$ matrices and A is an $n \times n$ matrix, $H(x, y)$ is a 1×1 matrix. We identify this matrix with its single entry. The bilinearity of H follows as in Example 1. For example, for $a \in F$ and $x_1, x_2, y \in V$, we have

$$\begin{aligned} H(ax_1 + x_2, y) &= (ax_1 + x_2)^t A y = (ax_1^t + x_2^t) A y \\ &= ax_1^t A y + x_2^t A y \\ &= aH(x_1, y) + H(x_2, y). \end{aligned} \quad \blacklozenge$$

We list several properties possessed by all bilinear forms. Their proofs are left to the reader (see Exercise 2).

For any bilinear form H on a vector space V over a field F , the following properties hold.

1. If, for any $x \in V$, the functions $L_x, R_x: V \rightarrow F$ are defined by

$$L_x(y) = H(x, y) \quad \text{and} \quad R_x(y) = H(y, x) \quad \text{for all } y \in V,$$

then L_x and R_x are linear.

2. $H(\theta, x) = H(x, \theta) = 0$ for all $x \in V$.
3. For all $x, y, z, w \in V$,

$$H(x + y, z + w) = H(x, z) + H(x, w) + H(y, z) + H(y, w).$$

4. If $J: V \times V \rightarrow F$ is defined by $J(x, y) = H(y, x)$, then J is a bilinear form.

Definitions. Let V be a vector space, let H_1 and H_2 be bilinear forms on V , and let a be a scalar. We define the **sum** $H_1 + H_2$ and the **scalar product** aH_1 by the equations

$$(H_1 + H_2)(x, y) = H_1(x, y) + H_2(x, y)$$

and

$$(aH_1)(x, y) = a(H_1(x, y)) \quad \text{for all } x, y \in V.$$

The following theorem is an immediate consequence of the definitions.

Theorem 6.31. For any vector space V , the sum of two bilinear forms and the product of a scalar and a bilinear form on V are again bilinear forms on V . Furthermore, $\mathcal{B}(V)$ is a vector space with respect to these operations.

Proof. Exercise. ■

Let $\beta = \{v_1, v_2, \dots, v_n\}$ be an ordered basis for an n -dimensional vector space V , and let $H \in \mathcal{B}(V)$. We can associate with H an $n \times n$ matrix A whose entry in row i and column j is defined by

$$A_{ij} = H(v_i, v_j) \quad \text{for } i, j = 1, 2, \dots, n.$$

Definition. The matrix A above is called the **matrix representation** of H with respect to the ordered basis β and is denoted by $\psi_\beta(H)$.

We can therefore regard ψ_β as a mapping from $\mathcal{B}(V)$ to $M_{n \times n}(F)$, where F is the field of scalars for V , that takes a bilinear form H into its matrix representation $\psi_\beta(H)$. We first consider an example and then show that ψ_β is an isomorphism.

Example 3

Consider the bilinear form H of Example 1, and let

$$\beta = \left\{ \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ -1 \end{pmatrix} \right\} \quad \text{and} \quad B = \psi_\beta(H).$$

Then

$$\begin{aligned} B_{11} &= H\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \end{pmatrix}\right) = 2 + 3 + 4 - 1 = 8, \\ B_{12} &= H\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ -1 \end{pmatrix}\right) = 2 - 3 + 4 + 1 = 4, \\ B_{21} &= H\left(\begin{pmatrix} 1 \\ -1 \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \end{pmatrix}\right) = 2 + 3 - 4 + 1 = 2, \end{aligned}$$

and

$$B_{22} = H\left(\begin{pmatrix} 1 \\ -1 \end{pmatrix}, \begin{pmatrix} 1 \\ -1 \end{pmatrix}\right) = 2 - 3 - 4 - 1 = -6.$$

So

$$\psi_\beta(H) = \begin{pmatrix} 8 & 4 \\ 2 & -6 \end{pmatrix}.$$

If γ is the standard ordered basis for \mathbb{R}^2 , the reader can verify that

$$\psi_\gamma(H) = \begin{pmatrix} 2 & 3 \\ 4 & -1 \end{pmatrix}. \quad \blacklozenge$$

Theorem 6.32. For any n -dimensional vector space V over F and any ordered basis β for V , $\psi_\beta: \mathcal{B}(V) \rightarrow M_{n \times n}(F)$ is an isomorphism.

Proof. We leave the proof that ψ_β is linear to the reader.

To show that ψ_β is one-to-one, suppose that $\psi_\beta(H) = O$ for some $H \in \mathcal{B}(V)$. Fix $v_i \in \beta$, and recall the mapping $L_{v_i}: V \rightarrow F$, which is linear by property 1 on page 420. By hypothesis, $L_{v_i}(v_j) = H(v_i, v_j) = 0$ for all $v_j \in \beta$. Hence L_{v_i} is the zero transformation from V to F . So

$$H(v_i, x) = L_{v_i}(x) = 0 \quad \text{for all } x \in V \text{ and } v_i \in \beta. \quad (7)$$

Next fix an arbitrary $y \in V$, and recall the linear mapping $R_y: V \rightarrow F$ defined in property 1 on page 420. By (7), $R_y(v_i) = H(v_i, y) = 0$ for all $v_i \in \beta$, and hence R_y is the zero transformation. So $H(x, y) = R_y(x) = 0$ for all $x, y \in V$. Thus H is the zero bilinear form, and therefore ψ_β is one-to-one.

To show that ψ_β is onto, consider any $A \in M_{n \times n}(F)$. Recall the isomorphism $\phi_\beta: V \rightarrow F^n$ defined in Section 2.4. For $x \in V$, we view $\phi_\beta(x) \in F^n$ as a column vector. Let $H: V \times V \rightarrow F$ be the mapping defined by

$$H(x, y) = [\phi_\beta(x)]^t A [\phi_\beta(y)] \quad \text{for all } x, y \in V.$$

A slight embellishment of the method of Example 2 can be used to prove that $H \in \mathcal{B}(V)$. We show that $\psi_\beta(H) = A$. Let $v_i, v_j \in \beta$. Then $\phi_\beta(v_i) = e_i$ and $\phi_\beta(v_j) = e_j$; hence, for any i and j ,

$$H(v_i, v_j) = [\phi_\beta(v_i)]^t A [\phi_\beta(v_j)] = e_i^t A e_j = A_{ij}.$$

We conclude that $\psi_\beta(H) = A$ and ψ_β is onto. ■

Corollary 1. For any n -dimensional vector space V , $\mathcal{B}(V)$ has dimension n^2 .

Proof. Exercise. ■

The following corollary is easily established by reviewing the proof of Theorem 6.32.

Corollary 2. Let V be an n -dimensional vector space over F with ordered basis β . If $H \in \mathcal{B}(V)$ and $A \in M_{n \times n}(F)$, then $\psi_\beta(H) = A$ if and only if $H(x, y) = [\phi_\beta(x)]^t A [\phi_\beta(y)]$ for all $x, y \in V$.

The following result is now an immediate consequence of Corollary 2.

Corollary 3. Let F be a field, n be a positive integer, and β be the standard ordered basis for F^n . Then for any $H \in \mathcal{B}(F^n)$, there exists a unique matrix $A \in M_{n \times n}(F)$, namely, $A = \psi_\beta(H)$, such that

$$H(x, y) = x^t A y \quad \text{for all } x, y \in F^n.$$

Example 4

Define a function $H: \mathbb{R}^2 \times \mathbb{R}^2 \rightarrow \mathbb{R}$ by

$$H\left(\begin{pmatrix} a_1 \\ a_2 \end{pmatrix}, \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}\right) = \det \begin{pmatrix} a_1 & b_1 \\ a_2 & b_2 \end{pmatrix} = a_1 b_2 - a_2 b_1 \quad \text{for } \begin{pmatrix} a_1 \\ a_2 \end{pmatrix}, \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} \in \mathbb{R}^2.$$

It can be shown that H is a bilinear form. We find the matrix A in Corollary 3 such that $H(x, y) = x^t A y$ for all $x, y \in \mathbb{R}^2$.

Since $A_{ij} = H(e_i, e_j)$ for all i and j , we have

$$\begin{aligned} A_{11} &= \det \begin{pmatrix} 1 & 1 \\ 0 & 0 \end{pmatrix} = 0 & A_{12} &= \det \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} = 1, \\ A_{21} &= \det \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} = -1 & \text{and} & A_{22} = \det \begin{pmatrix} 0 & 0 \\ 1 & 1 \end{pmatrix} = 0. \end{aligned}$$

Therefore $A = \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}$. ◆

There is an analogy between bilinear forms and linear operators on finite-dimensional vector spaces in that both are associated with unique square matrices and the correspondences depend on the choice of an ordered basis for the vector space. As in the case of linear operators, one can pose the following question: How does the matrix corresponding to a fixed bilinear form change when the ordered basis is changed? As we have seen, the corresponding question for matrix representations of linear operators leads to the definition of the similarity relation on square matrices. In the case of bilinear forms, the corresponding question leads to another relation on square matrices, the *congruence* relation.

Definition. Let $A, B \in M_{n \times n}(F)$. Then B is said to be **congruent** to A if there exists an invertible matrix $Q \in M_{n \times n}(F)$ such that $B = Q^t A Q$.

Observe that the relation of congruence is an equivalence relation (see Exercise 12).

The next theorem relates congruence to the matrix representation of a bilinear form.

Theorem 6.33. Let V be a finite-dimensional vector space with ordered bases $\beta = \{v_1, v_2, \dots, v_n\}$ and $\gamma = \{w_1, w_2, \dots, w_n\}$, and let Q be the change of coordinate matrix changing γ -coordinates into β -coordinates. Then, for any $H \in \mathcal{B}(V)$, we have $\psi_\gamma(H) = Q^t \psi_\beta(H) Q$. Therefore $\psi_\gamma(H)$ is congruent to $\psi_\beta(H)$.

Proof. There are essentially two proofs of this theorem. One involves a direct computation, while the other follows immediately from a clever observation. We give the more direct proof here, leaving the other proof for the exercises (see Exercise 13).

Suppose that $A = \psi_\beta(H)$ and $B = \psi_\gamma(H)$. Then for $1 \leq i, j \leq n$,

$$w_i = \sum_{k=1}^n Q_{ki} v_k \quad \text{and} \quad w_j = \sum_{r=1}^n Q_{rj} v_r.$$

Thus

$$\begin{aligned} B_{ij} &= H(w_i, w_j) = H\left(\sum_{k=1}^n Q_{ki} v_k, w_j\right) \\ &= \sum_{k=1}^n Q_{ki} H(v_k, w_j) \\ &= \sum_{k=1}^n Q_{ki} H\left(v_k, \sum_{r=1}^n Q_{rj} v_r\right) \\ &= \sum_{k=1}^n Q_{ki} \sum_{r=1}^n Q_{rj} H(v_k, v_r) \\ &= \sum_{k=1}^n Q_{ki} \sum_{r=1}^n Q_{rj} A_{kr} \\ &= \sum_{k=1}^n Q_{ki} \sum_{r=1}^n A_{kr} Q_{rj} \\ &= \sum_{k=1}^n Q_{ki} (AQ)_{kj} \\ &= \sum_{k=1}^n Q_{ik}^t (AQ)_{kj} = (Q^t AQ)_{ij}. \end{aligned}$$

Hence $B = Q^t AQ$. ■

The following result is the converse of Theorem 6.33.

Corollary. Let V be an n -dimensional vector space with ordered basis β , and let H be a bilinear form on V . For any $n \times n$ matrix B , if B is congruent to $\psi_\beta(H)$, then there exists an ordered basis γ for V such that $\psi_\gamma(H) = B$. In this case, if Q is a matrix such that $B = Q^t \psi_\beta(H) Q$, then Q changes γ -coordinates into β -coordinates.

Proof. Suppose that $B = Q^t \psi_\beta(H) Q$ for some invertible matrix Q and that $\beta = \{v_1, v_2, \dots, v_n\}$. Let $\gamma = \{w_1, w_2, \dots, w_n\}$, where

$$w_j = \sum_{i=1}^n Q_{ij} v_i \quad \text{for } 1 \leq j \leq n.$$

Since Q is invertible, γ is an ordered basis for V , and Q is the change of coordinate matrix that changes γ -coordinates into β -coordinates. Therefore, by Theorem 6.33,

$$B = Q^t \psi_\beta(H) Q = \psi_\gamma(H). \quad \blacksquare$$

Symmetric Bilinear Forms

Like the diagonalization problem for linear operators, there is an analogous *diagonalization* problem for bilinear forms, namely, the problem of determining those bilinear forms for which there are diagonal matrix representations. As we will see, there is a close relationship between *diagonalizable* bilinear forms and those that are called *symmetric*.

Definition. A bilinear form H on a vector space V is **symmetric** if $H(x, y) = H(y, x)$ for all $x, y \in V$.

As the name suggests, symmetric bilinear forms correspond to symmetric matrices.

Theorem 6.34. Let H be a bilinear form on a finite-dimensional vector space V , and let β be an ordered basis for V . Then H is symmetric if and only if $\psi_\beta(H)$ is symmetric.

Proof. Let $\beta = \{v_1, v_2, \dots, v_n\}$ and $B = \psi_\beta(H)$.

First assume that H is symmetric. Then for $1 \leq i, j \leq n$,

$$B_{ij} = H(v_i, v_j) = H(v_j, v_i) = B_{ji},$$

and it follows that B is symmetric.

Conversely, suppose that B is symmetric. Let $J: V \times V \rightarrow F$, where F is the field of scalars for V , be the mapping defined by $J(x, y) = H(y, x)$ for all $x, y \in V$. By property 4 on page 420, J is a bilinear form. Let $C = \psi_\beta(J)$. Then, for $1 \leq i, j \leq n$,

$$C_{ij} = J(v_i, v_j) = H(v_j, v_i) = B_{ji} = B_{ij}.$$

Thus $C = B$. Since ψ_β is one-to-one, we have $J = H$. Hence $H(y, x) = J(x, y) = H(x, y)$ for all $x, y \in V$, and therefore H is symmetric. \blacksquare

Definition. A bilinear form H on a finite-dimensional vector space V is called **diagonalizable** if there is an ordered basis β for V such that $\psi_\beta(H)$ is a diagonal matrix.

Corollary. Let H be a diagonalizable bilinear form on a finite-dimensional vector space V . Then H is symmetric.

Proof. Suppose that H is diagonalizable. Then there is an ordered basis β for V such that $\psi_\beta(H) = D$ is a diagonal matrix. Trivially, D is a symmetric matrix, and hence, by Theorem 6.34, H is symmetric. ■

Unfortunately, the converse is not true, as is illustrated by the following example.

Example 5

Let $F = Z_2$, $V = F^2$, and $H: V \times V \rightarrow F$ be the bilinear form defined by

$$H\left(\begin{pmatrix} a_1 \\ a_2 \end{pmatrix}, \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}\right) = a_1b_2 + a_2b_1.$$

Clearly H is symmetric. In fact, if β is the standard ordered basis for V , then

$$A = \psi_\beta(H) = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix},$$

a symmetric matrix. We show that H is not diagonalizable.

By way of contradiction, suppose that H is diagonalizable. Then there is an ordered basis γ for V such that $B = \psi_\gamma(H)$ is a diagonal matrix. So by Theorem 6.33, there exists an invertible matrix Q such that $B = Q^t A Q$. Since Q is invertible, it follows that $\text{rank}(B) = \text{rank}(A) = 2$, and consequently the diagonal entries of B are nonzero. Since the only nonzero scalar of F is 1,

$$B = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}.$$

Suppose that

$$Q = \begin{pmatrix} a & b \\ c & d \end{pmatrix}.$$

Then

$$\begin{aligned} \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} &= B = Q^t A Q \\ &= \begin{pmatrix} a & c \\ b & d \end{pmatrix} \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} a & b \\ c & d \end{pmatrix} = \begin{pmatrix} ac + ac & bc + ad \\ bc + ad & bd + bd \end{pmatrix}. \end{aligned}$$

But $p + p = 0$ for all $p \in F$; hence $ac + ac = 0$. Thus, comparing the row 1, column 1 entries of the matrices in the equation above, we conclude that $1 = 0$, a contradiction. Therefore H is not diagonalizable. ♦

The bilinear form of Example 5 is an anomaly. Its failure to be diagonalizable is due to the fact that the scalar field Z_2 is of characteristic two. Recall from Appendix C that a field F is of **characteristic two** if $1 + 1 = 0$ in F . If F is not of characteristic two, then $1 + 1 = 2$ has a multiplicative inverse, which we denote by $1/2$.

Before proving the converse of the corollary to Theorem 6.34 for scalar fields that are not of characteristic two, we establish the following lemma.

Lemma. Let H be a nonzero symmetric bilinear form on a vector space V over a field F not of characteristic two. Then there is a vector x in V such that $H(x, x) \neq 0$.

Proof. Since H is nonzero, we can choose vectors $u, v \in V$ such that $H(u, v) \neq 0$. If $H(u, u) \neq 0$ or $H(v, v) \neq 0$, there is nothing to prove. Otherwise, set $x = u + v$. Then

$$H(x, x) = H(u, u) + H(u, v) + H(v, u) + H(v, v) = 2H(u, v) \neq 0$$

because $2 \neq 0$ and $H(u, v) \neq 0$. ■

Theorem 6.35. Let V be a finite-dimensional vector space over a field F not of characteristic two. Then every symmetric bilinear form on V is diagonalizable.

Proof. We use mathematical induction on $n = \dim(V)$. If $n = 1$, then every element of $\mathcal{B}(V)$ is diagonalizable. Now suppose that the theorem is valid for vector spaces of dimension less than n for some fixed integer $n > 1$, and suppose that $\dim(V) = n$. If H is the zero bilinear form on V , then trivially H is diagonalizable; so suppose that H is a nonzero symmetric bilinear form on V . By the lemma, there exists a nonzero vector x in V such that $H(x, x) \neq 0$. Recall the function $L_x : V \rightarrow F$ defined by $L_x(y) = H(x, y)$ for all $y \in V$. By property 1 on page 420, L_x is linear. Furthermore, since $L_x(x) = H(x, x) \neq 0$, L_x is nonzero. Consequently, $\text{rank}(L_x) = 1$, and hence $\dim(N(L_x)) = n - 1$.

The restriction of H to $N(L_x)$ is obviously a symmetric bilinear form on a vector space of dimension $n - 1$. Thus, by the induction hypothesis, there exists an ordered basis $\{v_1, v_2, \dots, v_{n-1}\}$ for $N(L_x)$ such that $H(v_i, v_j) = 0$ for $i \neq j$ ($1 \leq i, j \leq n - 1$). Set $v_n = x$. Then $v_n \notin N(L_x)$, and so $\beta = \{v_1, v_2, \dots, v_n\}$ is an ordered basis for V . In addition, $H(v_i, v_n) = H(v_n, v_i) = 0$ for $i = 1, 2, \dots, n - 1$. We conclude that $\psi_\beta(H)$ is a diagonal matrix, and therefore H is diagonalizable. ■

Corollary. Let F be a field that is not of characteristic two. If $A \in M_{n \times n}(F)$ is a symmetric matrix, then A is congruent to a diagonal matrix.

Proof. Exercise. ■

Diagonalization of Symmetric Matrices

Let A be a symmetric $n \times n$ matrix with entries from a field F not of characteristic two. By the corollary to Theorem 6.35, there are matrices $Q, D \in M_{n \times n}(F)$ such that Q is invertible, D is diagonal, and $Q^t A Q = D$. We now give a method for computing Q and D . This method requires familiarity with elementary matrices and their properties, which the reader may wish to review in Section 3.1.

If E is an elementary $n \times n$ matrix, then AE can be obtained by performing an elementary column operation on A . By Exercise 21, $E^t A$ can be obtained by performing the same operation on the rows of A rather than on its columns. Thus $E^t AE$ can be obtained from A by performing an elementary operation on the columns of A and then performing the same operation on the rows of AE . (Note that the order of the operations can be reversed because of the associative property of matrix multiplication.) Suppose that Q is an invertible matrix and D is a diagonal matrix such that $Q^t AQ = D$. By Corollary 3 to Theorem 3.6 (p. 158), Q is a product of elementary matrices, say $Q = E_1 E_2 \cdots E_k$. Thus

$$D = Q^t AQ = E_k^t E_{k-1}^t \cdots E_1^t AE_1 E_2 \cdots E_k.$$

From the preceding equation, we conclude that *by means of several elementary column operations and the corresponding row operations, A can be transformed into a diagonal matrix D . Furthermore, if E_1, E_2, \dots, E_k are the elementary matrices corresponding to these elementary column operations indexed in the order performed, and if $Q = E_1 E_2 \cdots E_k$, then $Q^t AQ = D$.*

Example 6

Let A be the symmetric matrix in $M_{3 \times 3}(R)$ defined by

$$A = \begin{pmatrix} 1 & -1 & 3 \\ -1 & 2 & 1 \\ 3 & 1 & 1 \end{pmatrix}.$$

We use the procedure just described to find an invertible matrix Q and a diagonal matrix D such that $Q^t AQ = D$.

We begin by eliminating all of the nonzero entries in the first row and first column except for the entry in column 1 and row 1. To this end, we add the first column of A to the second column to produce a zero in row 1 and column 2. The elementary matrix that corresponds to this elementary column operation is

$$E_1 = \begin{pmatrix} 1 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}.$$

We perform the corresponding elementary operation on the rows of AE_1 to obtain

$$E_1^t AE_1 = \begin{pmatrix} 1 & 0 & 3 \\ 0 & 1 & 4 \\ 3 & 4 & 1 \end{pmatrix}.$$

We now use the first column of $E_1^t AE_1$ to eliminate the 3 in row 1 column 3, and follow this operation with the corresponding row operation. The corresponding elementary matrix E_2 and the result of the elementary operations

$E_2^t E_1^t A E_1 E_2$ are, respectively,

$$E_2 = \begin{pmatrix} 1 & 0 & -3 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \quad \text{and} \quad E_2^t E_1^t A E_1 E_2 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 4 \\ 0 & 4 & -8 \end{pmatrix}.$$

Finally, we subtract 4 times the second column of $E_2^t E_1^t A E_1 E_2$ from the third column and follow this with the corresponding row operation. The corresponding elementary matrix E_3 and the result of the elementary operations $E_3^t E_2^t E_1^t A E_1 E_2 E_3$ are, respectively,

$$E_3 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & -4 \\ 0 & 0 & 1 \end{pmatrix} \quad \text{and} \quad E_3^t E_2^t E_1^t A E_1 E_2 E_3 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & -24 \end{pmatrix}.$$

Since we have obtained a diagonal matrix, the process is complete. So we let

$$Q = E_1 E_2 E_3 = \begin{pmatrix} 1 & 1 & -7 \\ 0 & 1 & -4 \\ 0 & 0 & 1 \end{pmatrix} \quad \text{and} \quad D = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & -24 \end{pmatrix}$$

to obtain the desired diagonalization $Q^t A Q = D$. ◆

The reader should justify the following method for computing Q without recording each elementary matrix separately. The method is inspired by the algorithm for computing the inverse of a matrix developed in Section 3.2. We use a sequence of elementary column operations and corresponding row operations to change the $n \times 2n$ matrix $(A|I)$ into the form $(D|B)$, where D is a diagonal matrix and $B = Q^t$. It then follows that $D = Q^t A Q$.

Starting with the matrix A of the preceding example, this method produces the following sequence of matrices:

$$\begin{aligned} (A|I) &= \left(\begin{array}{ccc|ccc} 1 & -1 & 3 & 1 & 0 & 0 \\ -1 & 2 & 1 & 0 & 1 & 0 \\ 3 & 1 & 1 & 0 & 0 & 1 \end{array} \right) \rightarrow \left(\begin{array}{ccc|ccc} 1 & 0 & 3 & 1 & 0 & 0 \\ -1 & 1 & 1 & 0 & 1 & 0 \\ 3 & 4 & 1 & 0 & 0 & 1 \end{array} \right) \\ &\rightarrow \left(\begin{array}{ccc|ccc} 1 & 0 & 3 & 1 & 0 & 0 \\ 0 & 1 & 4 & 1 & 1 & 0 \\ 3 & 4 & 1 & 0 & 0 & 1 \end{array} \right) \rightarrow \left(\begin{array}{ccc|ccc} 1 & 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 4 & 1 & 1 & 0 \\ 3 & 4 & -8 & 0 & 0 & 1 \end{array} \right) \\ &\rightarrow \left(\begin{array}{ccc|ccc} 1 & 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 4 & 1 & 1 & 0 \\ 0 & 4 & -8 & -3 & 0 & 1 \end{array} \right) \rightarrow \left(\begin{array}{ccc|ccc} 1 & 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & 1 & 0 \\ 0 & 4 & -24 & -3 & 0 & 1 \end{array} \right) \\ &\rightarrow \left(\begin{array}{ccc|ccc} 1 & 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & 1 & 0 \\ 0 & 0 & -24 & -7 & -4 & 1 \end{array} \right) = (D|Q^t). \end{aligned}$$

Therefore

$$D = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & -24 \end{pmatrix}, \quad Q^t = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ -7 & -4 & 1 \end{pmatrix}, \quad \text{and} \quad Q = \begin{pmatrix} 1 & 1 & -7 \\ 0 & 1 & -4 \\ 0 & 0 & 1 \end{pmatrix}.$$

Quadratic Forms

Associated with symmetric bilinear forms are functions called *quadratic forms*.

Definition. Let V be a vector space over a field F not of characteristic two.² A function $K: V \rightarrow F$ is called a **quadratic form on V** if there exists a symmetric bilinear form $H \in \mathcal{B}(V)$ such that

$$K(x) = H(x, x) \quad \text{for all } x \in V. \quad (8)$$

There is a one-to-one correspondence between symmetric bilinear forms and quadratic forms given by (8). In fact, if K is a quadratic form on a vector space V over a field F not of characteristic two, and $K(x) = H(x, x)$ for some symmetric bilinear form H on V , then we can recover H from K because

$$H(x, y) = \frac{1}{2}[K(x + y) - K(x) - K(y)] \quad (9)$$

(See Exercise 16.)

Example 7

The classic example of a quadratic form is the homogeneous second-degree polynomial of several variables. Given the variables t_1, t_2, \dots, t_n that take values in a field F not of characteristic two and given (not necessarily distinct) scalars a_{ij} ($1 \leq i \leq j \leq n$), define the polynomial

$$f(t_1, t_2, \dots, t_n) = \sum_{i \leq j} a_{ij} t_i t_j.$$

Any such polynomial is a quadratic form. In fact, if β is the standard ordered basis for F^n , then the symmetric bilinear form H corresponding to the quadratic form f has the matrix representation $\psi_\beta(H) = A$, where

$$A_{ij} = A_{ji} = \begin{cases} a_{ii} & \text{if } i = j \\ \frac{1}{2}a_{ij} & \text{if } i \neq j. \end{cases}$$

²Although it is possible to define quadratic forms over a field of characteristic 2, our definition can't be used in this case.

To see this, apply (9) to obtain $H(e_i, e_j) = A_{ij}$ from the quadratic form K , and verify that f is computable from H by (8) using f in place of K .

For example, given the polynomial

$$f(t_1, t_2, t_3) = 2t_1^2 - t_2^2 + 6t_1t_2 - 4t_2t_3$$

with real coefficients, let

$$A = \begin{pmatrix} 2 & 3 & 0 \\ 3 & -1 & -2 \\ 0 & -2 & 0 \end{pmatrix}.$$

Setting $H(x, y) = x^t A y$ for all $x, y \in \mathbb{R}^3$, we see that

$$f(t_1, t_2, t_3) = (t_1, t_2, t_3) A \begin{pmatrix} t_1 \\ t_2 \\ t_3 \end{pmatrix} \quad \text{for } \begin{pmatrix} t_1 \\ t_2 \\ t_3 \end{pmatrix} \in \mathbb{R}^3. \quad \blacklozenge$$

Quadratic Forms on a Real Inner Product Space

Since symmetric matrices over R are *orthogonally diagonalizable* (see Theorem 6.20 p. 381), the theory of symmetric bilinear forms and quadratic forms on finite-dimensional vector spaces over R is especially nice. The following theorem and its corollary are useful.

Theorem 6.36. *Let V be a finite-dimensional real inner product space, and let H be a symmetric bilinear form on V . Then there exists an orthonormal basis β for V such that $\psi_\beta(H)$ is a diagonal matrix.*

Proof. Choose any orthonormal basis $\gamma = \{v_1, v_2, \dots, v_n\}$ for V , and let $A = \psi_\gamma(H)$. Since A is symmetric, there exists an orthogonal matrix Q and a diagonal matrix D such that $D = Q^t A Q$ by Theorem 6.20. Let $\beta = \{w_1, w_2, \dots, w_n\}$ be defined by

$$w_j = \sum_{i=1}^n Q_{ij} v_i \quad \text{for } 1 \leq j \leq n.$$

By Theorem 6.33, $\psi_\beta(H) = D$. Furthermore, since Q is orthogonal and γ is orthonormal, β is orthonormal by Exercise 30 of Section 6.5. ■

Corollary. *Let K be a quadratic form on a finite-dimensional real inner product space V . There exists an orthonormal basis $\beta = \{v_1, v_2, \dots, v_n\}$ for V and scalars $\lambda_1, \lambda_2, \dots, \lambda_n$ (not necessarily distinct) such that if $x \in V$ and*

$$x = \sum_{i=1}^n s_i v_i, \quad s_i \in R,$$

then

$$K(x) = \sum_{i=1}^n \lambda_i s_i^2.$$

In fact, if H is the symmetric bilinear form determined by K , then β can be chosen to be any orthonormal basis for V such that $\psi_\beta(H)$ is a diagonal matrix.

Proof. Let H be the symmetric bilinear form for which $K(x) = H(x, x)$ for all $x \in V$. By Theorem 6.36, there exists an orthonormal basis $\beta = \{v_1, v_2, \dots, v_n\}$ for V such that $\psi_\beta(H)$ is the diagonal matrix

$$D = \begin{pmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & \lambda_n \end{pmatrix}.$$

Let $x \in V$, and suppose that $x = \sum_{i=1}^n s_i v_i$. Then

$$K(x) = H(x, x) = [\phi_\beta(x)]^t D [\phi_\beta(x)] = (s_1, s_2, \dots, s_n) D \begin{pmatrix} s_1 \\ s_2 \\ \vdots \\ s_n \end{pmatrix} = \sum_{i=1}^n \lambda_i s_i^2. \quad \blacksquare$$

Example 8

For the homogeneous real polynomial of degree 2 defined by

$$f(t_1, t_2) = 5t_1^2 + 2t_2^2 + 4t_1 t_2, \quad (10)$$

we find an orthonormal basis $\gamma = \{x_1, x_2\}$ for \mathbb{R}^2 and scalars λ_1 and λ_2 such that if

$$\begin{pmatrix} t_1 \\ t_2 \end{pmatrix} \in \mathbb{R}^2 \quad \text{and} \quad \begin{pmatrix} t_1 \\ t_2 \end{pmatrix} = s_1 x_1 + s_2 x_2,$$

then $f(t_1, t_2) = \lambda_1 s_1^2 + \lambda_2 s_2^2$. We can think of s_1 and s_2 as the coordinates of (t_1, t_2) relative to γ . Thus the polynomial $f(t_1, t_2)$, as an expression involving the coordinates of a point with respect to the standard ordered basis for \mathbb{R}^2 , is transformed into a new polynomial $g(s_1, s_2) = \lambda_1 s_1^2 + \lambda_2 s_2^2$ interpreted as an expression involving the coordinates of a point relative to the new ordered basis γ .

Let H denote the symmetric bilinear form corresponding to the quadratic form defined by (10), let β be the standard ordered basis for \mathbb{R}^2 , and let $A = \psi_\beta(H)$. Then

$$A = \psi_\beta(H) = \begin{pmatrix} 5 & 2 \\ 2 & 2 \end{pmatrix}.$$

Next, we find an orthogonal matrix Q such that $Q^t A Q$ is a diagonal matrix. For this purpose, observe that $\lambda_1 = 6$ and $\lambda_2 = 1$ are the eigenvalues of A with corresponding orthonormal eigenvectors

$$v_1 = \frac{1}{\sqrt{5}} \begin{pmatrix} 2 \\ 1 \end{pmatrix} \quad \text{and} \quad v_2 = \frac{1}{\sqrt{5}} \begin{pmatrix} 1 \\ -2 \end{pmatrix}.$$

Let $\gamma = \{v_1, v_2\}$. Then γ is an orthonormal basis for \mathbb{R}^2 consisting of eigenvectors of A . Hence, setting

$$Q = \frac{1}{\sqrt{5}} \begin{pmatrix} 2 & 1 \\ 1 & -2 \end{pmatrix},$$

we see that Q is an orthogonal matrix and

$$Q^t A Q = \begin{pmatrix} 6 & 0 \\ 0 & 1 \end{pmatrix}.$$

Clearly Q is also a change of coordinate matrix. Consequently,

$$\psi_\gamma(H) = Q^t \psi_\beta(H) Q = Q^t A Q = \begin{pmatrix} 6 & 0 \\ 0 & 1 \end{pmatrix}.$$

Thus by the corollary to Theorem 6.36,

$$K(x) = 6s_1^2 + s_2^2$$

for any $x = s_1 v_1 + s_2 v_2 \in \mathbb{R}^2$. So $g(s_1, s_2) = 6s_1^2 + s_2^2$. \blacklozenge

The next example illustrates how the theory of quadratic forms can be applied to the problem of describing quadratic surfaces in \mathbb{R}^3 .

Example 9

Let \mathcal{S} be the surface in \mathbb{R}^3 defined by the equation

$$2t_1^2 + 6t_1 t_2 + 5t_2^2 - 2t_2 t_3 + 2t_3^2 + 3t_1 - 2t_2 - t_3 + 14 = 0. \quad (11)$$

Then (11) describes the points of \mathcal{S} in terms of their coordinates relative to β , the standard ordered basis for \mathbb{R}^3 . We find a new orthonormal basis γ for \mathbb{R}^3 so that the equation describing the coordinates of \mathcal{S} relative to γ is simpler than (11).

We begin with the observation that the terms of second degree on the left side of (11) add to form a quadratic form K on \mathbb{R}^3 :

$$K \begin{pmatrix} t_1 \\ t_2 \\ t_3 \end{pmatrix} = 2t_1^2 + 6t_1 t_2 + 5t_2^2 - 2t_2 t_3 + 2t_3^2.$$

Next, we diagonalize K . Let H be the symmetric bilinear form corresponding to K , and let $A = \psi_\beta(H)$. Then

$$A = \begin{pmatrix} 2 & 3 & 0 \\ 3 & 5 & -1 \\ 0 & -1 & 2 \end{pmatrix}.$$

The characteristic polynomial of A is $(-1)(t-2)(t-7)t$; hence A has the eigenvalues $\lambda_1 = 2$, $\lambda_2 = 7$, and $\lambda_3 = 0$. Corresponding unit eigenvectors are

$$v_1 = \frac{1}{\sqrt{10}} \begin{pmatrix} 1 \\ 0 \\ 3 \end{pmatrix}, \quad v_2 = \frac{1}{\sqrt{35}} \begin{pmatrix} 3 \\ 5 \\ -1 \end{pmatrix}, \quad \text{and} \quad v_3 = \frac{1}{\sqrt{14}} \begin{pmatrix} -3 \\ 2 \\ 1 \end{pmatrix}.$$

Set $\gamma = \{v_1, v_2, v_3\}$ and

$$Q = \begin{pmatrix} \frac{1}{\sqrt{10}} & \frac{3}{\sqrt{35}} & \frac{-3}{\sqrt{14}} \\ 0 & \frac{5}{\sqrt{35}} & \frac{2}{\sqrt{14}} \\ \frac{3}{\sqrt{10}} & \frac{-1}{\sqrt{35}} & \frac{1}{\sqrt{14}} \end{pmatrix}.$$

As in Example 8, Q is a change of coordinate matrix changing γ -coordinates to β -coordinates, and

$$\psi_\gamma(H) = Q^t \psi_\beta(H) Q = Q^t A Q = \begin{pmatrix} 2 & 0 & 0 \\ 0 & 7 & 0 \\ 0 & 0 & 0 \end{pmatrix}.$$

By the corollary to Theorem 6.36, if $x = s_1 v_1 + s_2 v_2 + s_3 v_3$, then

$$K(x) = 2s_1^2 + 7s_2^2. \tag{12}$$

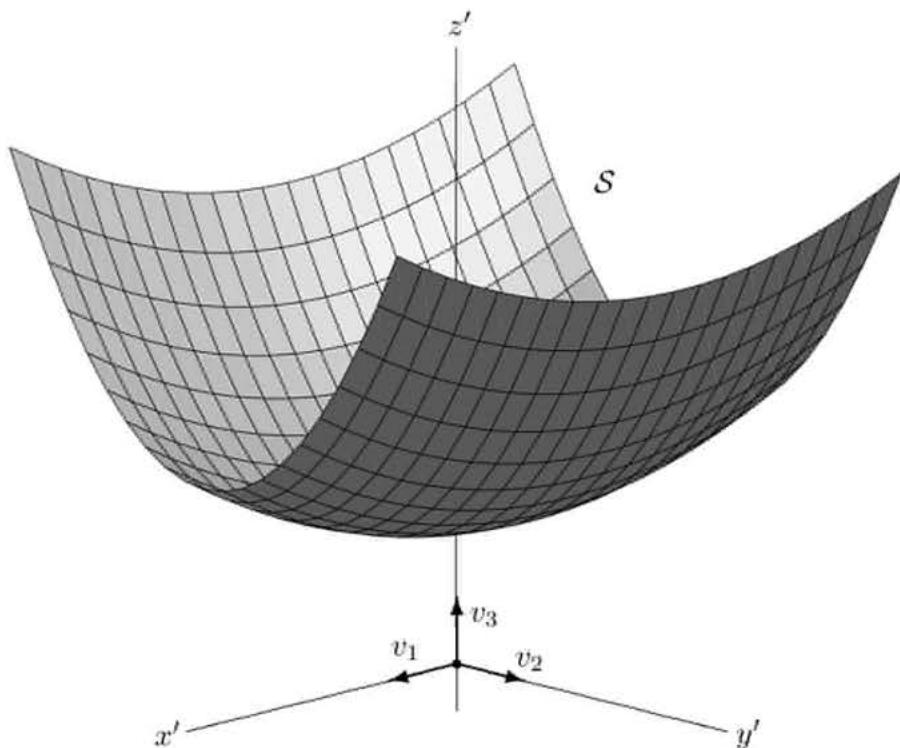
We are now ready to transform (11) into an equation involving coordinates relative to γ . Let $x = (t_1, t_2, t_3) \in \mathbb{R}^3$, and suppose that $x = s_1 v_1 + s_2 v_2 + s_3 v_3$. Then, by Theorem 2.22 (p. 112),

$$x = \begin{pmatrix} t_1 \\ t_2 \\ t_3 \end{pmatrix} = Q \begin{pmatrix} s_1 \\ s_2 \\ s_3 \end{pmatrix},$$

and therefore

$$t_1 = \frac{s_1}{\sqrt{10}} + \frac{3s_2}{\sqrt{35}} - \frac{3s_3}{\sqrt{14}},$$

$$t_2 = \frac{5s_2}{\sqrt{35}} + \frac{2s_3}{\sqrt{14}},$$



quadratic

Figure 6.7

and

$$t_3 = \frac{3s_1}{\sqrt{10}} - \frac{s_2}{\sqrt{35}} + \frac{s_3}{\sqrt{14}}.$$

Thus

$$3t_1 - 2t_2 - t_3 = -\frac{14s_3}{\sqrt{14}} = -\sqrt{14}s_3.$$

Combining (11), (12), and the preceding equation, we conclude that if $x \in \mathbb{R}^3$ and $x = s_1v_1 + s_2v_2 + s_3v_3$, then $x \in S$ if and only if

$$2s_1^2 + 7s_2^2 - \sqrt{14}s_3 + 14 = 0 \quad \text{or} \quad s_3 = \frac{\sqrt{14}}{7}s_1^2 + \frac{\sqrt{14}}{2}s_2^2 + \sqrt{14}.$$

Consequently, if we draw new axes x' , y' , and z' in the directions of v_1 , v_2 , and v_3 , respectively, the graph of the equation, rewritten as

$$z' = \frac{\sqrt{14}}{7}(x')^2 + \frac{\sqrt{14}}{2}(y')^2 + \sqrt{14},$$

coincides with the surface \mathcal{S} . We recognize \mathcal{S} to be an elliptic paraboloid.

Figure 6.7 is a sketch of the surface \mathcal{S} drawn so that the vectors v_1 , v_2 and v_3 are oriented to lie in the principal directions. For practical purposes, the scale of the z' axis has been adjusted so that the figure fits the page. ♦

The Second Derivative Test for Functions of Several Variables

We now consider an application of the theory of quadratic forms to multivariable calculus—the derivation of the second derivative test for local extrema of a function of several variables. We assume an acquaintance with the calculus of functions of several variables to the extent of Taylor's theorem. The reader is undoubtedly familiar with the one-variable version of Taylor's theorem. For a statement and proof of the multivariable version, consult, for example, *An Introduction to Analysis* 2d ed, by William R. Wade (Prentice Hall, Upper Saddle River, N.J., 2000).

Let $z = f(t_1, t_2, \dots, t_n)$ be a fixed real-valued function of n real variables for which all third-order partial derivatives exist and are continuous. The function f is said to have a **local maximum** at a point $p \in \mathbb{R}^n$ if there exists a $\delta > 0$ such that $f(p) \geq f(x)$ whenever $\|x - p\| < \delta$. Likewise, f has a **local minimum** at $p \in \mathbb{R}^n$ if there exists a $\delta > 0$ such that $f(p) \leq f(x)$ whenever $\|x - p\| < \delta$. If f has either a local minimum or a local maximum at p , we say that f has a **local extremum** at p . A point $p \in \mathbb{R}^n$ is called a **critical point** of f if $\partial f(p)/\partial t_i = 0$ for $i = 1, 2, \dots, n$. It is a well-known fact that if f has a local extremum at a point $p \in \mathbb{R}^n$, then p is a critical point of f . For, if f has a local extremum at $p = (p_1, p_2, \dots, p_n)$, then for any $i = 1, 2, \dots, n$ the function ϕ_i defined by $\phi_i(t) = f(p_1, p_2, \dots, p_{i-1}, t, p_{i+1}, \dots, p_n)$ has a local extremum at $t = p_i$. So, by an elementary single-variable argument,

$$\frac{\partial f(p)}{\partial t_i} = \frac{d\phi_i(p_i)}{dt} = 0.$$

Thus p is a critical point of f . But critical points are not necessarily local extrema.

The second-order partial derivatives of f at a critical point p can often be used to test for a local extremum at p . These partials determine a matrix $A(p)$ in which the row i , column j entry is

$$\frac{\partial^2 f(p)}{(\partial t_i)(\partial t_j)}.$$

This matrix is called the **Hessian matrix** of f at p . Note that if the third-order partial derivatives of f are continuous, then the mixed second-order partials of f at p are independent of the order in which they are taken, and hence $A(p)$ is a symmetric matrix. In this case, all of the eigenvalues of $A(p)$ are real.

Theorem 6.37 (The Second Derivative Test). Let $f(t_1, t_2, \dots, t_n)$ be a real-valued function in n real variables for which all third-order partial derivatives exist and are continuous. Let $p = (p_1, p_2, \dots, p_n)$ be a critical point of f , and let $A(p)$ be the Hessian of f at p .

- (a) If all eigenvalues of $A(p)$ are positive, then f has a local minimum at p .
- (b) If all eigenvalues of $A(p)$ are negative, then f has a local maximum at p .
- (c) If $A(p)$ has at least one positive and at least one negative eigenvalue, then f has no local extremum at p (p is called a **saddle-point** of f).
- (d) If $\text{rank}(A(p)) < n$ and $A(p)$ does not have both positive and negative eigenvalues, then the second derivative test is inconclusive.

Proof. If $p \neq 0$, we may define a function $g: \mathbb{R}^n \rightarrow \mathbb{R}$ by

$$g(t_1, t_2, \dots, t_n) = f(t_1 + p_1, t_2 + p_2, \dots, t_n + p_n) - f(p).$$

The following facts are easily verified.

1. The function f has a local maximum [minimum] at p if and only if g has a local maximum [minimum] at $0 = (0, 0, \dots, 0)$.
2. The partial derivatives of g at 0 are equal to the corresponding partial derivatives of f at p .
3. 0 is a critical point of g .
4. $A_{ij}(p) = \frac{\partial^2 g(0)}{\partial t_i \partial t_j}$ for all i and j .

In view of these facts, we may assume without loss of generality that $p = 0$ and $f(p) = 0$.

Now we apply Taylor's theorem to f to obtain the first-order approximation of f around 0 . We have

$$\begin{aligned} f(t_1, t_2, \dots, t_n) &= f(0) + \sum_{i=1}^n \frac{\partial f(0)}{\partial t_i} t_i + \frac{1}{2} \sum_{i,j=1}^n \frac{\partial^2 f(0)}{\partial t_i \partial t_j} t_i t_j + S(t_1, t_2, \dots, t_n) \\ &= \frac{1}{2} \sum_{i,j=1}^n \frac{\partial^2 f(0)}{\partial t_i \partial t_j} t_i t_j + S(t_1, t_2, \dots, t_n), \end{aligned} \tag{13}$$

where S is a real-valued function on \mathbb{R}^n such that

$$\lim_{x \rightarrow 0} \frac{S(x)}{\|x\|^2} = \lim_{(t_1, t_2, \dots, t_n) \rightarrow 0} \frac{S(t_1, t_2, \dots, t_n)}{t_1^2 + t_2^2 + \dots + t_n^2} = 0. \tag{14}$$

Let $K: \mathbb{R}^n \rightarrow \mathbb{R}$ be the quadratic form defined by

$$K \begin{pmatrix} t_1 \\ t_2 \\ \vdots \\ t_n \end{pmatrix} = \frac{1}{2} \sum_{i,j=1}^n \frac{\partial^2 f(0)}{\partial t_i \partial t_j} t_i t_j, \tag{15}$$

H be the symmetric bilinear form corresponding to K , and β be the standard ordered basis for \mathbb{R}^n . It is easy to verify that $\psi_\beta(H) = \frac{1}{2}A(p)$. Since $A(p)$ is symmetric, Theorem 6.20 (p. 381) implies that there exists an orthogonal matrix Q such that

$$Q^t A(p) Q = \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \dots & \lambda_n \end{pmatrix}$$

is a diagonal matrix whose diagonal entries are the eigenvalues of $A(p)$. Let $\gamma = \{v_1, v_2, \dots, v_n\}$ be the orthogonal basis for \mathbb{R}^n whose i th vector is the i th column of Q . Then Q is the change of coordinate matrix changing γ -coordinates into β -coordinates, and by Theorem 6.33

$$\psi_\gamma(H) = Q^t \psi_\beta(H) Q = \frac{1}{2} Q^t A(p) Q = \begin{pmatrix} \frac{\lambda_1}{2} & 0 & \dots & 0 \\ 0 & \frac{\lambda_2}{2} & \dots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \dots & \frac{\lambda_n}{2} \end{pmatrix}.$$

Suppose that $A(p)$ is not the zero matrix. Then $A(p)$ has nonzero eigenvalues. Choose $\epsilon > 0$ such that $\epsilon < |\lambda_i|/2$ for all $\lambda_i \neq 0$. By (14), there exists $\delta > 0$ such that for any $x \in \mathbb{R}^n$ satisfying $0 < \|x\| < \delta$, we have $|S(x)| < \epsilon \|x\|^2$. Consider any $x \in \mathbb{R}^n$ such that $0 < \|x\| < \delta$. Then, by (13) and (15),

$$|f(x) - K(x)| = |S(x)| < \epsilon \|x\|^2,$$

and hence

$$K(x) - \epsilon \|x\|^2 < f(x) < K(x) + \epsilon \|x\|^2. \quad (16)$$

Suppose that $x = \sum_{i=1}^n s_i v_i$. Then

$$\|x\|^2 = \sum_{i=1}^n s_i^2 \quad \text{and} \quad K(x) = \frac{1}{2} \sum_{i=1}^n \lambda_i s_i^2.$$

Combining these equations with (16), we obtain

$$\sum_{i=1}^n \left(\frac{1}{2} \lambda_i - \epsilon \right) s_i^2 < f(x) < \sum_{i=1}^n \left(\frac{1}{2} \lambda_i + \epsilon \right) s_i^2. \quad (17)$$

Now suppose that all eigenvalues of $A(p)$ are positive. Then $\frac{1}{2}\lambda_i - \epsilon > 0$ for all i , and hence, by the left inequality in (17),

$$f(\theta) = 0 \leq \sum_{i=1}^n \left(\frac{1}{2}\lambda_i - \epsilon \right) s_i^2 < f(x).$$

Thus $f(\theta) \leq f(x)$ for $\|x\| < \delta$, and so f has a local minimum at θ . By a similar argument using the right inequality in (17), we have that if all of the eigenvalues of $A(p)$ are negative, then f has a local maximum at θ . This establishes (a) and (b) of the theorem.

Next, suppose that $A(p)$ has both a positive and a negative eigenvalue, say, $\lambda_i > 0$ and $\lambda_j < 0$ for some i and j . Then $\frac{1}{2}\lambda_i - \epsilon > 0$ and $\frac{1}{2}\lambda_j + \epsilon < 0$. Let s be any real number such that $0 < |s| < \delta$. Substituting $x = sv_i$ and $x = sv_j$ into the left inequality and the right inequality of (17), respectively, we obtain

$$f(\theta) = 0 < \left(\frac{1}{2}\lambda_i - \epsilon \right) s^2 < f(sv_i) \quad \text{and} \quad f(sv_j) < \left(\frac{1}{2}\lambda_j + \epsilon \right) s^2 < 0 = f(\theta).$$

Thus f attains both positive and negative values arbitrarily close to θ ; so f has neither a local maximum nor a local minimum at θ . This establishes (c).

To show that the second-derivative test is inconclusive under the conditions stated in (d), consider the functions

$$f(t_1, t_2) = t_1^2 - t_2^4 \quad \text{and} \quad g(t_1, t_2) = t_1^2 + t_2^4$$

at $p = \theta$. In both cases, the function has a critical point at p , and

$$A(p) = \begin{pmatrix} 2 & 0 \\ 0 & 0 \end{pmatrix}.$$

However, f does not have a local extremum at θ , whereas g has a local minimum at θ . ■

Sylvester's Law of Inertia

Any two matrix representations of a bilinear form have the same rank because rank is preserved under congruence. We can therefore define the **rank** of a bilinear form to be the rank of any of its matrix representations. If a matrix representation is a diagonal matrix, then the rank is equal to the number of nonzero diagonal entries of the matrix.

We confine our analysis to symmetric bilinear forms on finite-dimensional real vector spaces. Each such form has a diagonal matrix representation in which the diagonal entries may be positive, negative, or zero. Although these entries are not unique, we show that the number of entries that are positive and the number that are negative are unique. That is, they are independent

of the choice of diagonal representation. This result is called *Sylvester's law of inertia*. We prove the law and apply it to describe the equivalence classes of congruent symmetric real matrices.

Theorem 6.38 (Sylvester's Law of Inertia). *Let H be a symmetric bilinear form on a finite-dimensional real vector space V . Then the number of positive diagonal entries and the number of negative diagonal entries in any diagonal matrix representation of H are each independent of the diagonal representation.*

Proof. Suppose that β and γ are ordered bases for V that determine diagonal representations of H . Without loss of generality, we may assume that β and γ are ordered so that on each diagonal the entries are in the order of positive, negative, and zero. It suffices to show that both representations have the same number of positive entries because the number of negative entries is equal to the difference between the rank and the number of positive entries. Let p and q be the number of positive diagonal entries in the matrix representations of H with respect to β and γ , respectively. We suppose that $p \neq q$ and arrive at a contradiction. Without loss of generality, assume that $p < q$. Let

$$\beta = \{v_1, v_2, \dots, v_p, \dots, v_r, \dots, v_n\} \text{ and } \gamma = \{w_1, w_2, \dots, w_q, \dots, w_r, \dots, w_n\},$$

where r is the rank of H and $n = \dim(V)$. Let $L: V \rightarrow \mathbb{R}^{p+r-q}$ be the mapping defined by

$$L(x) = (H(x, v_1), H(x, v_2), \dots, H(x, v_p), H(x, w_{q+1}), \dots, H(x, w_r)).$$

It is easily verified that L is linear and $\text{rank}(L) \leq p + r - q$. Hence

$$\text{nullity}(L) \geq n - (p + r - q) > n - r.$$

So there exists a nonzero vector v_0 such that $v_0 \notin \text{span}(\{v_{r+1}, v_{r+2}, \dots, v_n\})$, but $v_0 \in N(L)$. Since $v_0 \in N(L)$, it follows that $H(v_0, v_i) = 0$ for $i \leq p$ and $H(v_0, w_i) = 0$ for $q < i \leq r$. Suppose that

$$v_0 = \sum_{j=1}^n a_j v_j = \sum_{j=1}^n b_j w_j.$$

For any $i \leq p$,

$$H(v_0, v_i) = H\left(\sum_{j=1}^n a_j v_j, v_i\right) = \sum_{j=1}^n a_j H(v_j, v_i) = a_i H(v_i, v_i).$$

But for $i \leq p$, we have $H(v_i, v_i) > 0$ and $H(v_0, v_i) = 0$, so that $a_i = 0$. Similarly, $b_i = 0$ for $q+1 \leq i \leq r$. Since v_0 is not in the span of $\{v_{r+1}, v_{r+2}, \dots, v_n\}$, it follows that $a_i \neq 0$ for some $p < i \leq r$. Thus

$$H(v_0, v_0) = H\left(\sum_{j=1}^n a_j v_j, \sum_{i=1}^n a_i v_i\right) = \sum_{j=1}^n a_j^2 H(v_j, v_j) = \sum_{j=p+1}^r a_j^2 H(v_j, v_j) < 0.$$

Furthermore,

$$H(v_0, v_0) = H\left(\sum_{j=1}^n b_j w_j, \sum_{i=1}^n b_i w_i\right) = \sum_{j=1}^n b_j^2 H(w_j, w_j) = \sum_{j=p+1}^r b_j^2 H(w_j, w_j) \geq 0.$$

So $H(v_0, v_0) < 0$ and $H(v_0, v_0) \geq 0$, which is a contradiction. We conclude that $p = q$. ■

Definitions. The number of positive diagonal entries in a diagonal representation of a symmetric bilinear form on a real vector space is called the **index** of the form. The difference between the number of positive and the number of negative diagonal entries in a diagonal representation of a symmetric bilinear form is called the **signature** of the form. The three terms *rank*, *index*, and *signature* are called the **invariants** of the bilinear form because they are invariant with respect to matrix representations. These same terms apply to the associated quadratic form. Notice that the values of any two of these invariants determine the value of the third.

Example 10

The bilinear form corresponding to the quadratic form K of Example 9 has a 3×3 diagonal matrix representation with diagonal entries of 2, 7, and 0. Therefore the rank, index, and signature of K are each 2. ◆

Example 11

The matrix representation of the bilinear form corresponding to the quadratic form $K(x, y) = x^2 - y^2$ on \mathbb{R}^2 with respect to the standard ordered basis is the diagonal matrix with diagonal entries of 1 and -1 . Therefore the rank of K is 2, the index of K is 1, and the signature of K is 0. ◆

Since the congruence relation is intimately associated with bilinear forms, we can apply Sylvester's law of inertia to study this relation on the set of real symmetric matrices. Let A be an $n \times n$ real symmetric matrix, and suppose that D and E are each diagonal matrices congruent to A . By Corollary 3 to Theorem 6.32, A is the matrix representation of the bilinear form H on \mathbb{R}^n defined by $H(x, y) = x^t A y$ with respect to the standard ordered basis for \mathbb{R}^n . Therefore Sylvester's law of inertia tells us that D and E have the same number of positive and negative diagonal entries. We can state this result as the matrix version of Sylvester's law.

Corollary 1 (Sylvester's Law of Inertia for Matrices). Let A be a real symmetric matrix. Then the number of positive diagonal entries and the number of negative diagonal entries in any diagonal matrix congruent to A is independent of the choice of the diagonal matrix.

Definitions. Let A be a real symmetric matrix, and let D be a diagonal matrix that is congruent to A . The number of positive diagonal entries of D is called the **index** of A . The difference between the number of positive diagonal entries and the number of negative diagonal entries of D is called the **signature** of A . As before, the **rank**, **index**, and **signature** of a matrix are called the **invariants** of the matrix, and the values of any two of these invariants determine the value of the third.

Any two of these invariants can be used to determine an equivalence class of congruent real symmetric matrices.

Corollary 2. Two real symmetric $n \times n$ matrices are congruent if and only if they have the same invariants.

Proof. If A and B are congruent $n \times n$ symmetric matrices, then they are both congruent to the same diagonal matrix, and it follows that they have the same invariants.

Conversely, suppose that A and B are $n \times n$ symmetric matrices with the same invariants. Let D and E be diagonal matrices congruent to A and B , respectively, chosen so that the diagonal entries are in the order of positive, negative, and zero. (Exercise 23 allows us to do this.) Since A and B have the same invariants, so do D and E . Let p and r denote the index and the rank, respectively, of both D and E . Let d_i denote the i th diagonal entry of D , and let Q be the $n \times n$ diagonal matrix whose i th diagonal entry q_i is given by

$$q_i = \begin{cases} \frac{1}{\sqrt{d_i}} & \text{if } 1 \leq i \leq p \\ \frac{1}{\sqrt{-d_i}} & \text{if } p < i \leq r \\ 1 & \text{if } r < i. \end{cases}$$

Then $Q^t D Q = J_{pr}$, where

$$J_{pr} = \begin{pmatrix} I_p & O & O \\ O & -I_{r-p} & O \\ O & O & O \end{pmatrix}.$$

It follows that A is congruent to J_{pr} . Similarly, B is congruent to J_{pr} , and hence A is congruent to B . ■

The matrix J_{pr} acts as a canonical form for the theory of real symmetric matrices. The next corollary, whose proof is contained in the proof of Corollary 2, describes the role of J_{pr} .

Corollary 3. A real symmetric $n \times n$ matrix A has index p and rank r if and only if A is congruent to J_{pr} (as just defined).

Example 12

Let

$$A = \begin{pmatrix} 1 & 1 & -3 \\ -1 & 2 & 1 \\ 3 & 1 & 1 \end{pmatrix}, \quad B = \begin{pmatrix} 1 & 2 & 1 \\ 2 & 3 & 2 \\ 1 & 2 & 1 \end{pmatrix}, \quad \text{and} \quad C = \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & 2 \\ 1 & 2 & 1 \end{pmatrix}.$$

We apply Corollary 2 to determine which pairs of the matrices A , B , and C are congruent.

The matrix A is the 3×3 matrix of Example 6, where it is shown that A is congruent to a diagonal matrix with diagonal entries 1, 1, and -24 . Therefore A has rank 3 and index 2. Using the methods of Example 6 (it is not necessary to compute Q), it can be shown that B and C are congruent, respectively, to the diagonal matrices

$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & -1 \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & -4 \end{pmatrix}.$$

It follows that both A and C have rank 3 and index 2, while B has rank 3 and index 1. We conclude that A and C are congruent but that B is congruent to neither A nor C . ◆

EXERCISES

1. Label the following statements as true or false.
 - (a) Every quadratic form is a bilinear form.
 - (b) If two matrices are congruent, they have the same eigenvalues.
 - (c) Symmetric bilinear forms have symmetric matrix representations.
 - (d) Any symmetric matrix is congruent to a diagonal matrix.
 - (e) The sum of two symmetric bilinear forms is a symmetric bilinear form.
 - (f) If two symmetric matrices over a field not of characteristic two have the same characteristic polynomial, then they are matrix representations of the same bilinear form.
 - (g) There exists a bilinear form H such that $H(x, y) \neq 0$ for all x and y .

- (h) If V is a vector space of dimension n , then $\dim(\mathcal{B}(V)) = 2n$.
- (i) Let H be a bilinear form on a finite-dimensional vector space V with $\dim(V) > 1$. For any $x \in V$, there exists $y \in V$ such that $y \neq 0$, but $H(x, y) = 0$.
- (j) If H is any bilinear form on a finite-dimensional real inner product space V , then there exists an ordered basis β for V such that $\psi_\beta(H)$ is a diagonal matrix.
2. Prove properties 1, 2, 3, and 4 on page 420.
3. (a) Prove that the sum of two bilinear forms is a bilinear form.
 (b) Prove that the product of a scalar and a bilinear form is a bilinear form.
 (c) Prove Theorem 6.31.
4. Determine which of the mappings that follow are bilinear forms. Justify your answers.
- (a) Let $V = C[0, 1]$ be the space of continuous real-valued functions on the closed interval $[0, 1]$. For $f, g \in V$, define
- $$H(f, g) = \int_0^1 f(t)g(t)dt.$$
- (b) Let V be a vector space over F , and let $J \in \mathcal{B}(V)$ be nonzero. Define $H: V \times V \rightarrow F$ by
- $$H(x, y) = [J(x, y)]^2 \quad \text{for all } x, y \in V.$$
- (c) Define $H: R \times R \rightarrow R$ by $H(t_1, t_2) = t_1 + 2t_2$.
 (d) Consider the vectors of R^2 as column vectors, and let $H: R^2 \rightarrow R$ be the function defined by $H(x, y) = \det(x, y)$, the determinant of the 2×2 matrix with columns x and y .
 (e) Let V be a real inner product space, and let $H: V \times V \rightarrow R$ be the function defined by $H(x, y) = \langle x, y \rangle$ for $x, y \in V$.
 (f) Let V be a complex inner product space, and let $H: V \times V \rightarrow C$ be the function defined by $H(x, y) = \langle x, y \rangle$ for $x, y \in V$.
5. Verify that each of the given mappings is a bilinear form. Then compute its matrix representation with respect to the given ordered basis β .
- (a) $H: R^3 \times R^3 \rightarrow R$, where
- $$H\left(\begin{pmatrix} a_1 \\ a_2 \\ a_3 \end{pmatrix}, \begin{pmatrix} b_1 \\ b_2 \\ b_3 \end{pmatrix}\right) = a_1b_1 - 2a_1b_2 + a_2b_1 - a_3b_3$$

and

$$\beta = \left\{ \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \right\}.$$

- (b) Let $V = M_{2 \times 2}(R)$ and

$$\beta = \left\{ \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}, \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}, \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix}, \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} \right\}.$$

Define $H: V \times V \rightarrow R$ by $H(A, B) = \text{tr}(A) \cdot \text{tr}(B)$.

- (c) Let $\beta = \{\cos t, \sin t, \cos 2t, \sin 2t\}$. Then β is an ordered basis for $V = \text{span}(\beta)$, a four-dimensional subspace of the space of all continuous functions on R . Let $H: V \times V \rightarrow R$ be the function defined by $H(f, g) = f'(0) \cdot g''(0)$.

6. Let $H: R^2 \times R^2 \rightarrow R$ be the function defined by

$$H \left(\begin{pmatrix} a_1 \\ a_2 \end{pmatrix}, \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} \right) = a_1 b_2 + a_2 b_1 \quad \text{for } \begin{pmatrix} a_1 \\ a_2 \end{pmatrix}, \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} \in R^2.$$

- (a) Prove that H is a bilinear form.
 (b) Find the 2×2 matrix A such that $H(x, y) = x^t A y$ for all $x, y \in R^2$.
7. Let V and W be vector spaces over the same field, and let $T: V \rightarrow W$ be a linear transformation. For any $H \in \mathcal{B}(W)$, define $\widehat{T}(H): V \times V \rightarrow F$ by $\widehat{T}(H)(x, y) = H(T(x), T(y))$ for all $x, y \in V$. Prove the following results.
- (a) If $H \in \mathcal{B}(W)$, then $\widehat{T}(H) \in \mathcal{B}(V)$.
 (b) $\widehat{T}: \mathcal{B}(W) \rightarrow \mathcal{B}(V)$ is a linear transformation.
 (c) If T is an isomorphism, then so is \widehat{T} .
8. Assume the notation of Theorem 6.32.
- (a) Prove that for any ordered basis β , ψ_β is linear.
 (b) Let β be an ordered basis for an n -dimensional space V over F , and let $\phi_\beta: V \rightarrow F^n$ be the standard representation of V with respect to β . For $A \in M_{n \times n}(F)$, define $H: V \times V \rightarrow F$ by $H(x, y) = [\phi_\beta(x)]^t A [\phi_\beta(y)]$. Prove that $H \in \mathcal{B}(V)$. Can you establish this as a corollary to Exercise 7?
 (c) Prove the converse of (b): Let H be a bilinear form on V . If $A = \psi_\beta(H)$, then $H(x, y) = [\phi_\beta(x)]^t A [\phi_\beta(y)]$.
9. (a) Prove Corollary 1 to Theorem 6.32.
 (b) For a finite-dimensional vector space V , describe a method for finding an ordered basis for $\mathcal{B}(V)$.

10. Prove Corollary 2 to Theorem 6.32.
11. Prove Corollary 3 to Theorem 6.32.
12. Prove that the relation of congruence is an equivalence relation.
13. Use Corollary 2 to Theorem 6.32 and Theorem 2.22(b) to obtain an alternate proof to Theorem 6.33.
14. Let V be a finite-dimensional vector space and $H \in \mathcal{B}(V)$. Prove that, for any ordered bases β and γ of V , $\text{rank}(\psi_\beta(H)) = \text{rank}(\psi_\gamma(H))$.
15. Prove the following results.
 - (a) Any square diagonal matrix is symmetric.
 - (b) Any matrix congruent to a diagonal matrix is symmetric.
 - (c) the corollary to Theorem 6.35
16. Let V be a vector space over a field F not of characteristic two, and let H be a symmetric bilinear form on V . Prove that if $K(x) = H(x, x)$ is the quadratic form associated with H , then, for all $x, y \in V$,

$$H(x, y) = \frac{1}{2}[K(x + y) - K(x) - K(y)].$$

17. For each of the given quadratic forms K on a real inner product space V , find a symmetric bilinear form H such that $K(x) = H(x, x)$ for all $x \in V$. Then find an orthonormal basis β for V such that $\psi_\beta(H)$ is a diagonal matrix.
 - (a) $K: \mathbb{R}^2 \rightarrow \mathbb{R}$ defined by $K \begin{pmatrix} t_1 \\ t_2 \end{pmatrix} = -2t_1^2 + 4t_1t_2 + t_2^2$
 - (b) $K: \mathbb{R}^2 \rightarrow \mathbb{R}$ defined by $K \begin{pmatrix} t_1 \\ t_2 \end{pmatrix} = 7t_1^2 - 8t_1t_2 + t_2^2$
 - (c) $K: \mathbb{R}^3 \rightarrow \mathbb{R}$ defined by $K \begin{pmatrix} t_1 \\ t_2 \\ t_3 \end{pmatrix} = 3t_1^2 + 3t_2^2 + 3t_3^2 - 2t_1t_3$
18. Let S be the set of all $(t_1, t_2, t_3) \in \mathbb{R}^3$ for which

$$3t_1^2 + 3t_2^2 + 3t_3^2 - 2t_1t_3 + 2\sqrt{2}(t_1 + t_3) + 1 = 0.$$

Find an orthonormal basis β for \mathbb{R}^3 for which the equation relating the coordinates of points of S relative to β is simpler. Describe S geometrically.

19. Prove the following refinement of Theorem 6.37(d).

- (a) If $0 < \text{rank}(A) < n$ and A has no negative eigenvalues, then f has no local maximum at p .
- (b) If $0 < \text{rank}(A) < n$ and A has no positive eigenvalues, then f has no local minimum at p .
20. Prove the following variation of the second-derivative test for the case $n = 2$: Define

$$D = \left[\frac{\partial^2 f(p)}{\partial t_1^2} \right] \left[\frac{\partial^2 f(p)}{\partial t_2^2} \right] - \left[\frac{\partial^2 f(p)}{\partial t_1 \partial t_2} \right]^2.$$

- (a) If $D > 0$ and $\partial^2 f(p)/\partial t_1^2 > 0$, then f has a local minimum at p .
- (b) If $D > 0$ and $\partial^2 f(p)/\partial t_1^2 < 0$, then f has a local maximum at p .
- (c) If $D < 0$, then f has no local extremum at p .
- (d) If $D = 0$, then the test is inconclusive.

Hint: Observe that, as in Theorem 6.37, $D = \det(A) = \lambda_1 \lambda_2$, where λ_1 and λ_2 are the eigenvalues of A .

21. Let A and E be in $M_{n \times n}(F)$, with E an elementary matrix. In Section 3.1, it was shown that AE can be obtained from A by means of an elementary column operation. Prove that $E^t A$ can be obtained by means of the same elementary operation performed on the rows rather than on the columns of A . *Hint:* Note that $E^t A = (A^t E)^t$.
22. For each of the following matrices A with entries from R , find a diagonal matrix D and an invertible matrix Q such that $Q^t A Q = D$.

$$(a) \begin{pmatrix} 1 & 3 \\ 3 & 2 \end{pmatrix} \quad (b) \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \quad (c) \begin{pmatrix} 3 & 1 & 2 \\ 1 & 4 & 0 \\ 2 & 0 & -1 \end{pmatrix}$$

Hint for (b): Use an elementary operation other than interchanging columns.

23. Prove that if the diagonal entries of a diagonal matrix are permuted, then the resulting diagonal matrix is congruent to the original one.
24. Let T be a linear operator on a real inner product space V , and define $H: V \times V \rightarrow R$ by $H(x, y) = \langle x, T(y) \rangle$ for all $x, y \in V$.
- (a) Prove that H is a bilinear form.
- (b) Prove that H is symmetric if and only if T is self-adjoint.
- (c) What properties must T have for H to be an inner product on V ?
- (d) Explain why H may fail to be a bilinear form if V is a complex inner product space.
25. Prove the converse to Exercise 24(a): Let V be a finite-dimensional real inner product space, and let H be a bilinear form on V . Then there

exists a unique linear operator T on V such that $H(x, y) = \langle x, T(y) \rangle$ for all $x, y \in V$. Hint: Choose an orthonormal basis β for V , let $A = \psi_\beta(H)$, and let T be the linear operator on V such that $[T]_\beta = A$. Visit goo.gl/bGAfSy for a solution.

26. Prove that the number of distinct equivalence classes of congruent $n \times n$ real symmetric matrices is

$$\frac{(n+1)(n+2)}{2}.$$

6.9* EINSTEIN'S SPECIAL THEORY OF RELATIVITY

As a consequence of physical experiments performed in the latter half of the nineteenth century (most notably the Michelson–Morley experiment of 1887), physicists concluded that *the results obtained in measuring the speed of light c are independent of the velocity of the instrument used to measure it*. For example, suppose that while on Earth, an experimenter measures the speed of light emitted from the sun and finds it to be 186,000 miles per second. Now suppose that the experimenter places the measuring equipment in a spaceship that leaves Earth traveling at 100,000 miles per second in a direction away from the sun. A repetition of the same experiment from the spaceship yields the same result: Light is traveling at 186,000 miles per second relative to the spaceship, rather than 86,000 miles per second as one might expect!

This revelation led to a new way of relating coordinate systems used to locate events in space–time. The result was Albert Einstein's *special theory of relativity*. In this section, we develop via a linear algebra viewpoint the essence of Einstein's theory.

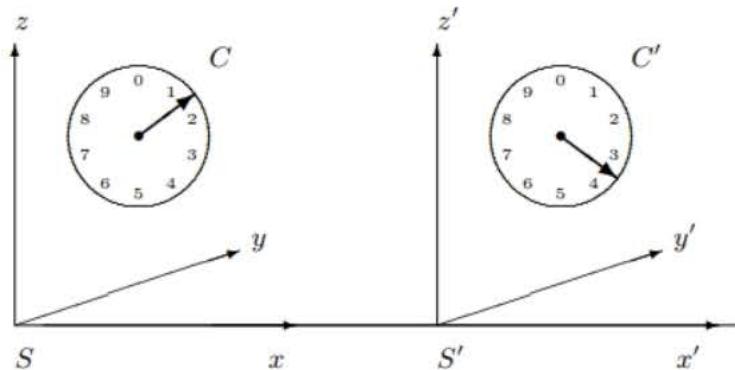


Figure 6.8

The basic problem is to compare two different inertial (nonaccelerating) coordinate systems S and S' that are in motion relative to each other under

the assumption that the speed of light is the same when measured in either system. We assume that S' moves at a constant velocity in relation to S as measured from S .

To simplify matters, we assume that the two coordinate systems have parallel axes and share the same x -axis, and that the motion of S' relative to S is along this common axis. (See Figure 6.8.)

We also suppose that there are two clocks, C and C' , placed in space so that C is stationary relative to S and C' is stationary relative to S' . These clocks give readings that are real numbers in units of seconds. They are calibrated so that at the instant the origins of S and S' coincide, both give the reading of 0.

Given any event p (something whose position and time of occurrence can be described), we may assign a set of *space-time coordinates* to it. For example, if p is an event that occurs at position (x, y, z) relative to S and at time t as read on clock C , we can assign to p the set of coordinates

$$\begin{pmatrix} x \\ y \\ z \\ t \end{pmatrix}.$$

This ordered 4-tuple is called the **space-time coordinates** of p relative to S and C . Likewise, p has a set of space-time coordinates

$$\begin{pmatrix} x' \\ y' \\ z' \\ t' \end{pmatrix}$$

relative to S' and C' .

Because motion is along a common x -axis, which lies in a common xy -plane, the third component of the space-time coordinates of p is always zero. Thus we consider only the first, second, and fourth coordinates of p , and write

$$\begin{pmatrix} x \\ y \\ t \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} x' \\ y' \\ t' \end{pmatrix}$$

to denote the space-time coordinates of an event p relative to S and S' , respectively.

As we have mentioned, our unit of time is the second. Our measure of an object's velocity v is the ratio of its velocity (expressed in miles per second) to the speed of light expressed in the same units (which is approximately 186,000 miles per second). For example, if S' is moving at 18,600 miles per second relative to S and the speed of light, c , is 186,000 miles per second, the

velocity of S' relative to S , v , would have a value of $v = 0.1$. For this reason, the speed of light c has the value 1.

For a fixed velocity v , let $\mathbf{T}_v: \mathbb{R}^3 \rightarrow \mathbb{R}^3$ be the mapping defined by

$$\mathbf{T}_v \begin{pmatrix} x \\ y \\ t \end{pmatrix} = \begin{pmatrix} x' \\ y' \\ t' \end{pmatrix},$$

where

$$\begin{pmatrix} x \\ y \\ t \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} x' \\ y' \\ t' \end{pmatrix}$$

are the space-time coordinates of the same event with respect to S and C and with respect to S' and C' , respectively.

In what follows, we make four assumptions:

1. The origin of S' moves in the positive direction of the common x -axis relative to S at a constant velocity of $v > 0$.
2. The origin of S moves in the negative direction of the common x -axis relative to S' at the constant velocity of $-v < 0$.
3. \mathbf{T}_v is a linear isomorphism.
4. The speed of any light beam, when measured in either S or S' , using the clocks C in S and C' in S' , is always $c = 1$.

Since motion is strictly along the x -axis and we assume that the y -axis is unaffected, we have that for any x , y , and t , there exist x' and t' such that

$$\mathbf{T}_v \begin{pmatrix} x \\ y \\ t \end{pmatrix} = \begin{pmatrix} x' \\ y \\ t' \end{pmatrix} \quad \text{and} \quad \mathbf{T}_v \begin{pmatrix} 0 \\ y \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ y \\ 0 \end{pmatrix}.$$

Our goal in this section is to calculate the matrix representation of \mathbf{T}_v with respect to the standard basis for \mathbb{R}^3 .

Theorem 6.39. Consider $\{e_1, e_2, e_3\}$, the standard ordered basis for \mathbb{R}^3 . Then

- (a) $\mathbf{T}_v(e_2) = e_2$.
- (b) \mathbf{T}_v maps $\text{span}\{e_1, e_3\}$ into itself.
- (c) \mathbf{T}_v^* maps $\text{span}\{e_1, e_3\}$ into itself.

Proof. Parts (a) and (b) follow immediately from the equations above.

For $i = 1$ and $i = 3$,

$$\langle \mathbf{T}_v^*(e_i), e_2 \rangle = \langle e_i, \mathbf{T}_v(e_2) \rangle = \langle e_i, e_2 \rangle = 0,$$

and hence (c) follows. ■

Suppose that, at the instant the origins of S and S' coincide, a light flash is emitted from their common origin. The event of the light flash when measured either relative to S and C or relative to S' and C' has space-time coordinates

$$\begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}.$$

Let P be the set of all events in the xy -plane whose space-time coordinates

$$\begin{pmatrix} x \\ y \\ t \end{pmatrix}$$

relative to S and C are such that the flash is observable in the common xy -plane at the point (x, y) (as measured relative to S) at the time t (as measured on C). Let us characterize P in terms of x, y , and t . Since the speed of light is 1, at any time $t \geq 0$ the light flash is observable from any point in the plane whose distance to the origin of S (as measured on S) is $t \cdot 1 = t$. These are precisely the points in the xy -plane with $x^2 + y^2 = t^2$, or $x^2 + y^2 - t^2 = 0$. Hence an event lies in P if and only if its space-time coordinates

$$\begin{pmatrix} x \\ y \\ t \end{pmatrix} \quad (t \geq 0)$$

relative to S and C satisfy the equation $x^2 + y^2 - t^2 = 0$. Since the speed of light when measured in either coordinate system is the same, we can characterize P in terms of the space-time coordinates relative to S' and C' similarly: An event lies in P if and only if, relative to S' and C' , its space-time coordinates

$$\begin{pmatrix} x' \\ y \\ t' \end{pmatrix} \quad (t' \geq 0)$$

satisfy the equation $(x')^2 + y^2 - (t')^2 = 0$.

Let

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

Theorem 6.40. *If $\langle L_A(w), w \rangle = 0$ for some $w \in \mathbb{R}^3$, then*

$$\langle T_v^* L_A T_v(w), w \rangle = 0.$$

Proof. Let

$$w = \begin{pmatrix} x \\ y \\ t \end{pmatrix} \in \mathbb{R}^3,$$

and suppose that $\langle \mathbf{L}_A(w), w \rangle = 0$.

CASE 1. $t \geq 0$. Since $\langle \mathbf{L}_A(w), w \rangle = x^2 + y^2 - t^2$, the vector w gives the coordinates of an event in P relative to S and C . Because

$$\begin{pmatrix} x' \\ y \\ t' \end{pmatrix}$$

are the space-time coordinates of the same event relative to S' and C' , the discussion preceding Theorem 6.40 yields

$$(x')^2 + y^2 - (t')^2 = 0.$$

Thus $\langle \mathbf{T}_v^* \mathbf{L}_A \mathbf{T}_v(w), w \rangle = \langle \mathbf{L}_A \mathbf{T}_v(w), \mathbf{T}_v(w) \rangle = (x')^2 + y^2 - (t')^2 = 0$, and the conclusion follows. ■

CASE 2. $t < 0$. The proof follows by applying case 1 to $-w$. ■

We now proceed to deduce information about \mathbf{T}_v . Let

$$w_1 = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \quad \text{and} \quad w_2 = \begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix}.$$

Clearly $\{w_1, w_2\}$ is an orthogonal basis for the span of $\{e_1, e_3\}$. The next result tells us even more.

Theorem 6.41. *There exists a nonzero scalar a such that $\mathbf{T}_v^* \mathbf{L}_A \mathbf{T}_v(w_1) = aw_2$ and $\mathbf{T}_v^* \mathbf{L}_A \mathbf{T}_v(w_2) = aw_1$.*

Proof. Because $\langle \mathbf{L}_A(w_1), w_1 \rangle = 0$, $\langle \mathbf{T}_v^* \mathbf{L}_A \mathbf{T}_v(w_1), w_1 \rangle = 0$ by Theorem 6.40. Thus $\mathbf{T}_v^* \mathbf{L}_A \mathbf{T}_v(w_1)$ is orthogonal to w_1 . Since $\{w_1, w_2\}$ is an orthogonal basis for $\text{span}\{e_1, e_3\}$ and each of \mathbf{T}_v^* , \mathbf{L}_A , and \mathbf{T}_v maps this span into itself, it follows that $\mathbf{T}_v^* \mathbf{L}_A \mathbf{T}_v(w_1)$ must be a multiple of w_2 , that is, $\mathbf{T}_v^* \mathbf{L}_A \mathbf{T}_v(w_1) = aw_2$ for some scalar a . Since \mathbf{T}_v and A are invertible, so is $\mathbf{T}_v^* \mathbf{L}_A \mathbf{T}_v$. Thus $a \neq 0$.

Similarly, there exists a nonzero scalar b such that $\mathbf{T}_v^* \mathbf{L}_A \mathbf{T}_v(w_2) = bw_1$.

Finally, we show that $a = b$. Since $\mathbf{T}_v^* \mathbf{L}_A \mathbf{T}_v(w_1) = aw_2$, we have

$$2a = \langle \mathbf{T}_v^* \mathbf{L}_A \mathbf{T}_v(w_1), w_2 \rangle = \langle w_1, \mathbf{T}_v^* \mathbf{L}_A \mathbf{T}_v(w_2) \rangle = \langle w_1, bw_1 \rangle = 2b.$$

So $a = b$. ■

Actually, $a = b = 1$, as we see in the following result.

For the rest of this section let $B_v = [\mathbf{T}_v]_\beta$, where β is the standard ordered basis for \mathbb{R}^3 .

Theorem 6.42. *Given $B_v = [\mathbf{T}_v]_\beta$, as defined above,*

- (a) $B_v^*AB_v = A$.
 (b) $T_v^*L_A T_v = L_A$.

Proof. Since

$$e_1 = \frac{1}{2}(w_1 + w_2) \quad \text{and} \quad e_3 = \frac{1}{2}(w_1 - w_2),$$

it follows from Theorem 6.41 that

$$T_v^*L_A T_v(e_1) = ae_1 \quad \text{and} \quad T_v^*L_A T_v(e_3) = -ae_3.$$

Furthermore, $T_v^*L_A T_v(e_2) = e_2$, and hence

$$B_v^*AB_v = \begin{pmatrix} a & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & -a \end{pmatrix}.$$

Let

$$w = \begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix}.$$

Then $\langle L_A(w), w \rangle = 0$, and hence by Theorem 6.40

$$0 = \langle T_v^*L_A T_v(w), w \rangle = \langle B_v^*AB_v w, w \rangle = 1 - a.$$

Thus $a = 1$. As a consequence, $B_v^*AB_v = A$. This proves (a). Part (b) now follows. ■

Now consider the situation 1 second after the origins of S and S' have coincided as measured by the clock C . Since the origin of S' is moving along the x -axis at a velocity v as measured in S , its space-time coordinates relative to S and C are

$$\begin{pmatrix} v \\ 0 \\ 1 \end{pmatrix}.$$

Similarly, the space-time coordinates for the origin of S' relative to S' and C' must be

$$\begin{pmatrix} 0 \\ 0 \\ t' \end{pmatrix}$$

for some $t' > 0$. Thus we have

$$T_v \begin{pmatrix} v \\ 0 \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ t' \end{pmatrix} \quad \text{for some } t' > 0. \quad (18)$$

By Theorem 6.42

$$\left\langle \mathsf{T}_v^* \mathsf{L}_A \mathsf{T}_v \begin{pmatrix} v \\ 0 \\ 1 \end{pmatrix}, \begin{pmatrix} v \\ 0 \\ 1 \end{pmatrix} \right\rangle = \left\langle \mathsf{L}_A \begin{pmatrix} v \\ 0 \\ 1 \end{pmatrix}, \begin{pmatrix} v \\ 0 \\ 1 \end{pmatrix} \right\rangle = v^2 - 1. \quad (19)$$

But also

$$\begin{aligned} \left\langle \mathsf{T}_v^* \mathsf{L}_A \mathsf{T}_v \begin{pmatrix} v \\ 0 \\ 1 \end{pmatrix}, \begin{pmatrix} v \\ 0 \\ 1 \end{pmatrix} \right\rangle &= \left\langle \mathsf{L}_A \mathsf{T}_v \begin{pmatrix} v \\ 0 \\ 1 \end{pmatrix}, \mathsf{T}_v \begin{pmatrix} v \\ 0 \\ 1 \end{pmatrix} \right\rangle \\ &= \left\langle \mathsf{L}_A \begin{pmatrix} 0 \\ 0 \\ t' \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ t' \end{pmatrix} \right\rangle = -(t')^2. \end{aligned} \quad (20)$$

Combining (19) and (20), we conclude that $v^2 - 1 = -(t')^2$, or

$$t' = \sqrt{1 - v^2}. \quad (21)$$

Thus, from (18) and (21), we obtain

$$\mathsf{T}_v \begin{pmatrix} v \\ 0 \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ \sqrt{1 - v^2} \end{pmatrix}. \quad (22)$$

Next, recall that the origin of S moves in the negative direction of the x' -axis of S' at the constant velocity $-v < 0$ as measured from S' . Consequently, 1 second after the origins of S and S' have coincided as measured on clock C , there exists a time $t'' > 0$ as measured on clock C' such that

$$\mathsf{T}_v \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} = \begin{pmatrix} -vt'' \\ 0 \\ t'' \end{pmatrix}. \quad (23)$$

From (23), it follows in a manner similar to the derivation of (22) that

$$t'' = \frac{1}{\sqrt{1 - v^2}}. \quad (24)$$

Hence, from (23) and (24),

$$\mathsf{T}_v \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} = \begin{pmatrix} \frac{-v}{\sqrt{1 - v^2}} \\ 0 \\ \frac{1}{\sqrt{1 - v^2}} \end{pmatrix}. \quad (25)$$

The following result is now easily proved using (22), (25), and Theorem 6.39.

Theorem 6.43. Let β be the standard ordered basis for \mathbb{R}^3 . Then

$$[\mathbf{T}_v]_{\beta} = B_v = \begin{pmatrix} \frac{1}{\sqrt{1-v^2}} & 0 & \frac{-v}{\sqrt{1-v^2}} \\ 0 & 1 & 0 \\ \frac{-v}{\sqrt{1-v^2}} & 0 & \frac{1}{\sqrt{1-v^2}} \end{pmatrix}.$$

Time Contraction

A most curious and paradoxical conclusion follows if we accept Einstein's theory. Suppose that an astronaut leaves our solar system in a space vehicle traveling at a fixed velocity v as measured relative to our solar system. It follows from Einstein's theory that, at the end of time t as measured on Earth, the time that passes on the space vehicle is only $t\sqrt{1-v^2}$. To establish this result, consider the coordinate systems S and S' and clocks C and C' that we have been studying. Suppose that the origin of S' coincides with the space vehicle and the origin of S coincides with a point in the solar system (stationary relative to the sun) so that the origins of S and S' coincide and clocks C and C' read zero at the moment the astronaut embarks on the trip.

As viewed from S , the space-time coordinates of the vehicle at any time $t > 0$ as measured by C are

$$\begin{pmatrix} vt \\ 0 \\ t \end{pmatrix},$$

whereas, as viewed from S' , the space-time coordinates of the vehicle at any time $t' > 0$ as measured by C' are

$$\begin{pmatrix} 0 \\ 0 \\ t' \end{pmatrix}.$$

But if two sets of space-time coordinates

$$\begin{pmatrix} vt \\ 0 \\ t \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} 0 \\ 0 \\ t' \end{pmatrix}$$

are to describe the same event, it must follow that

$$\mathbf{T}_v \begin{pmatrix} vt \\ 0 \\ t \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ t' \end{pmatrix}.$$

Thus

$$\begin{pmatrix} \frac{1}{\sqrt{1-v^2}} & 0 & \frac{-v}{\sqrt{1-v^2}} \\ 0 & 1 & 0 \\ \frac{-v}{\sqrt{1-v^2}} & 0 & \frac{1}{\sqrt{1-v^2}} \end{pmatrix} \begin{pmatrix} vt \\ 0 \\ t \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ t' \end{pmatrix}.$$

From the preceding equation, we obtain $\frac{-v^2t}{\sqrt{1-v^2}} + \frac{t}{\sqrt{1-v^2}} = t'$, or

$$t' = t\sqrt{1-v^2}. \quad (26)$$

This is the desired result.

A dramatic consequence of time contraction is that distances are contracted along the line of motion (see Exercise 7).

Let us make one additional point. Suppose that we choose units of distance and time commonly used in the study of motion, such as the mile, the kilometer, and the second. Recall that the velocity v we have been using is actually the ratio of the velocity using these units with the speed of light c , using the same units. For this reason, we can replace v in any of the equations given in this section with the ratio v/c , where v and c are given using the same units of measurement. Thus, for example, given a set of units of distance and time, (26) becomes

$$t' = t\sqrt{1 - \frac{v^2}{c^2}}.$$

EXERCISES

1. Complete the proof of Theorem 6.40 for the case $t < 0$.

2. For

$$w_1 = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \quad \text{and} \quad w_2 = \begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix},$$

show that

- (a) $\{w_1, w_2\}$ is an orthogonal basis for $\text{span}(\{e_1, e_3\})$;
(b) $\text{span}(\{e_1, e_3\})$ is $T_v^* L_A T_v$ -invariant.

3. Derive (24), and prove that

$$T_v \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} = \begin{pmatrix} \frac{-v}{\sqrt{1-v^2}} \\ 0 \\ \frac{1}{\sqrt{1-v^2}} \end{pmatrix}. \quad (25)$$

Hint: Use a technique similar to the derivation of (22).

4. Consider three coordinate systems S , S' , and S'' with the corresponding axes (x, x', x'') , (y, y', y'') , and (z, z', z'') parallel and such that the x -, x' -, and x'' -axes coincide. Suppose that S' is moving past S at a velocity $v_1 > 0$ (as measured on S), S'' is moving past S' at a velocity $v_2 > 0$ (as measured on S'), and S'' is moving past S at a velocity $v_3 > 0$ (as measured on S), and that there are three clocks C , C' , and C'' such that C is stationary relative to S , C' is stationary relative to S' , and C'' is stationary relative to S'' . Suppose that when measured on any of the three clocks, all the origins of S , S' , and S'' coincide at time 0. Assuming that $T_{v_3} = T_{v_2} T_{v_1}$ (i.e., $B_{v_3} = B_{v_2} B_{v_1}$), prove that

$$v_3 = \frac{v_1 + v_2}{1 + v_1 v_2}.$$

Note that substituting $v_2 = 1$ in this equation yields $v_3 = 1$. This tells us that the speed of light as measured in S or S' is the same. Why would we be surprised if this were not the case?

5. Compute $(B_v)^{-1}$. Show $(B_v)^{-1} = B_{(-v)}$. Conclude that if S' moves at a negative velocity v relative to S , then $[T_v]_\beta = B_v$, where B_v is of the form given in Theorem 6.43. Visit goo.gl/9gWNYu for a solution.
6. Suppose that an astronaut left Earth in the year 2000 and traveled to a star 99 light years away from Earth at 99% of the speed of light and that upon reaching the star immediately turned around and returned to Earth at the same speed. Assuming Einstein's special theory of relativity, show that if the astronaut was 20 years old at the time of departure, then he or she would return to Earth at age 48.2 in the year 2200. Explain the use of Exercise 4 in solving this problem.
7. Recall the moving space vehicle considered in the study of time contraction. Suppose that the vehicle is moving toward a fixed star located on the x -axis of S at a distance b units from the origin of S . If the space vehicle moves toward the star at velocity v , Earthlings (who remain "almost" stationary relative to S) compute the time it takes for the vehicle to reach the star as $t = b/v$. Due to the phenomenon of time contraction, the astronaut perceives a time span of $t' = t\sqrt{1 - v^2} = (b/v)\sqrt{1 - v^2}$. A paradox appears in that the astronaut perceives a time span inconsistent with a distance of b and a velocity of v . The paradox is resolved by observing that the distance from the solar system to the star as measured by the astronaut is less than b .

Assuming that the coordinate systems S and S' and clocks C and C' are as in the discussion of time contraction, prove the following results.

- (a) At time t (as measured on C), the space-time coordinates of the star relative to S and C are

$$\begin{pmatrix} b \\ 0 \\ t \end{pmatrix}.$$

- (b) At time t (as measured on C), the space-time coordinates of the star relative to S' and C' are

$$\begin{pmatrix} \frac{b-vt}{\sqrt{1-v^2}} \\ 0 \\ \frac{t-bv}{\sqrt{1-v^2}} \end{pmatrix}.$$

- (c) For

$$x' = \frac{b-vt}{\sqrt{1-v^2}} \quad \text{and} \quad t' = \frac{t-bv}{\sqrt{1-v^2}},$$

we have $x' = b\sqrt{1-v^2} - t'v$.

This result may be interpreted to mean that at time t' as measured by the astronaut, the distance from the astronaut to the star, as measured by the astronaut, (see Figure 6.9) is

$$b\sqrt{1-v^2} - t'v.$$

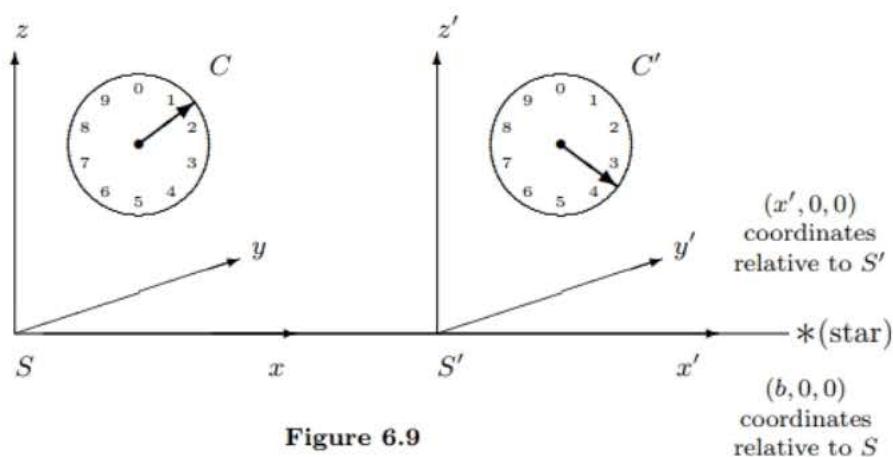


Figure 6.9

- (d) Conclude from the preceding equation that

- (1) the speed of the space vehicle relative to the star, as measured by the astronaut, is v ;

- (2) the distance from Earth to the star, as measured by the astronaut, is $b\sqrt{1 - v^2}$.

Thus distances along the line of motion of the space vehicle appear to be contracted by a factor of $\sqrt{1 - v^2}$.

6.10* CONDITIONING AND THE RAYLEIGH QUOTIENT

In Section 3.4, we studied specific techniques that allow us to solve systems of linear equations in the form $Ax = b$, where A is an $m \times n$ matrix and b is an $m \times 1$ vector. Such systems often arise in applications to the real world. The coefficients in the system are frequently obtained from experimental data, and, in many cases, both m and n are so large that a computer must be used in the calculation of the solution. Thus two types of errors must be considered. First, experimental errors arise in the collection of data since no instruments can provide completely accurate measurements. Second, computers introduce roundoff errors. One might intuitively feel that small relative changes in the coefficients of the system cause small relative errors in the solution. A system that has this property is called **well-conditioned**; otherwise, the system is called **ill-conditioned**.

We now consider several examples of these types of errors, concentrating primarily on changes in b rather than on changes in the entries of A . In addition, we assume that A is a square, complex (or real), invertible matrix since this is the case most frequently encountered in applications.

Example 1

Consider the system

$$\begin{aligned}x_1 + x_2 &= 5 \\x_1 - x_2 &= 1.\end{aligned}$$

The solution to this system is

$$\begin{pmatrix} 3 \\ 2 \end{pmatrix}.$$

Now suppose that we change the system somewhat and consider the new system

$$\begin{aligned}x_1 + x_2 &= 5 \\x_1 - x_2 &= 1.0001.\end{aligned}$$

This modified system has the solution

$$\begin{pmatrix} 3.00005 \\ 1.99995 \end{pmatrix}.$$

We see that a change of 10^{-4} in one coefficient has caused a change of less than 10^{-4} in each coordinate of the new solution. More generally, the system

$$\begin{aligned}x_1 + x_2 &= 5 \\x_1 - x_2 &= 1 + \delta\end{aligned}$$

has the solution

$$\begin{pmatrix} 3 + \delta/2 \\ 2 - \delta/2 \end{pmatrix}.$$

Hence small changes in b introduce small changes in the solution. Of course, we are really interested in *relative changes* since a change in the solution of, say, 10, is considered large if the original solution is of the order 10^{-2} , but small if the original solution is of the order 10^6 .

We use the notation δb to denote the vector $b' - b$, where b is the vector in the original system and b' is the vector in the modified system. Thus we have

$$\delta b = \begin{pmatrix} 5 \\ 1+h \end{pmatrix} - \begin{pmatrix} 5 \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ h \end{pmatrix}.$$

We now define the **relative change** in b to be the scalar $\|\delta b\|/\|b\|$, where $\|\cdot\|$ denotes the standard norm on C^n (or R^n); that is, $\|b\| = \sqrt{\langle b, b \rangle}$. Most of what follows, however, is true for any norm. Similar definitions hold for the **relative change** in x . In this example,

$$\frac{\|\delta b\|}{\|b\|} = \frac{|h|}{\sqrt{26}} \quad \text{and} \quad \frac{\|\delta x\|}{\|x\|} = \frac{\left\| \begin{pmatrix} 3 + h/2 \\ 2 - h/2 \end{pmatrix} - \begin{pmatrix} 3 \\ 2 \end{pmatrix} \right\|}{\left\| \begin{pmatrix} 3 \\ 2 \end{pmatrix} \right\|} = \frac{|h|}{\sqrt{26}}.$$

Thus the relative change in x equals, coincidentally, the relative change in b ; so the system is well-conditioned. ♦

Example 2

Consider the system

$$\begin{array}{rcl} x_1 & + & x_2 = 3 \\ x_1 & + & 1.00001x_2 = 3.00001, \end{array}$$

which has

$$\begin{pmatrix} 2 \\ 1 \end{pmatrix}$$

as its solution. The solution to the related system

$$\begin{array}{rcl} x_1 + x_2 = 3 \\ x_1 + 1.00001x_2 = 3.00001 + \delta \end{array}$$

is

$$\begin{pmatrix} 2 - (10^5)h \\ 1 + (10^5)h \end{pmatrix}.$$

Hence,

$$\frac{\|\delta x\|}{\|x\|} = 10^5 \sqrt{2/5} |h| \geq 10^4 |h|,$$

while

$$\frac{\|\delta b\|}{\|b\|} \approx \frac{|h|}{3\sqrt{2}}.$$

Thus the relative change in x is at least 10^4 times the relative change in b ! This system is very ill-conditioned. Observe that the lines defined by the two equations are nearly coincident. So a small change in either line could greatly alter the point of intersection, that is, the solution to the system. ♦

To apply the full strength of the theory of self-adjoint matrices to the study of conditioning, we need the notion of the norm of a matrix. (See Exercises 26–30 of Section 6.1 for further results about norms.)

Definition. Let A be a complex (or real) $n \times n$ matrix. Define the **norm** of A by

$$\|A\|_E = \max_{x \neq 0} \frac{\|Ax\|}{\|x\|},$$

where $x \in \mathbb{C}^n$ or $x \in \mathbb{R}^n$.

Intuitively, $\|A\|_E$ represents the maximum *magnification* of a vector by the matrix A . The question of whether or not this maximum exists, as well as the problem of how to compute it, can be answered by the use of the so-called *Rayleigh quotient*.

Definition. Let B be an $n \times n$ self-adjoint matrix. The **Rayleigh quotient** for $x \neq 0$ is defined to be the scalar $R(x) = \langle Bx, x \rangle / \|x\|^2$.

The following result characterizes the extreme values of the Rayleigh quotient of a self-adjoint matrix.

Theorem 6.44. For a self-adjoint matrix $B \in M_{n \times n}(F)$, we have that $\max_{x \neq 0} R(x)$ is the largest eigenvalue of B and $\min_{x \neq 0} R(x)$ is the smallest eigenvalue of B .

Proof. By Theorems 6.19 (p. 381) and 6.20 (p. 381), we may choose an orthonormal basis $\{v_1, v_2, \dots, v_n\}$ of eigenvectors of B such that $Bv_i = \lambda_i v_i$ ($1 \leq i \leq n$), where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$. (Recall that by the lemma to Theorem 6.17, p. 370, the eigenvalues of B are real.) Now, for $x \in F^n$, there exist scalars a_1, a_2, \dots, a_n such that

$$x = \sum_{i=1}^n a_i v_i.$$

Hence

$$R(x) = \frac{\langle Bx, x \rangle}{\|x\|^2} = \frac{\left\langle \sum_{i=1}^n a_i \lambda_i v_i, \sum_{j=1}^n a_j v_j \right\rangle}{\|x\|^2}$$

$$= \frac{\sum_{i=1}^n \lambda_i |a_i|^2}{\|x\|^2} \leq \frac{\lambda_1 \sum_{i=1}^n |a_i|^2}{\|x\|^2} = \frac{\lambda_1 \|x\|^2}{\|x\|^2} = \lambda_1.$$

It is easy to see that $R(v_1) = \lambda_1$, so we have demonstrated the first half of the theorem. The second half is proved similarly. ■

Corollary 1. *For any square matrix A , $\|A\|_E$ is finite and, in fact, equals $\sqrt{\lambda}$, where λ is the largest eigenvalue of A^*A .*

Proof. Let B be the self-adjoint matrix A^*A , and let λ be the largest eigenvalue of B . Since, for $x \neq 0$,

$$0 \leq \frac{\|Ax\|^2}{\|x\|^2} = \frac{\langle Ax, Ax \rangle}{\|x\|^2} = \frac{\langle A^*Ax, x \rangle}{\|x\|^2} = \frac{\langle Bx, x \rangle}{\|x\|^2} = R(x),$$

it follows from Theorem 6.44 that $\|A\|_E^2 = \lambda$. ■

Observe that the proof of Corollary 1 shows that all the eigenvalues of A^*A are nonnegative. For our next result, we need the following lemma.

Lemma. *For any square matrix A , λ is an eigenvalue of A^*A if and only if λ is an eigenvalue of AA^* .*

Proof. Let λ be an eigenvalue of A^*A . If $\lambda = 0$, then A^*A is not invertible. Hence A and A^* are not invertible, so that λ is also an eigenvalue of AA^* . The proof of the converse is similar.

Suppose now that $\lambda \neq 0$. Then there exists $x \neq 0$ such that $A^*Ax = \lambda x$. Apply A to both sides to obtain $(AA^*)(Ax) = \lambda(Ax)$. Since $Ax \neq 0$ (lest $\lambda x = 0$), we have that λ is an eigenvalue of AA^* . The proof of the converse is left as an exercise. ■

Corollary 2. *Let A be an invertible matrix. Then $\|A^{-1}\|_E = 1/\sqrt{\lambda}$, where λ is the smallest eigenvalue of A^*A .*

Proof. Recall that λ is an eigenvalue of an invertible matrix if and only if λ^{-1} is an eigenvalue of its inverse.

Now let $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$ be the eigenvalues of A^*A , which by the lemma are the eigenvalues of AA^* . Then $\|A^{-1}\|_E^2$ equals the largest eigenvalue of $(A^{-1})^*A^{-1} = (AA^*)^{-1}$, which equals $1/\lambda_n$. ■

For many applications, it is only the largest and smallest eigenvalues that are of interest. For example, in the case of vibration problems, the smallest eigenvalue represents the lowest frequency at which vibrations can occur.

We see the role of both of these eigenvalues in our study of conditioning.

Example 3

Let

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

Then

$$B = A^*A = \begin{pmatrix} 2 & -1 & 1 \\ -1 & 2 & 1 \\ 1 & 1 & 2 \end{pmatrix}.$$

The eigenvalues of B are 3, 3, and 0. Therefore $\|A\|_E = \sqrt{3}$. For any

$$x = \begin{pmatrix} a \\ b \\ c \end{pmatrix} \neq 0,$$

we may compute $R(x)$ for the matrix B as

$$3 \geq R(x) = \frac{\langle Bx, x \rangle}{\|x\|^2} = \frac{2(a^2 + b^2 + c^2 - ab + ac + bc)}{a^2 + b^2 + c^2}. \quad \blacklozenge$$

Now that we know $\|A\|_E$ exists for every square matrix A , we can make use of the inequality $\|Ax\| \leq \|A\|_E \cdot \|x\|$, which holds for every x .

Assume in what follows that A is invertible, $b \neq 0$, and $Ax = b$. For a given δb , let δx be the vector that satisfies $A(x + \delta x) = b + \delta b$. Then $A(\delta x) = \delta b$, and so $\delta x = A^{-1}(\delta b)$. Hence

$$\|b\| = \|Ax\| \leq \|A\|_E \cdot \|x\| \quad \text{and} \quad \|\delta x\| = \|A^{-1}(\delta b)\| \leq \|A^{-1}\|_E \cdot \|\delta b\|.$$

Thus

$$\frac{\|\delta x\|}{\|x\|} \leq \frac{\|x\|^2}{\|b\|/\|A\|_E} \leq \frac{\|A^{-1}\|_E \cdot \|\delta b\| \cdot \|A\|_E}{\|b\|} = \|A\|_E \cdot \|A^{-1}\|_E \cdot \left(\frac{\|\delta b\|}{\|b\|} \right).$$

Similarly (see Exercise 9),

$$\frac{1}{\|A\|_E \cdot \|A^{-1}\|_E} \left(\frac{\|\delta b\|}{\|b\|} \right) \leq \frac{\|\delta x\|}{\|x\|}.$$

The number $\|A\|_E \cdot \|A^{-1}\|_E$ is called the **condition number** of A and is denoted $\text{cond}(A)$. We summarize these results in the following theorem.

Theorem 6.45. *For the system $Ax = b$, where A is invertible and $b \neq 0$, the following statements are true.*

- (a) We have $\frac{1}{\text{cond}(A)} \frac{\|\delta b\|}{\|b\|} \leq \frac{\|\delta x\|}{\|x\|} \leq \text{cond}(A) \frac{\|\delta b\|}{\|b\|}$.

- (b) If λ_1 and λ_n are the largest and smallest eigenvalues, respectively, of A^*A , then $\text{cond}(A) = \sqrt{\lambda_1/\lambda_n}$.

Proof. Statement (a) follows from the previous inequalities, and (b) follows from Corollaries 1 and 2 to Theorem 6.44. ■

It should be noted that the definition of $\text{cond}(A)$ depends on how the norm of A is defined. There are many reasonable ways of defining the norm of a matrix. In fact, the only property needed to establish Theorem 6.45(a) and the two displayed inequalities preceding it is that $\|Ax\| \leq \|A\|_E \cdot \|x\|$ for all x .

It is clear from Theorem 6.45(a) that $\text{cond}(A) \geq 1$. It is left as an exercise to prove that $\text{cond}(A) = 1$ if and only if A is a scalar multiple of a unitary or orthogonal matrix. Moreover, it can be shown with some work that equality can be obtained in (a) by an appropriate choice of b and δb .

We can see immediately from (a) that if $\text{cond}(A)$ is close to 1, then a small relative error in b forces a small relative error in x . If $\text{cond}(A)$ is large, however, then the relative error in x may be small even though the relative error in b is large, or the relative error in x may be large even though the relative error in b is small! In short, $\text{cond}(A)$ merely indicates the *potential* for large relative errors.

We have so far considered only errors in the vector b . If there is an error δA in the coefficient matrix of the system $Ax = b$, the situation is more complicated. For example, $A + \delta A$ may fail to be invertible. But under the appropriate assumptions, it can be shown that a bound for the relative error in x can be given in terms of $\text{cond}(A)$. For example, Charles Cullen (Charles G. Cullen, *An Introduction to Numerical Linear Algebra*, PWS Publishing Co., Boston 1994, p. 60) shows that if $A + \delta A$ is invertible, then

$$\frac{\|\delta x\|}{\|x + \delta x\|} \leq \text{cond}(A) \frac{\|\delta A\|_E}{\|A\|_E}.$$

It should be mentioned that, in practice, one never computes $\text{cond}(A)$ from its definition, for it would be an unnecessary waste of time to compute A^{-1} merely to determine its norm. In fact, if a computer is used to find A^{-1} , the computed inverse of A in all likelihood only approximates A^{-1} , and the error in the computed inverse is affected by the size of $\text{cond}(A)$. So we are caught in a vicious circle! There are, however, some situations in which a usable approximation of $\text{cond}(A)$ can be found. Thus, in most cases, the estimate of the relative error in x is based on an estimate of $\text{cond}(A)$.

EXERCISES

1. Label the following statements as true or false.

- (a) If $Ax = b$ is well-conditioned, then $\text{cond}(A)$ is small.
 (b) If $\text{cond}(A)$ is large, then $Ax = b$ is ill-conditioned.
 (c) If $\text{cond}(A)$ is small, then $Ax = b$ is well-conditioned.
 (d) The norm of A equals the Rayleigh quotient.
 (e) The norm of A always equals the largest eigenvalue of A .
2. Compute the norms of the following matrices.
- (a) $\begin{pmatrix} 4 & 0 \\ 1 & 3 \end{pmatrix}$ (b) $\begin{pmatrix} 5 & 3 \\ -3 & 3 \end{pmatrix}$ (c) $\begin{pmatrix} 1 & \frac{-2}{\sqrt{3}} & 0 \\ 0 & \frac{-2}{\sqrt{3}} & 1 \\ 0 & \frac{2}{\sqrt{3}} & 1 \end{pmatrix}$
3. Prove that if B is symmetric, then $\|B\|_E$ is the largest eigenvalue of B .
4. Let A and A^{-1} be as follows:

$$A = \begin{pmatrix} 6 & 13 & -17 \\ 13 & 29 & -38 \\ -17 & -38 & 50 \end{pmatrix} \quad \text{and} \quad A^{-1} = \begin{pmatrix} 6 & -4 & 1 \\ -4 & 11 & 7 \\ -1 & 7 & 5 \end{pmatrix}.$$

The eigenvalues of A are approximately 84.74, 0.2007, and 0.0588.

- (a) Approximate $\|A\|_E$, $\|A^{-1}\|_E$, and $\text{cond}(A)$. (Note Exercise 3.)
 (b) Suppose that we have vectors x and \tilde{x} such that $Ax = b$ and $\|b - A\tilde{x}\| \leq 0.001$. Use (a) to determine upper bounds for $\|\tilde{x} - A^{-1}b\|$ (the absolute error) and $\|\tilde{x} - A^{-1}b\|/\|A^{-1}b\|$ (the relative error).
5. Suppose that x is the actual solution of $Ax = b$ and that a computer arrives at an approximate solution \tilde{x} . If $\text{cond}(A) = 100$, $\|b\| = 1$, and $\|b - A\tilde{x}\| = 0.1$, obtain upper and lower bounds for $\|x - \tilde{x}\|/\|x\|$.
6. Let

$$B = \begin{pmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{pmatrix}.$$

Compute

$$R \begin{pmatrix} 1 \\ -2 \\ 3 \end{pmatrix}, \quad \|B\|_E, \quad \text{and} \quad \text{cond}(B).$$

7. Let B be a symmetric matrix. Prove that $\min_{x \neq 0} R(x)$ equals the smallest eigenvalue of B .
8. Prove that if λ is an eigenvalue of AA^* , then λ is an eigenvalue of A^*A . This completes the proof of the lemma to Corollary 2 to Theorem 6.44.

9. Prove that if A is an invertible matrix and $Ax = b$, then

$$\frac{1}{\|A\|_E \cdot \|A^{-1}\|_E} \left(\frac{\|\delta b\|}{\|b\|} \right) \leq \frac{\|\delta x\|}{\|x\|}.$$

10. Prove the left inequality of (a) in Theorem 6.45.
11. Prove that $\text{cond}(A) = 1$ if and only if A is a scalar multiple of a unitary or orthogonal matrix.
12. (a) Let A and B be square matrices that are unitarily equivalent. Prove that $\|A\|_E = \|B\|_E$.
- (b) Let T be a linear operator on a finite-dimensional inner product space V . Define

$$\|T\|_E = \max_{x \neq 0} \frac{\|T(x)\|}{\|x\|}.$$

Prove that $\|T\|_E = \|[\mathbf{T}]_\beta\|_E$, where β is any orthonormal basis for V .

- (c) Let V be an infinite-dimensional inner product space with an orthonormal basis $\{v_1, v_2, \dots\}$. Let T be the linear operator on V such that $T(v_k) = kv_k$. Prove that $\|T\|_E$ (defined in (b)) does not exist.

Visit goo.gl/B8Uw33 for a solution.

The next exercise assumes the definitions of *singular value* and *pseudoinverse* and the results of Section 6.7.

13. Let A be an $n \times n$ matrix of rank r with the nonzero singular values $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r$. Prove each of the following results.
- (a) $\|A\|_E = \sigma_1$.
- (b) $\|A^\dagger\|_E = \frac{1}{\sigma_r}$.
- (c) If A is invertible (and hence $r = n$), then $\text{cond}(A) = \frac{\sigma_1}{\sigma_n}$.

6.11* THE GEOMETRY OF ORTHOGONAL OPERATORS

By Theorem 6.22 (p. 383), any rigid motion on a finite-dimensional real inner product space is the composite of an orthogonal operator and a translation. Thus, to understand the geometry of rigid motions thoroughly, we must analyze the structure of orthogonal operators. In this section, we show that any orthogonal operator on a finite-dimensional real inner product space can be described in terms of rotations and reflections.

This material assumes familiarity with the results about direct sums developed at the end of Section 5.2 and with the definition of the determinant of a linear operator given in Section 5.1 as well as elementary properties of the determinant in Exercise 8 of Section 5.1.

We now extend our earlier definitions of *rotation* and *reflection* on \mathbb{R}^2 to all 2-dimensional real inner product spaces.

Definitions. Let T be a linear operator on a two-dimensional real inner product space V .

We call T a **rotation** if there exists an orthonormal basis $\beta = \{x_1, x_2\}$ for V and a real number θ such that

$$T(x_1) = (\cos \theta)x_1 + (\sin \theta)x_2 \quad \text{and} \quad T(x_2) = (-\sin \theta)x_1 + (\cos \theta)x_2.$$

We call T a **reflection** if there exists a one-dimensional subspace W of V such that $T(x) = -x$ for all $x \in W$ and $T(y) = y$ for all $y \in W^\perp$. In this context, T is called a **reflection of V about W^\perp** .

As a convenience, we define *rotations* and *reflections* on a 1-dimensional inner product space.

Definitions. A linear operator T on a 1-dimensional inner product space V is called a **rotation** if T is the identity and a **reflection** if $T(x) = -x$ for all $x \in V$.

Trivially, rotations and reflections on 1-dimensional inner product spaces are orthogonal operators. It should be noted that rotations and reflections on 2-dimensional real inner product spaces (or composites of these) are orthogonal operators (see Exercise 2).

Example 1

Some Typical Reflections

(a) Define $T: \mathbb{R}^2 \rightarrow \mathbb{R}^2$ by $T(a, b) = (-a, b)$, and let $W = \text{span}(\{e_1\})$. Then $T(x) = -x$ for all $x \in W$, and $T(y) = y$ for all $y \in W^\perp$. Thus T is a reflection of \mathbb{R}^2 about $W^\perp = \text{span}(\{e_2\})$, the y -axis.

(b) Define $T: \mathbb{R}^2 \rightarrow \mathbb{R}^2$ by $T(a, b) = (b, a)$, and let $W = \text{span}(\{(1, 1)\})$. Clearly $T(w) = w$ for all $w \in W$. Let $(a, b) \in W^\perp$. Then (a, b) is orthogonal to $(1, 1)$ and hence, $a + b = 0$. So $b = -a$. Thus $W^\perp = \text{span}(\{(1, -1)\})$. It follows that $T(a, b) = (a, -a) = -(-a, a) = -(b, a)$. Hence T is the reflection of \mathbb{R}^2 about W . ♦

The next theorem characterizes all orthogonal operators on a two-dimensional real inner product space V . The proof follows from Theorem 6.23 (p. 384) since all two-dimensional real inner product spaces are structurally identical. For a rigorous justification, apply Theorem 2.21 (p. 105), where β is an orthonormal basis for V . By Exercise 15 of Section 6.2, the resulting isomorphism $\phi_\beta: V \rightarrow \mathbb{R}^2$ preserves inner products. (See Exercise 8.)

Theorem 6.46. Let T be an orthogonal operator on a two-dimensional real inner product space V . Then T is either a rotation or a reflection. Furthermore, T is a rotation if and only if $\det(T) = 1$, and T is a reflection if and only if $\det(T) = -1$.

A complete description of the reflections of \mathbb{R}^2 is given in Section 6.5.

Corollary. Let V be a two-dimensional real inner product space.

- (a) The composite of a reflection and a rotation on V is a reflection on V .
- (b) The composite of two reflections on V is a rotation on V .
- (c) The product of two rotations on V is a rotation on V .

Proof. If T_1 is a reflection on V and T_2 is a rotation on V , then by Theorem 6.46, $\det(T_1) = 1$ and $\det(T_2) = -1$. Let $T = T_2T_1$ be the composite. Since T_2 and T_1 are orthogonal, so is T . Moreover, $\det(T) = \det(T_2) \cdot \det(T_1) = -1$. Thus, by Theorem 6.46, T is a reflection. The proof for T_1T_2 is similar.

The proofs of (b) and (c) are similar to that of (a). ■

We now study orthogonal operators on spaces of higher dimension.

Lemma. If T is a linear operator on a nonzero finite-dimensional real vector space V , then there exists a T -invariant subspace W of V such that $1 \leq \dim(W) \leq 2$.

Proof. Fix an ordered basis $\beta = \{y_1, y_2, \dots, y_n\}$ for V , and let $A = [T]_\beta$. Let $\phi_\beta: V \rightarrow \mathbb{R}^n$ be the linear transformation defined by $\phi_\beta(y_i) = e_i$ for $i = 1, 2, \dots, n$. Then ϕ_β is an isomorphism, and, as we have seen in Section 2.4, the diagram in Figure 6.10 commutes, that is, $L_A \phi_\beta = \phi_\beta T$. As a consequence, it suffices to show that there exists an L_A -invariant subspace Z of \mathbb{R}^n such that $1 \leq \dim(Z) \leq 2$. If we then define $W = \phi_\beta^{-1}(Z)$, it follows that W satisfies the conclusions of the lemma (see Exercise 12).

$$\begin{array}{ccc} V & \xrightarrow{T} & V \\ \downarrow \phi_\beta & & \downarrow \phi_\beta \\ \mathbb{R}^n & \xrightarrow{L_A} & \mathbb{R}^n \end{array}$$

Figure 6.10

The matrix A can be considered as an $n \times n$ matrix over C and, as such, can be used to define a linear operator U on C^n by $U(v) = Av$. Since U is a linear operator on a finite-dimensional vector space over C , it has an

eigenvalue $\lambda \in C$. Let $x \in C^n$ be an eigenvector corresponding to λ . We may write $\lambda = \lambda_1 + i\lambda_2$, where λ_1 and λ_2 are real, and

$$x = \begin{pmatrix} a_1 + ib_1 \\ a_2 + ib_2 \\ \vdots \\ a_n + ib_n \end{pmatrix},$$

where the a_i 's and b_i 's are real. Thus, setting

$$x_1 = \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{pmatrix} \quad \text{and} \quad x_2 = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix},$$

we have $x = x_1 + ix_2$, where x_1 and x_2 have real entries. Note that at least one of x_1 or x_2 is nonzero since $x \neq 0$. Hence

$$U(x) = \lambda x = (\lambda_1 + i\lambda_2)(x_1 + ix_2) = (\lambda_1 x_1 - \lambda_2 x_2) + i(\lambda_1 x_2 + \lambda_2 x_1).$$

Similarly,

$$U(x) = A(x_1 + ix_2) = Ax_1 + iAx_2.$$

Comparing the real and imaginary parts of these two expressions for $U(x)$, we conclude that

$$Ax_1 = \lambda_1 x_1 - \lambda_2 x_2 \quad \text{and} \quad Ax_2 = \lambda_1 x_2 + \lambda_2 x_1.$$

Finally, let $Z = \text{span}(\{x_1, x_2\})$, the span being taken as a subspace of R^n . Since $x_1 \neq 0$ or $x_2 \neq 0$, Z is a nonzero subspace. Thus $1 \leq \dim(Z) \leq 2$, and the preceding pair of equations shows that Z is L_A -invariant. ■

Theorem 6.47. *Let T be an orthogonal operator on a nonzero finite-dimensional real inner product space V . Then there exists a collection of pairwise orthogonal T -invariant subspaces $\{W_1, W_2, \dots, W_m\}$ of V such that*

- (a) $1 \leq \dim(W_i) \leq 2$ for $i = 1, 2, \dots, m$.
- (b) $V = W_1 \oplus W_2 \oplus \dots \oplus W_m$.

Proof. The proof is by mathematical induction on $\dim(V)$. If $\dim(V) = 1$, the result is obvious. So assume that the result is true whenever $\dim(V) < n$ for some fixed integer $n > 1$.

Suppose $\dim(V) = n$. By the lemma, there exists a T -invariant subspace W_1 of V such that $1 \leq \dim(W) \leq 2$. If $W_1 = V$, the result is established. Otherwise, $W_1^\perp \neq \{0\}$. By Exercise 13, W_1^\perp is T -invariant and the restriction

of T to W_1^\perp is orthogonal. Since $\dim(W_1^\perp) < n$, we may apply the induction hypothesis to $T_{W_1^\perp}$ and conclude that there exists a collection of pairwise orthogonal T -invariant subspaces $\{W_1, W_2, \dots, W_m\}$ of W_1^\perp such that $1 \leq \dim(W_i) \leq 2$ for $i = 2, 3, \dots, m$ and $W_1^\perp = W_2 \oplus W_3 \oplus \dots \oplus W_m$. Thus $\{W_1, W_2, \dots, W_m\}$ is pairwise orthogonal, and by Exercise 13(d) of Section 6.2,

$$V = W_1 \oplus W_1^\perp = W_1 \oplus W_2 \oplus \dots \oplus W_m. \quad \blacksquare$$

Applying Theorem 6.46 in the context of Theorem 6.47, we conclude that the restriction of T to W_i is either a rotation or a reflection for each $i = 1, 2, \dots, m$. Thus, in some sense, T is composed of rotations and reflections. Unfortunately, very little can be said about the uniqueness of the decomposition of V in Theorem 6.47. For example, the W_i 's, the number m of W_i 's, and the number of W_i 's for which T_{W_i} is a reflection are not unique. Although the number of W_i 's for which T_{W_i} is a reflection is not unique, whether this number is even or odd is an intrinsic property of T . Moreover, we can always decompose V so that T_{W_i} is a reflection for at most one W_i . These facts are established in the following result.

Theorem 6.48. Let T, V, W_1, \dots, W_m be as in Theorem 6.47.

- (a) The number of W_i 's for which T_{W_i} is a reflection is even or odd according to whether $\det(T) = 1$ or $\det(T) = -1$.
- (b) It is always possible to decompose V as in Theorem 6.47 so that the number of W_i 's for which T_{W_i} is a reflection is zero or one according to whether $\det(T) = 1$ or $\det(T) = -1$. Furthermore, if T_{W_i} is a reflection, then $\dim(W_i) = 1$.

Proof. (a) Let r denote the number of W_i 's in the decomposition for which T_{W_i} is a reflection. Then, by Exercise 14,

$$\det(T) = \det(T_{W_1}) \cdot \det(T_{W_2}) \cdot \dots \cdot \det(T_{W_m}) = (-1)^r,$$

proving (a).

(b) Let $E = \{x \in V : T(x) = -x\}$; then E is a T -invariant subspace of V . Let $W = E^\perp$; then W is T -invariant. So by applying Theorem 6.47 to T_W , we obtain a collection of pairwise orthogonal T -invariant subspaces $\{W_1, W_2, \dots, W_k\}$ of W such that $W = W_1 \oplus W_2 \oplus \dots \oplus W_k$, and for $1 \leq i \leq k$, the dimension of each W_i is either 1 or 2. Observe that, for each $i = 1, 2, \dots, k$, T_{W_i} is a rotation. For otherwise, if T_{W_i} is a reflection, there exists a nonzero $x \in W_i$ for which $T(x) = -x$. But then, $x \in W_i \cap E \subseteq E^\perp \cap E = \{0\}$, a contradiction. If $E = \{0\}$, the result follows. Otherwise, choose an orthonormal basis β for E containing p vectors ($p > 0$). It is possible to decompose β into a pairwise disjoint union $\beta = \beta_1 \cup \beta_2 \cup \dots \cup \beta_r$ such that each β_i contains exactly two vectors for $i < r$, and β_r contains

two vectors if p is even and one vector if p is odd. For each $i = 1, 2, \dots, r$, let $W_{k+i} = \text{span}(\beta_i)$. Then, clearly, $\{W_1, W_2, \dots, W_k, \dots, W_{k+r}\}$ is pairwise orthogonal, and

$$V = W_1 \oplus W_2 \oplus \cdots \oplus W_k \oplus \cdots \oplus W_{k+r}. \quad (27)$$

Moreover, if any β_i contains two vectors, then

$$\det(T_{W_{k+i}}) = \det([T_{W_{k+i}}]_{\beta_i}) = \det \begin{pmatrix} -1 & 0 \\ 0 & -1 \end{pmatrix} = 1.$$

So $T_{W_{k+i}}$ is a rotation, and hence T_{W_j} is a rotation for $j < k + r$. If β_r consists of one vector, then $\dim(W_{k+r}) = 1$ and

$$\det(T_{W_{k+r}}) = \det([T_{W_{k+r}}]_{\beta_r}) = \det(-1) = -1.$$

Thus $T_{W_{k+r}}$ is a reflection by Theorem 6.47, and we conclude that the decomposition in (27) satisfies (b). ■

Example 2

Orthogonal Operators on a Three-Dimensional Real Inner Product Space

Let T be an orthogonal operator on a three-dimensional real inner product space V . Then, by Theorem 6.48(b), V can be decomposed into a direct sum of T -invariant orthogonal subspaces so that the restriction of T to each is either a rotation or a reflection, with at most one reflection. Let

$$V = W_1 \oplus W_2 \oplus \cdots \oplus W_m$$

be such a decomposition. Clearly, $m = 2$ or $m = 3$.

If $m = 2$, then $V = W_1 \oplus W_2$. Without loss of generality, suppose that $\dim(W_1) = 1$ and $\dim(W_2) = 2$. Thus T_{W_1} is a reflection or the identity on W_1 , and T_{W_2} is a rotation.

If $m = 3$, then $V = W_1 \oplus W_2 \oplus W_3$ and $\dim(W_i) = 1$ for all i . If T_{W_i} is not a reflection, then it is the identity on W_i . If no T_{W_i} is a reflection, then T is the identity operator. ♦

EXERCISES

- Label the following statements as true or false. Assume that the underlying vector spaces are one or two-dimensional real inner product spaces.
 - Any orthogonal operator is either a rotation or a reflection.
 - The composite of any two rotations is a rotation.
 - The identity operator is a rotation.

- (d) The composite of two reflections is a reflection.
 (e) Any orthogonal operator is a composite of rotations.
 (f) For any orthogonal operator T , if $\det(T) = -1$, then T is a reflection.
 (g) Reflections always have eigenvalues.
 (h) Rotations always have eigenvalues.
 (i) If T is an operator on a 2-dimensional space V and W is a subspace of dimension 1 such that T is a reflection of V about W^\perp , then W is the eigenspace of T corresponding to the eigenvalue $\lambda = -1$.
 (j) The composite of an orthogonal operator and a translation is an orthogonal operator.
2. Prove that rotations, reflections, and composites of rotations and reflections are orthogonal operators.
3. Let
- $$A = \begin{pmatrix} \frac{1}{2} & \frac{\sqrt{3}}{2} \\ \frac{\sqrt{3}}{2} & -\frac{1}{2} \end{pmatrix} \quad \text{and} \quad B = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}.$$
- (a) Prove that L_A is a reflection.
 (b) Find the subspace of \mathbb{R}^2 on which L_A acts as the identity.
 (c) Prove that L_{AB} and L_{BA} are rotations.
4. For any real number ϕ , let
- $$A = \begin{pmatrix} \cos \phi & \sin \phi \\ \sin \phi & -\cos \phi \end{pmatrix}.$$
- (a) Prove that L_A is a reflection.
 (b) Find the axis in \mathbb{R}^2 about which L_A reflects.
5. For any real number ϕ , define $T_\phi = L_A$, where
- $$A = \begin{pmatrix} \cos \phi & -\sin \phi \\ \sin \phi & \cos \phi \end{pmatrix}.$$
- (a) Prove that any rotation on \mathbb{R}^2 is of the form T_ϕ for some ϕ .
 (b) Prove that $T_\phi T_\psi = T_{(\phi+\psi)}$ for any $\phi, \psi \in \mathbb{R}$.
 (c) Deduce that any two rotations on \mathbb{R}^2 commute.
6. Prove that if T is a rotation on a 2-dimensional inner product space, then $-T$ is also a rotation.
7. Prove that if T is a reflection on a 2-dimensional inner product space, then T^2 is the identity operator.

8. Prove Theorem 6.46 using the hints preceding the statement of the theorem.
9. Prove that no orthogonal operator can be both a rotation and a reflection.
10. Prove that if V is a two-dimensional real inner product space, then the composite of two reflections on V is a rotation of V .
11. Let V be a one- or a two-dimensional real inner product space. Define $T: V \rightarrow V$ by $T(x) = -x$. Prove that T is a rotation if and only if $\dim(V) = 2$.
12. Complete the proof of the lemma to Theorem 6.47 by showing that $W = \phi_{\beta}^{-1}(Z)$ satisfies the required conditions.
13. Let T be an orthogonal [unitary] operator on a finite-dimensional real [complex] inner product space V . If W is a T -invariant subspace of V , prove the following results.
 - (a) T_W is an orthogonal [unitary] operator on W .
 - (b) W^\perp is a T -invariant subspace of V . Hint: Use the fact that T_W is one-to-one and onto to conclude that, for any $y \in W$, $T^*(y) = T^{-1}(y) \in W$.
 - (c) T_{W^\perp} is an orthogonal [unitary] operator on W .
14. Let T be a linear operator on a finite-dimensional vector space V , where V is a direct sum of T -invariant subspaces, say, $V = W_1 \oplus W_2 \oplus \dots \oplus W_k$. Prove that $\det(T) = \det(T_{W_1}) \cdot \det(T_{W_2}) \cdot \dots \cdot \det(T_{W_k})$.
15. Complete the proof of the corollary to Theorem 6.48.
16. Let V be a real inner product space of dimension 2. For any $x, y \in V$ such that $x \neq y$ and $\|x\| = \|y\| = 1$, show that there exists a unique rotation T on V such that $T(x) = y$. Visit goo.gl/ahQT67 for a solution.

INDEX OF DEFINITIONS FOR CHAPTER 6

Adjoint of a linear operator	355	Critical point	436
Adjoint of a matrix	329	Diagonalizable bilinear form	425
Bilinear form	419	Fourier coefficients of a vector relative to an orthonormal set	345
Complex inner product space	330	Frobenius inner product	330
Condition number	463	Gram–Schmidt orthogonalization process	342
Congruent matrices	423	Hessian matrix	436
Conjugate transpose (adjoint) of a matrix	329		

- Index of a bilinear form 441
Index of a matrix 442
Inner product 327
Inner product space 330
Invariants of a bilinear form 441
Invariants of a matrix 442
Least squares line 358
Legendre polynomials 344
Local extremum 436
Local maximum 436
Local minimum 436
Matrix representation of a bilinear form 421
Minimal solution of a system of equations 361
Norm of a matrix 461
Norm of a vector 331
Normal matrix 367
Normal operator 367
Normalizing a vector 333
Orthogonal complement of a subset of an inner product space 347
Orthogonally equivalent matrices 381
Orthogonal matrix 379
Orthogonal operator 376
Orthogonal projection 395
Orthogonal projection on a subspace 348
Orthogonal subset of an inner product space 333
Orthogonal vectors 333
Orthonormal basis 339
Orthonormal set 333
Penrose conditions 418
Polar decomposition of a matrix 409
Pseudoinverse of a linear transformation 410
Pseudoinverse of a matrix 411
Quadratic form 430
Rank of a bilinear form 439
Rayleigh quotient 461
Real inner product space 330
Reflection 467
Resolution of the identity operator induced by a linear transformation 399
Rigid motion 383
Rotation 467
Self-adjoint matrix 370
Self-adjoint operator 370
Signature of a form 441
Signature of a matrix 442
Singular value decomposition of a matrix 407
Singular value of a linear transformation 404
Singular value of a matrix 406
Space-time coordinates 449
Spectral decomposition of a linear operator 399
Spectrum of a linear operator 399
Standard inner product 328
Symmetric bilinear form 425
Translation 383
Trigonometric polynomial 396
Unitarily equivalent matrices 381
Unitary matrix 379
Unitary operator 376
Unit vector 333

Canonical Forms

-
- 7.1 The Jordan Canonical Form I
 - 7.2 The Jordan Canonical Form II
 - 7.3 The Minimal Polynomial
 - 7.4* The Rational Canonical Form

As we learned in Chapter 5, the advantage of a diagonalizable linear operator lies in the simplicity of its description. Such an operator has a diagonal matrix representation, or, equivalently, there is an ordered basis for the underlying vector space consisting of eigenvectors of the operator. However, not every linear operator is diagonalizable, even if its characteristic polynomial splits. Example 3 of Section 5.2 describes such an operator.

It is the purpose of this chapter to consider alternative matrix representations for nondiagonalizable operators. These representations are called *canonical forms*. There are different kinds of canonical forms, and their advantages and disadvantages depend on how they are applied. Every canonical form of an operator is obtained by an appropriate choice of an ordered basis. Naturally, the canonical forms of a linear operator are not diagonal matrices if the linear operator is not diagonalizable.

In this chapter, we treat two common canonical forms. The first of these, the *Jordan canonical form*, requires that the characteristic polynomial of the operator splits. This form is always available if the underlying field is algebraically closed, that is, if every polynomial with coefficients from the field splits. For example, the field of complex numbers is algebraically closed by the fundamental theorem of algebra (see Appendix D). The first two sections deal with this form. The *rational canonical form*, treated in Section 7.4, does not require such a factorization.

7.1 THE JORDAN CANONICAL FORM I

Let T be a linear operator on a finite-dimensional vector space V , and suppose that the characteristic polynomial of T splits. Recall from Section 5.2 that the diagonalizability of T depends on whether the union of ordered bases for the distinct eigenspaces of T is an ordered basis for V . So a lack of diagonalizability means that at least one eigenspace of T is too “small.”

In this section, we extend the definition of eigenspace to *generalized eigenspace*. From these subspaces, we select ordered bases whose union is an ordered basis β for V such that

$$[T]_{\beta} = \begin{pmatrix} A_1 & O & \cdots & O \\ O & A_2 & \cdots & O \\ \vdots & \vdots & & \vdots \\ O & O & \cdots & A_k \end{pmatrix},$$

where each O is a zero matrix, and each A_i is a square matrix of the form (λ) or

$$\begin{pmatrix} \lambda & 1 & 0 & \cdots & 0 & 0 \\ 0 & \lambda & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & \lambda & 1 \\ 0 & 0 & 0 & \cdots & 0 & \lambda \end{pmatrix}$$

for some eigenvalue λ of T . Such a matrix A_i is called a **Jordan block** corresponding to λ , and the matrix $[T]_{\beta}$ is called a **Jordan canonical form** of T . We also say that the ordered basis β is a **Jordan canonical basis** for T . Observe that each Jordan block A_i is “almost” a diagonal matrix—in fact, $[T]_{\beta}$ is a diagonal matrix if and only if each A_i is of the form (λ) .

Example 1

Suppose that T is a linear operator on \mathbb{C}^8 , and $\beta = \{v_1, v_2, \dots, v_8\}$ is an ordered basis for \mathbb{C}^8 such that

$$J = [T]_{\beta} = \left(\begin{array}{ccc|ccccc} 2 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 2 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 3 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{array} \right)$$

is a Jordan canonical form of T . Notice that the characteristic polynomial of T is $\det(J - tI) = (t - 2)^4(t - 3)^2t^2$, and hence the multiplicity of each eigenvalue is the number of times that the eigenvalue appears on the diagonal of J . Also observe that v_1, v_4, v_5 , and v_7 are the only vectors in β that are eigenvectors of T . These are the vectors corresponding to the columns of J with no 1 above the diagonal entry. ♦

In Sections 7.1 and 7.2, we prove that every linear operator whose characteristic polynomial splits has a Jordan canonical form that is unique up to

the order of its Jordan blocks. Nevertheless, it is not the case that the Jordan canonical form is completely determined by the characteristic polynomial of the operator. For example, let T' be the linear operator on \mathbb{C}^8 such that $[T']_\beta = J'$, where β is the ordered basis in Example 1 and

$$J' = \begin{pmatrix} 2 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 3 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}.$$

Then the characteristic polynomial of T' is also $(t - 2)^4(t - 3)^2t^2$. But the operator T' has the Jordan canonical form J' , which is different from J , the Jordan canonical form of the linear operator T of Example 1.

Consider again the matrix J and the ordered basis β of Example 1. Notice that $T(v_2) = v_1 + 2v_2$, and so $(T - 2I)(v_2) = v_1$. Similarly, $(T - 2I)(v_3) = v_2$. Since v_1 and v_4 are eigenvectors of T corresponding to $\lambda = 2$, it follows that $(T - 2I)^3(v_i) = 0$ for $i = 1, 2, 3$, and 4. Similarly $(T - 3I)^2(v_i) = 0$ for $i = 5, 6$, and $(T - 0I)^2(v_i) = 0$ for $i = 7, 8$.

Because of the structure of each Jordan block in a Jordan canonical form, we can generalize these observations: *If v lies in a Jordan canonical basis for a linear operator T and is associated with a $k \times k$ Jordan block with diagonal entry λ , then $(T - \lambda I)^k(v) = 0$.* Eigenvectors satisfy this condition for $k = 1$.

Definition. Let T be a linear operator on a vector space V , and let λ be a scalar. A nonzero vector x in V is called a **generalized eigenvector of T corresponding to λ** if $(T - \lambda I)^p(x) = 0$ for some positive integer p .

Notice that if x is a generalized eigenvector of T corresponding to λ , and p is the smallest positive integer for which $(T - \lambda I)^p(x) = 0$, then $(T - \lambda I)^{p-1}(x)$ is an eigenvector of T corresponding to λ . Therefore λ is an eigenvalue of T .

In the context of Example 1, each vector in β is a generalized eigenvector of T . In fact, v_1, v_2, v_3 , and v_4 correspond to the eigenvalue 2, v_5 and v_6 correspond to the eigenvalue 3, and v_7 and v_8 correspond to the eigenvalue 0.

Just as eigenvectors lie in eigenspaces, generalized eigenvectors lie in “generalized eigenspaces.”

Definition. Let T be a linear operator on a vector space V , and let λ be an eigenvalue of T . The **generalized eigenspace of T corresponding to λ** , denoted K_λ , is the subset of V defined by

$$K_\lambda = \{x \in V : (T - \lambda I)^p(x) = 0 \text{ for some positive integer } p\}.$$

Note that K_λ consists of the zero vector and all generalized eigenvectors corresponding to λ .

Recall that a subspace W of V is T -invariant for a linear operator T if $T(W) \subseteq W$. In the development that follows, we assume the results of Exercises 3 and 4 of Section 5.4. In particular, for any polynomial $g(x)$, if W is T -invariant, then it is also $g(T)$ -invariant. Furthermore, the range of a linear operator T is T -invariant.

Theorem 7.1. *Let T be a linear operator on a vector space V , and let λ be an eigenvalue of T . Then*

- (a) K_λ is a T -invariant subspace of V containing E_λ (the eigenspace of T corresponding to λ).
- (b) For any eigenvalue μ of T such that $\mu \neq \lambda$,

$$K_\lambda \cap K_\mu = \{0\}.$$

- (c) For any scalar $\mu \neq \lambda$, the restriction of $T - \mu I$ to K_λ is one-to-one and onto.

Proof. (a) Clearly, $0 \in K_\lambda$. Suppose that x and y are in K_λ . Then there exist positive integers p and q such that

$$(T - \lambda I)^p(x) = (T - \lambda I)^q(y) = 0.$$

Therefore

$$\begin{aligned} (T - \lambda I)^{p+q}(x + y) &= (T - \lambda I)^{p+q}(x) + (T - \lambda I)^{p+q}(y) \\ &= (T - \lambda I)^q(0) + (T - \lambda I)^p(0) \\ &= 0, \end{aligned}$$

and hence $x + y \in K_\lambda$. The proof that K_λ is closed under scalar multiplication is straightforward.

To show that K_λ is T -invariant, consider any $x \in K_\lambda$. Choose a positive integer p such that $(T - \lambda I)^p(x) = 0$. Then

$$(T - \lambda I)^p T(x) = T(T - \lambda I)^p(x) = T(0) = 0.$$

Therefore $T(x) \in K_\lambda$.

Finally, it is a simple observation that E_λ is contained in K_λ .

(b) Suppose that $w \in K_\lambda \cap K_\mu$. Then there exist positive integers p and q such that

$$(T - \lambda I)^p(w) = (T - \mu I)^q(w) = 0.$$

As polynomials in t , $(t - \lambda)^p$ and $(t - \mu)^q$ are relatively prime, and hence by Theorem E.2 in Appendix E there exist polynomials $q_1(t)$ and $q_2(t)$ such that

$$q_1(t)(t - \lambda)^p + q_2(t)(t - \mu)^q = 1.$$

It follows that

$$q_1(T)(T - \lambda I)^p(w) + q_2(T)(T - \mu I)^q(w) = w,$$

and hence $\theta = w$.

(c) Let μ be any scalar such that $\mu \neq \lambda$. Since K_λ is T -invariant, it is also $(T - \mu I)$ -invariant. If $(T - \mu I)(w) = \theta$ for some $w \in K_\lambda$, then $w \in E_\mu \subset K_\mu$. Hence $w = \theta$ by (b), and we conclude that the restriction of $T - \mu I$ to K_λ is one-to-one.

We now show that the restriction of $T - \mu I$ to K_λ is onto. Let $x \in K_\lambda$, p be the smallest positive integer for which $(T - \lambda I)^p(x) = \theta$, and

$$W = \text{span}\{x, (T - \lambda I)(x), \dots, (T - \lambda I)^{p-1}(x)\}.$$

It is easily shown that W is a T -invariant subspace of K_λ , and hence it is also $(T - \mu I)$ -invariant. So $T - \mu I$ maps W to W . Since $T - \mu I$ is one-to-one on K_λ , it must be one-to-one on W . Thus, because W is finite-dimensional, $T - \mu I$ maps W onto W . So there exists a $y \in W$ such that $(T - \mu I)(y) = x$. Therefore $T - \mu I$ maps K_λ onto K_λ . ■

Theorem 7.2. *Let T be a linear operator on a finite-dimensional vector space V such that the characteristic polynomial of T splits. Suppose that λ is an eigenvalue of T with multiplicity m . Then*

- (a) $\dim(K_\lambda) \leq m$.
- (b) $K_\lambda = N((T - \lambda I)^m)$.

Proof. (a) Let $W = K_\lambda$, and let $h(t)$ be the characteristic polynomial of T_W . By Theorem 5.20 (p. 312), $h(t)$ divides the characteristic polynomial of T , and by Theorem 7.1(b), λ is the only eigenvalue of T_W . Hence $h(t) = (-1)^d(t - \lambda)^d$, where $d = \dim(W)$, and $d \leq m$.

(b) Clearly $N((T - \lambda I)^m) \subseteq K_\lambda$. The characteristic polynomial of T is of the form $f(t) = (t - \lambda)^m g(t)$, where $g(t)$ is a product of powers of the form $(t - \mu)$ for eigenvalues $\mu \neq \lambda$. By Theorem 7.1, $g(T)$ is one-to-one on K_λ . Since K_λ is finite-dimensional, it is also onto. Let $x \in K_\lambda$. Then there exists $y \in K_\lambda$ such that $g(T)(y) = x$. Hence

$$(T - \lambda I)^m(x) = (T - \lambda I)^m g(T)(y) = f(T)(y) = 0$$

by the Cayley-Hamilton Theorem. Therefore $x \in N((T - \lambda I)^m)$, and thus $K_\lambda \subseteq N((T - \lambda I)^m)$. ■

Theorem 7.3. *Let T be a linear operator on a finite-dimensional vector space V such that the characteristic polynomial of T splits, and let $\lambda_1, \lambda_2, \dots, \lambda_k$ be the distinct eigenvalues of T . Then, for every $x \in V$, there exist unique vectors $v_i \in K_{\lambda_i}$, for $i = 1, 2, \dots, k$, such that*

$$x = v_1 + v_2 + \cdots + v_k.$$

Proof. Suppose the characteristic polynomial of \mathbf{T} is

$$f(t) = (t - \lambda_1)^{n_1}(t - \lambda_2)^{n_2} \dots (t - \lambda_k)^{n_k}.$$

For $j = 1, 2, \dots, k$, let

$$f_j(t) = \prod_{\substack{i \neq j \\ i=1}}^k (t - \lambda_i)^{n_i}.$$

First, we prove the existence of the vectors v_i . Because the polynomials $f_j(t)$ are relatively prime for $j = 1, 2, \dots, k$, by Theorem E.2 there exist polynomials $q_j(t)$ for $j = 1, 2, \dots, k$ such that

$$q_1(t)f_1(t) + q_2(t)f_2(t) + \dots + q_k(t)f_k(t) = 1.$$

It follows that

$$v = q_1(\mathbf{T})f_1(\mathbf{T})(v) + q_2(\mathbf{T})f_2(\mathbf{T})(v) + \dots + q_k(\mathbf{T})f_k(\mathbf{T})(v).$$

For $i = 1, 2, \dots, k$, let $v_i = q_i(\mathbf{T})f_i(\mathbf{T})(v) = f_i(\mathbf{T})q_i(\mathbf{T})(v)$. Then for each i ,

$$(\mathbf{T} - \lambda_i \mathbf{I})^{n_i}(v_i) = f(\mathbf{T})(q_i(\mathbf{T})(v)) = 0$$

by the Cayley-Hamilton theorem. It follows that $v_i \in \mathbf{K}_{\lambda_i}$ for $i = 1, 2, \dots, k$.

Now we prove the uniqueness of the vectors v_i . Since each \mathbf{K}_{λ_j} is \mathbf{T} -invariant, each \mathbf{K}_{λ_j} is invariant under the operator $f_j(\mathbf{T})$ for $j = 1, 2, \dots, k$. Furthermore, $f_j(\mathbf{T})$ is one-to-one on \mathbf{K}_{λ_j} by Theorem 7.1(c) and maps \mathbf{K}_{λ_i} to $\{0\}$ for $i \neq j$ by Theorem 7.2(b).

Suppose that

$$v = \sum_{i=1}^k v_i = \sum_{i=1}^k w_i,$$

where $v_i, w_i \in \mathbf{K}_{\lambda_i}$ for $i = 1, 2, \dots, k$. Then for $j = 1, 2, \dots, k$,

$$f_j(\mathbf{T})(v) = \sum_{i=1}^k f_j(\mathbf{T})(v_i) = \sum_{i=1}^k f_j(\mathbf{T})(w_i),$$

and hence $f_j(\mathbf{T})(v_j) = f_j(\mathbf{T})(w_j)$ since $f_j(\mathbf{T})(v_i) = f_j(\mathbf{T})(w_i) = 0$ for $i \neq j$. Because $f_j(\mathbf{T})$ is one-to-one on \mathbf{K}_{λ_j} , it follows that $v_j = w_j$. ■

The next result extends Theorem 5.8(b) (p. 267) to all linear operators whose characteristic polynomials split. In this case, the eigenspaces are replaced by generalized eigenspaces.

Theorem 7.4. Let \mathbf{T} be a linear operator on a finite-dimensional vector space \mathbf{V} whose characteristic polynomial $(t - \lambda_1)^{m_1}(t - \lambda_2)^{m_2} \dots (t - \lambda_k)^{m_k}$ splits. For $i = 1, 2, \dots, k$, let β_i be an ordered basis for \mathbf{K}_{λ_i} . Then

- (a) $\beta_i \cap \beta_j = \emptyset$ for $i \neq j$.
- (b) $\beta = \beta_1 \cup \beta_2 \cup \dots \cup \beta_k$ is an ordered basis for V .
- (c) $\dim(K_{\lambda_i}) = m_i$ for all i .

Proof. (a) This is a direct consequence of Theorem 7.1(b).

(b) First, we prove that β is linearly independent. Consider a linear combination of vectors in β equal to θ . For each i , let v_i be the sum of the terms of this linear combination involving the vectors in β_i . Then

$$v_1 + v_2 + \dots + v_k = \theta = \theta + \theta + \dots + \theta,$$

and hence $v_i = \theta$ for all i by Theorem 7.3. Since each β_i is linearly independent, it follows that all of the coefficients in the linear combination are zeros. Hence β is linearly independent.

We now prove that β is a spanning set for V . Consider any vector $v \in V$. By Theorem 7.3, for $i = 1, 2, \dots, k$ there exist vectors $v_i \in K_{\lambda_i}$ such that $v = v_1 + v_2 + \dots + v_k$. Since each v_i is a linear combination of the vectors in β_i , v is a linear combination of the vectors in β . Hence β spans V .

(c) Since $\dim(V)$ is equal to the degree of the characteristic polynomial of T , it follows that $\dim(V) = \sum_{i=1}^k m_i$. By (b), $\dim(V) = \sum_{i=1}^k \dim(K_{\lambda_i})$. Hence $\sum_{i=1}^m (m_i - \dim(K_{\lambda_i})) = 0$. By Theorem 7.2(a), $m_i - \dim(K_{\lambda_i}) \geq 0$ for all i . Consequently, $m_i - \dim(K_{\lambda_i}) = 0$ for all i , which is the desired result. ■

Corollary. Let T be a linear operator on a finite-dimensional vector space V such that the characteristic polynomial of T splits. Then T is diagonalizable if and only if $E_\lambda = K_\lambda$ for every eigenvalue λ of T .

Proof. Combining Theorems 7.4 and 5.8(a) (p. 267), we see that T is diagonalizable if and only if $\dim(E_\lambda) = \dim(K_\lambda)$ for each eigenvalue λ of T . But $E_\lambda \subseteq K_\lambda$, and hence these subspaces have the same dimension if and only if they are equal. ■

We now focus our attention on the problem of selecting suitable bases for the generalized eigenspaces of a linear operator so that we may use Theorem 7.4 to obtain a Jordan canonical basis for the operator. For this purpose, we consider again the basis β of Example 1. We have seen that the first four vectors of β lie in the generalized eigenspace K_2 . Observe that the vectors in β that determine the first Jordan block of J are of the form

$$\{v_1, v_2, v_3\} = \{(T - 2I)^2(v_3), (T - 2I)(v_3), v_3\}.$$

Furthermore, observe that $(T - 2I)^3(v_3) = \theta$. The relation between these vectors is the key to finding Jordan canonical bases. This leads to the following definitions.

Definitions. Let T be a linear operator on a vector space V , and let x be a generalized eigenvector of T corresponding to the eigenvalue λ . Suppose that p is the smallest positive integer for which $(T - \lambda I)^p(x) = 0$. Then the ordered set

$$\{(T - \lambda I)^{p-1}(x), (T - \lambda I)^{p-2}(x), \dots, (T - \lambda I)(x), x\}$$

is called a **cycle of generalized eigenvectors** of T corresponding to λ . The vectors $(T - \lambda I)^{p-1}(x)$ and x are called the **initial vector** and the **end vector** of the cycle, respectively. We say that the **length** of the cycle is p .

Notice that the initial vector of a cycle of generalized eigenvectors of a linear operator T is the only eigenvector of T in the cycle. Also observe that if x is an eigenvector of T corresponding to the eigenvalue λ , then the set $\{x\}$ is a cycle of generalized eigenvectors of T corresponding to λ of length 1.

In Example 1, the subsets $\beta_1 = \{v_1, v_2, v_3\}$, $\beta_2 = \{v_4\}$, $\beta_3 = \{v_5, v_6\}$, and $\beta_4 = \{v_7, v_8\}$ are the cycles of generalized eigenvectors of T that occur in β . Notice that β is a disjoint union of these cycles. Furthermore, setting $W_i = \text{span}(\beta_i)$ for $i = 1, 2, 3, 4$, we see that β_i is a basis for W_i and $[T_{W_i}]_{\beta_i}$ is the i th Jordan block of the Jordan canonical form of T . This is precisely the condition that is required for a Jordan canonical basis.

Theorem 7.5. Let T be a linear operator on a finite-dimensional vector space V whose characteristic polynomial splits, and suppose that β is a basis for V such that β is a disjoint union of cycles of generalized eigenvectors of T . Then the following statements are true.

- (a) For each cycle γ of generalized eigenvectors contained in β , $W = \text{span}(\gamma)$ is T -invariant, and $[T_W]_\gamma$ is a Jordan block.
- (b) β is a Jordan canonical basis for V .

Proof. (a) Suppose that γ corresponds to the eigenvalue λ , γ has length p , and x is the end vector of γ . Then $\gamma = \{v_1, v_2, \dots, v_p\}$, where

$$v_i = (T - \lambda I)^{p-i}(x) \text{ for } i < p \quad \text{and} \quad v_p = x.$$

So

$$(T - \lambda I)(v_1) = (T - \lambda I)^p(x) = 0,$$

and hence $T(v_1) = \lambda v_1$. For $i > 1$,

$$(T - \lambda I)(v_i) = (T - \lambda I)^{p-(i-1)}(x) = v_{i-1}.$$

This shows that $T - \lambda I$ maps W into itself, and so T does also. Thus, by the preceding equations, we see that $[T_W]_\gamma$ is a Jordan block.

For (b), simply repeat the arguments of (a) for each cycle in β in order to obtain $[T]_\beta$. We leave the details as an exercise. ■

Theorem 7.7. Let T be a linear operator on a finite-dimensional vector space V , and let λ be an eigenvalue of T . Then K_λ has an ordered basis consisting of a union of disjoint cycles of generalized eigenvectors corresponding to λ .

Proof. The proof is by mathematical induction on $n = \dim(K_\lambda)$. The result is clear for $n = 1$. So suppose that for some integer $n > 1$ the result is true whenever $\dim(K_\lambda) < n$, and assume that $\dim(K_\lambda) = n$. Let U denote the restriction of $T - \lambda I$ to K_λ . Then $R(U)$ is a subspace of K_λ of lesser dimension, and $R(U)$ is the space of generalized eigenvectors corresponding to λ for the restriction of T to $R(U)$. Therefore, by the induction hypothesis, there exist disjoint cycles $\gamma_1, \gamma_2, \dots, \gamma_q$ of generalized eigenvectors of this restriction, and hence of T itself, corresponding to λ for which $\gamma = \bigcup_{i=1}^q \gamma_i$ is a basis for $R(U)$. For $i = 1, 2, \dots, q$, the end vector of γ_i is the image under U of a vector $v_i \in K_\lambda$, and so we can extend each γ_i to a larger cycle $\tilde{\gamma}_i = \gamma_i \cup \{v_i\}$ of generalized eigenvectors of T corresponding to λ . For $i = 1, 2, \dots, q$, let w_i be the initial vector of $\tilde{\gamma}_i$ (and hence of γ_i). Since $\{w_1, w_2, \dots, w_q\}$ is a linearly independent subset of E_λ , this set can be extended to a basis $\{w_1, w_2, \dots, w_q, u_1, u_2, \dots, u_s\}$ for E_λ . Then $\tilde{\gamma}_1, \tilde{\gamma}_2, \dots, \tilde{\gamma}_q, \{u_1\}, \{u_2\}, \dots, \{u_s\}$ are disjoint cycles of generalized eigenvectors of T corresponding to λ such that the initial vectors of these cycles are linearly independent. Therefore their union $\tilde{\gamma}$ is a linearly independent subset of K_λ by Theorem 7.6.

Finally, we show that $\tilde{\gamma}$ is a basis for K_λ . Suppose that γ consists of $r = \text{rank}(U)$ vectors. Then $\tilde{\gamma}$ consists of $r + q + s$ vectors. Furthermore, since $\{w_1, w_2, \dots, w_q, u_1, u_2, \dots, u_s\}$ is a basis for $E_\lambda = N(U)$, it follows that $\text{nullity}(U) = q + s$. Therefore

$$\dim(K_\lambda) = \text{rank}(U) + \text{nullity}(U) = r + q + s.$$

So $\tilde{\gamma}$ is a linearly independent subset of K_λ containing $\dim(K_\lambda)$ vectors, and thus $\tilde{\gamma}$ is a basis for K_λ . ■

Corollary 1. Let T be a linear operator on a finite-dimensional vector space V whose characteristic polynomial splits. Then T has a Jordan canonical form.

Proof. Let $\lambda_1, \lambda_2, \dots, \lambda_k$ be the distinct eigenvalues of T . By Theorem 7.7, for each i there is an ordered basis β_i consisting of a disjoint union of cycles of generalized eigenvectors corresponding to λ_i . Let $\beta = \beta_1 \cup \beta_2 \cup \dots \cup \beta_k$. Then, by Theorem 7.5(b), β is a Jordan canonical basis for V . ■

The Jordan canonical form also can be studied from the viewpoint of matrices.

Definition. Let $A \in M_{n \times n}(F)$ be such that the characteristic polynomial of A (and hence of L_A) splits. Then the **Jordan canonical form** of A is defined to be the Jordan canonical form of the linear operator L_A on F^n .

The next result is an immediate consequence of this definition and Corollary 1.

Corollary 2. Let A be an $n \times n$ matrix whose characteristic polynomial splits. Then A has a Jordan canonical form J , and A is similar to J . ■

Proof. Exercise.

We can now compute the Jordan canonical forms of matrices and linear operators in some simple cases, as is illustrated in the next two examples. The tools necessary for computing the Jordan canonical forms in general are developed in the next section.

Example 2

Let

$$A = \begin{pmatrix} 3 & 1 & -2 \\ -1 & 0 & 5 \\ -1 & -1 & 4 \end{pmatrix} \in M_{3 \times 3}(R).$$

To find the Jordan canonical form for A , we need to find a Jordan canonical basis for $T = L_A$.

The characteristic polynomial of A is

$$f(t) = \det(A - tI) = -(t - 3)(t - 2)^2.$$

Hence $\lambda_1 = 3$ and $\lambda_2 = 2$ are the eigenvalues of A with multiplicities 1 and 2, respectively. By Theorem 7.4, $\dim(K_{\lambda_1}) = 1$, and $\dim(K_{\lambda_2}) = 2$. By Theorem 7.2, $K_{\lambda_1} = N(T - 3I)$, and $K_{\lambda_2} = N((T - 2I)^2)$. Since $E_{\lambda_1} = N(T - 3I)$, we have that $E_{\lambda_1} = K_{\lambda_1}$. Observe that $(-1, 2, 1)$ is an eigenvector of T corresponding to $\lambda_1 = 3$; therefore

$$\beta_1 = \left\{ \begin{pmatrix} -1 \\ 2 \\ 1 \end{pmatrix} \right\}$$

is a basis for K_{λ_1} .

Since $\dim(K_{\lambda_2}) = 2$ and a generalized eigenspace has a basis consisting of a union of cycles, this basis is either a union of two cycles of length 1 or a single cycle of length 2. Since the rank of $A - 2I$ is 2 it follows that $\dim(E_{\lambda_2}) = 1$. Hence the former is not the case. It can easily be shown that $w = (1, -3, -1)$ is an eigenvector of A corresponding to λ_2 . This vector can serve as the initial

vector of the cycle. Any solution to the equation $(A - 2I)v = w$ will serve as the end vector. For example, we can take $v = (1, -2, 0)$. Thus

$$\beta_2 = \{(A - 2I)v, v\} = \left\{ \begin{pmatrix} 1 \\ -3 \\ -1 \end{pmatrix}, \begin{pmatrix} -1 \\ 2 \\ 0 \end{pmatrix} \right\}$$

is a cycle of generalized eigenvectors for $\lambda_2 = 2$. Finally, we take the union of these two bases to obtain

$$\beta = \beta_1 \cup \beta_2 = \left\{ \begin{pmatrix} -1 \\ 2 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ -3 \\ -1 \end{pmatrix}, \begin{pmatrix} -1 \\ 2 \\ 0 \end{pmatrix} \right\},$$

which is a Jordan canonical basis for A . Therefore,

$$J = [\mathbf{T}]_{\beta} = \left(\begin{array}{c|cc} 3 & 0 & 0 \\ \hline 0 & 2 & 1 \\ 0 & 0 & 2 \end{array} \right)$$

is a Jordan canonical form for A . Notice that A is similar to J . In fact, $J = Q^{-1}AQ$, where Q is the matrix whose columns are the vectors in β .

◆

Example 3

Let \mathbf{T} be the linear operator on $\mathbb{P}_2(R)$ defined by $\mathbf{T}(g(x)) = -g(x) - g'(x)$. We find a Jordan canonical form of \mathbf{T} and a Jordan canonical basis for \mathbf{T} .

Let β be the standard ordered basis for $\mathbb{P}_2(R)$. Then

$$[\mathbf{T}]_{\beta} = \begin{pmatrix} -1 & -1 & 0 \\ 0 & -1 & -2 \\ 0 & 0 & -1 \end{pmatrix},$$

which has the characteristic polynomial $f(t) = -(t + 1)^3$. Thus $\lambda = -1$ is the only eigenvalue of \mathbf{T} , and hence $K_{\lambda} = \mathbb{P}_2(R)$ by Theorem 7.4. So β is a basis for K_{λ} . Now

$$\dim(E_{\lambda}) = 3 - \text{rank}(A + I) = 3 - \text{rank} \begin{pmatrix} 0 & -1 & 0 \\ 0 & 0 & -2 \\ 0 & 0 & 0 \end{pmatrix} = 3 - 2 = 1.$$

Therefore a basis for K_{λ} cannot be a union of two or three cycles because the initial vector of each cycle is an eigenvector, and there do not exist two or more linearly independent eigenvectors. So the desired basis must consist

of a single cycle of length 3. If γ is such a cycle, then γ determines a single Jordan block

$$[T]_{\gamma} = \begin{pmatrix} -1 & 1 & 0 \\ 0 & -1 & 1 \\ 0 & 0 & -1 \end{pmatrix},$$

which is a Jordan canonical form of T .

The end vector $h(x)$ of such a cycle must satisfy $(T + I)^2(h(x)) \neq 0$. In any basis for K_{λ} , there must be a vector that satisfies this condition, or else no vector in K_{λ} satisfies this condition, contrary to our reasoning. Testing the vectors in β , we see that $h(x) = x^2$ is acceptable. Therefore

$$\gamma = \{(T + I)^2(x^2), (T + I)(x^2), x^2\} = \{2, -2x, x^2\}$$

is a Jordan canonical basis for T . \blacklozenge

In the next section, we develop a computational approach for finding a Jordan canonical form and a Jordan canonical basis. In the process, we prove that Jordan canonical forms are unique up to the order of the Jordan blocks.

Let T be a linear operator on a finite-dimensional vector space V , and suppose that the characteristic polynomial of T splits. By Theorem 5.10 (p. 277), T is diagonalizable if and only if V is the direct sum of the eigenspaces of T . If T is diagonalizable, then the eigenspaces and the generalized eigenspaces coincide.

For those familiar with the material on direct sums in Section 5.2, we conclude by stating a generalization of Theorem 5.10 for nondiagonalizable operators.

Theorem 7.8. *Let T be a linear operator on a finite-dimensional vector space V whose characteristic polynomial splits. Then V is the direct sum of the generalized eigenspaces of T .*

Proof. Exercise. ■

EXERCISES

1. Label the following statements as true or false.
 - Eigenvectors of a linear operator T are also generalized eigenvectors of T .
 - It is possible for a generalized eigenvector of a linear operator T to correspond to a scalar that is not an eigenvalue of T .
 - Any linear operator on a finite-dimensional vector space has a Jordan canonical form.
 - A cycle of generalized eigenvectors is linearly independent.

- (e) There is exactly one cycle of generalized eigenvectors corresponding to each eigenvalue of a linear operator on a finite-dimensional vector space.
- (f) Let T be a linear operator on a finite-dimensional vector space whose characteristic polynomial splits, and let $\lambda_1, \lambda_2, \dots, \lambda_k$ be the distinct eigenvalues of T . If, for each i , β_i is a basis for K_{λ_i} , then $\beta_1 \cup \beta_2 \cup \dots \cup \beta_k$ is a Jordan canonical basis for T .
- (g) For any Jordan block J , the operator L_J has Jordan canonical form J .
- (h) Let T be a linear operator on an n -dimensional vector space whose characteristic polynomial splits. Then, for any eigenvalue λ of T , $K_\lambda = N((T - \lambda I)^n)$.
2. For each matrix A , find a basis for each generalized eigenspace of L_A consisting of a union of disjoint cycles of generalized eigenvectors. Then find a Jordan canonical form J of A .
- (a) $A = \begin{pmatrix} 1 & 1 \\ -1 & 3 \end{pmatrix}$
- (b) $A = \begin{pmatrix} 1 & 2 \\ 3 & 2 \end{pmatrix}$
- (c) $A = \begin{pmatrix} 11 & -4 & -5 \\ 21 & -8 & -11 \\ 3 & -1 & 0 \end{pmatrix}$
- (d) $A = \begin{pmatrix} 2 & 1 & 0 & 0 \\ 0 & 2 & 1 & 0 \\ 0 & 0 & 3 & 0 \\ 0 & 1 & -1 & 3 \end{pmatrix}$
3. For each linear operator T , find a basis for each generalized eigenspace of T consisting of a union of disjoint cycles of generalized eigenvectors. Then find a Jordan canonical form J of T .
- (a) Define T on $P_2(R)$ by $T(f(x)) = 2f(x) - f'(x)$.
- (b) V is the real vector space of functions spanned by the set of real-valued functions $\{1, t, t^2, e^t, te^t\}$, and T is the linear operator on V defined by $T(f) = f'$.
- (c) T is the linear operator on $M_{2 \times 2}(R)$ defined for all $A \in M_{2 \times 2}(R)$ by $T(A) = BA$, where $B = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}$.
- (d) $T(A) = 2A + A^t$ for all $A \in M_{2 \times 2}(R)$.
- 4.[†] Let T be a linear operator on a vector space V , and let γ be a cycle of generalized eigenvectors that corresponds to the eigenvalue λ . Prove that $\text{span}(\gamma)$ is a T -invariant subspace of V . Visit goo.gl/Lw4ahY for a solution.
5. Let $\gamma_1, \gamma_2, \dots, \gamma_p$ be cycles of generalized eigenvectors of a linear operator T corresponding to an eigenvalue λ . Prove that if the initial eigenvectors are distinct, then the cycles are disjoint.

6. Let $T: V \rightarrow W$ be a linear transformation. Prove the following results.
- $N(T) = N(-T)$.
 - $N(T^k) = N((-T)^k)$.
 - If $V = W$ (so that T is a linear operator on V) and λ is an eigenvalue of T , then for any positive integer k

$$N((T - \lambda I_V)^k) = N((\lambda I_V - T)^k).$$

7. Let U be a linear operator on a finite-dimensional vector space V . Prove the following results.
- $N(U) \subseteq N(U^2) \subseteq \cdots \subseteq N(U^k) \subseteq N(U^{k+1}) \subseteq \cdots$.
 - If $\text{rank}(U^m) = \text{rank}(U^{m+1})$ for some positive integer m , then $\text{rank}(U^m) = \text{rank}(U^k)$ for any positive integer $k \geq m$.
 - If $\text{rank}(U^m) = \text{rank}(U^{m+1})$ for some positive integer m , then $N(U^m) = N(U^k)$ for any positive integer $k \geq m$.
 - Let T be a linear operator on V , and let λ be an eigenvalue of T . Prove that if $\text{rank}((T - \lambda I)^m) = \text{rank}((T - \lambda I)^{m+1})$ for some integer m , then $K_\lambda = N((T - \lambda I)^m)$.
 - Second Test for Diagonalizability.* Let T be a linear operator on V whose characteristic polynomial splits, and let $\lambda_1, \lambda_2, \dots, \lambda_k$ be the distinct eigenvalues of T . Then T is diagonalizable if and only if $\text{rank}(T - \lambda_i I) = \text{rank}((T - \lambda_i I)^2)$ for $1 \leq i \leq k$.
 - Use (e) to obtain a simpler proof of Exercise 24 of Section 5.4: If T is a diagonalizable linear operator on a finite-dimensional vector space V and W is a T -invariant subspace of V , then T_W is diagonalizable.
8. Let T be a linear operator on a finite-dimensional vector space V whose characteristic polynomial splits and has Jordan canonical form J . Prove that for any nonzero scalar c , cJ is a Jordan canonical form for cT .
9. Let T be a linear operator on a finite-dimensional vector space V whose characteristic polynomial splits.
- Prove Theorem 7.5(b).
 - Suppose that β is a Jordan canonical basis for T , and let λ be an eigenvalue of T . Let $\beta' = \beta \cap K_\lambda$. Prove that β' is a basis for K_λ .
10. Let T be a linear operator on a finite-dimensional vector space whose characteristic polynomial splits, and let λ be an eigenvalue of T .
- Suppose that γ is a basis for K_λ consisting of the union of q disjoint cycles of generalized eigenvectors. Prove that $q \leq \dim(E_\lambda)$.
 - Let β be a Jordan canonical basis for T , and suppose that $J = [T]_\beta$ has q Jordan blocks with λ in the diagonal positions. Prove that $q \leq \dim(E_\lambda)$.

11. Prove Corollary 2 to Theorem 7.7.
12. Let T be a linear operator on a finite-dimensional vector space V , and let λ be an eigenvalue of T with corresponding eigenspace and generalized eigenspace E_λ and K_λ , respectively. Let U be an invertible linear operator on V that commutes with T (i.e., $TU = UT$). Prove that $U(E_\lambda) = E_\lambda$ and $U(K_\lambda) = K_\lambda$.

Exercises 13 and 14 are concerned with direct sums of matrices, defined in Section 5.4 on page 318.

13. Prove Theorem 7.8.
14. Let T be a linear operator on a finite-dimensional vector space V such that the characteristic polynomial of T splits, and let $\lambda_1, \lambda_2, \dots, \lambda_k$ be the distinct eigenvalues of T . For each i , let J_i be the Jordan canonical form of the restriction of T to K_{λ_i} . Prove that

$$J = J_1 \oplus J_2 \oplus \cdots \oplus J_k$$

is the Jordan canonical form of J .

7.2 THE JORDAN CANONICAL FORM II

For the purposes of this section, we fix a linear operator T on an n -dimensional vector space V such that the characteristic polynomial of T splits. Let $\lambda_1, \lambda_2, \dots, \lambda_k$ be the distinct eigenvalues of T .

By Theorem 7.7 (p. 484), each generalized eigenspace K_{λ_i} contains an ordered basis β_i consisting of a union of disjoint cycles of generalized eigenvectors corresponding to λ_i . So by Theorems 7.4(b) (p. 480) and 7.5 (p. 482),

the union $\beta = \bigcup_{i=1}^k \beta_i$ is a Jordan canonical basis for T . For each i , let T_i be the restriction of T to K_{λ_i} , and let $A_i = [T_i]_{\beta_i}$. Then A_i is the Jordan canonical form of T_i , and

$$J = [T]_\beta = \begin{pmatrix} A_1 & O & \cdots & O \\ O & A_2 & \cdots & O \\ \vdots & \vdots & & \vdots \\ O & O & \cdots & A_k \end{pmatrix}$$

is the Jordan canonical form of T . In this matrix, each O is a zero matrix of appropriate size.

In this section, we compute the matrices A_i and the bases β_i , thereby computing J and β as well. While developing a method for finding J , it becomes evident that in some sense the matrices A_i are unique.

To aid in formulating the uniqueness theorem for J , we adopt the following convention: The basis β_i for K_{λ_i} will henceforth be ordered in such a way that the cycles appear in order of decreasing length. That is, if β_i is a disjoint union of cycles $\gamma_1, \gamma_2, \dots, \gamma_{n_i}$ and if the length of the cycle γ_j is p_j , we index the cycles so that $p_1 \geq p_2 \geq \dots \geq p_{n_i}$. This ordering of the cycles limits the possible orderings of vectors in β_i , which in turn determines the matrix A_i . It is in this sense that A_i is unique. It then follows that the Jordan canonical form for T is unique up to an ordering of the eigenvalues of T . As we will see, there is no uniqueness theorem for the bases β_i or for β . However, we show that for each i , the number n_i of cycles that form β_i , and the length p_j ($j = 1, 2, \dots, n_i$) of each cycle, is completely determined by T .

Example 1

To illustrate the discussion above, suppose that, for some i , the ordered basis β_i for K_{λ_i} is the union of four cycles $\beta_i = \gamma_1 \cup \gamma_2 \cup \gamma_3 \cup \gamma_4$ with respective lengths $p_1 = 3, p_2 = 3, p_3 = 2$, and $p_4 = 1$. Then

$$A_i = \left(\begin{array}{ccc|cccccc} \lambda_i & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & \lambda_i & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & \lambda_i & 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & \lambda_i & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & \lambda_i & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & \lambda_i & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & \lambda_i & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \lambda_i \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{array} \right) . \quad \blacklozenge$$

To help us visualize each of the matrices A_i and ordered bases β_i , we use an array of dots called a **dot diagram** of T_i , where T_i is the restriction of T to K_{λ_i} . Suppose that β_i is a disjoint union of cycles of generalized eigenvectors $\gamma_1, \gamma_2, \dots, \gamma_{n_i}$ with lengths $p_1 \geq p_2 \geq \dots \geq p_{n_i}$, respectively. The dot diagram of T_i contains one dot for each vector in β_i , and the dots are configured according to the following rules.

1. The array consists of n_i columns (one column for each cycle).
2. Counting from left to right, the j th column consists of the p_j dots that correspond to the vectors of γ_j starting with the initial vector at the top and continuing down to the end vector.

Denote the end vectors of the cycles by v_1, v_2, \dots, v_{n_i} . In the following dot diagram of T_i , each dot is labeled with the name of the vector in β_i to

which it corresponds.

$$\begin{array}{llll}
 \bullet (\mathbf{T} - \lambda_i \mathbf{l})^{p_1-1}(v_1) & \bullet (\mathbf{T} - \lambda_i \mathbf{l})^{p_2-1}(v_2) & \cdots & \bullet (\mathbf{T} - \lambda_i \mathbf{l})^{p_{n_i}-1}(v_{n_i}) \\
 \bullet (\mathbf{T} - \lambda_i \mathbf{l})^{p_1-2}(v_1) & \bullet (\mathbf{T} - \lambda_i \mathbf{l})^{p_2-2}(v_2) & \cdots & \bullet (\mathbf{T} - \lambda_i \mathbf{l})^{p_{n_i}-2}(v_{n_i}) \\
 & & \vdots & \\
 & & \bullet (\mathbf{T} - \lambda_i \mathbf{l})(v_{n_i}) & \\
 & \vdots & & \bullet v_{n_i} \\
 & \bullet (\mathbf{T} - \lambda_i \mathbf{l})(v_2) & & \\
 & \bullet v_2 & & \\
 & \bullet (\mathbf{T} - \lambda_i \mathbf{l})(v_1) & & \\
 & \bullet v_1 & &
 \end{array}$$

Notice that the dot diagram of \mathbf{T}_i has n_i columns (one for each cycle) and p_1 rows. Since $p_1 \geq p_2 \geq \cdots \geq p_{n_i}$, the columns of the dot diagram become shorter (or at least not longer) as we move from left to right.

Now let r_j denote the number of dots in the j th row of the dot diagram. Observe that $r_1 \geq r_2 \geq \cdots \geq r_{p_1}$. Furthermore, the diagram can be reconstructed from the values of the r_i 's. The proofs of these facts, which are combinatorial in nature, are treated in Exercise 9.

In Example 1, with $n_i = 4$, $p_1 = p_2 = 3$, $p_3 = 2$, and $p_4 = 1$, the dot diagram of \mathbf{T}_i is as follows:

$$\begin{array}{cccc}
 \bullet & \bullet & \bullet & \bullet \\
 \bullet & \bullet & \bullet & \\
 \bullet & \bullet & & \\
 \bullet & \bullet & &
 \end{array}$$

Here $r_1 = 4$, $r_2 = 3$, and $r_3 = 2$.

We now devise a method for computing the dot diagram of \mathbf{T}_i using the ranks of linear operators determined by \mathbf{T} and λ_i . It will follow that the dot diagram is completely determined by \mathbf{T} , from which it will follow that it is unique. On the other hand, β_i is not unique. For example, see Exercise 8. (It is for this reason that we associate the dot diagram with \mathbf{T}_i rather than with β_i .)

To determine the dot diagram of \mathbf{T}_i , we devise a method for computing each r_j , the number of dots in the j th row of the dot diagram, using only \mathbf{T} and λ_i . The next three results give us the required method. To facilitate our arguments, we fix a basis β_i for \mathbf{K}_{λ_i} so that β_i is a disjoint union of n_i cycles of generalized eigenvectors with lengths $p_1 \geq p_2 \geq \cdots \geq p_{n_i}$.

Theorem 7.9. *For any positive integer r , the vectors in β_i that are associated with the dots in the first r rows of the dot diagram of \mathbf{T}_i constitute a basis for $\mathbf{N}((\mathbf{T} - \lambda_i \mathbf{l})^r)$. Hence the number of dots in the first r rows of the dot diagram equals $\text{nullity}((\mathbf{T} - \lambda_i \mathbf{l})^r)$.*

Proof. Clearly, $N((T - \lambda_i I)^r) \subseteq K_{\lambda_i}$, and K_{λ_i} is invariant under $(T - \lambda_i I)^r$. Let U denote the restriction of $(T - \lambda_i I)^r$ to K_{λ_i} . By the preceding remarks, $N((T - \lambda_i I)^r) = N(U)$, and hence it suffices to establish the theorem for U . Now define

$$S_1 = \{x \in \beta_i : U(x) = 0\} \quad \text{and} \quad S_2 = \{x \in \beta_i : U(x) \neq 0\}.$$

Let a and b denote the number of vectors in S_1 and S_2 , respectively, and let $m_i = \dim(K_{\lambda_i})$. Then $a + b = m_i$. For any $x \in \beta_i$, $x \in S_1$ if and only if x is one of the first r vectors of a cycle, and this is true if and only if x corresponds to a dot in the first r rows of the dot diagram. Hence a is the number of dots in the first r rows of the dot diagram. For any $x \in S_2$, the effect of applying U to x is to move the dot corresponding to x exactly r places up its column to another dot. It follows that U maps S_2 in a one-to-one fashion into β_i . Thus $\{U(x) : x \in S_2\}$ is a basis for $R(U)$ consisting of b vectors. Hence $\text{rank}(U) = b$, and so $\text{nullity}(U) = m_i - b = a$. But S_1 is a linearly independent subset of $N(U)$ consisting of a vectors; therefore S_1 is a basis for $N(U)$. ■

In the case that $r = 1$, Theorem 7.9 yields the following corollary.

Corollary. *The dimension of E_{λ_i} is n_i . Hence in a Jordan canonical form of T , the number of Jordan blocks corresponding to λ_i equals the dimension of E_{λ_i} .*

Proof. Exercise. ■

We are now able to devise a method for describing the dot diagram in terms of the ranks of operators.

Theorem 7.10. *Let r_j denote the number of dots in the j th row of the dot diagram of T_i , the restriction of T to K_{λ_i} . Then the following statements are true.*

- (a) $r_1 = \dim(V) - \text{rank}(T - \lambda_i I)$.
- (b) $r_j = \text{rank}((T - \lambda_i I)^{j-1}) - \text{rank}((T - \lambda_i I)^j) \quad \text{if } j > 1$.

Proof. By Theorem 7.9, for $j = 1, 2, \dots, p-1$, we have

$$\begin{aligned} r_1 + r_2 + \cdots + r_j &= \text{nullity}((T - \lambda_i I)^j) \\ &= \dim(V) - \text{rank}((T - \lambda_i I)^j). \end{aligned}$$

Hence

$$r_1 = \dim(V) - \text{rank}(T - \lambda_i I),$$

and for $j > 1$,

$$r_j = (r_1 + r_2 + \cdots + r_j) - (r_1 + r_2 + \cdots + r_{j-1})$$

$$\begin{aligned}
 &= [\dim(V) - \text{rank}((T - \lambda_i I)^j)] - [\dim(V) - \text{rank}((T - \lambda_i I)^{j-1})] \\
 &= \text{rank}((T - \lambda_i I)^{j-1}) - \text{rank}((T - \lambda_i I)^j).
 \end{aligned}$$

■

Theorem 7.10 shows that the dot diagram of T_i is completely determined by T and λ_i . Hence we have proved the following result.

Corollary. *For any eigenvalue λ_i of T , the dot diagram of T_i is unique. Thus, subject to the convention that the cycles of generalized eigenvectors for the bases of each generalized eigenspace are listed in order of decreasing length, the Jordan canonical form of a linear operator or a matrix is unique up to the ordering of the eigenvalues.*

We apply these results to find the Jordan canonical forms of two matrices and a linear operator.

Example 2

Let

$$A = \begin{pmatrix} 2 & -1 & 0 & 1 \\ 0 & 3 & -1 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & -1 & 0 & 3 \end{pmatrix}.$$

We find the Jordan canonical form of A and a Jordan canonical basis for the linear operator $T = L_A$. The characteristic polynomial of A is

$$\det(A - tI) = (t - 2)^3(t - 3).$$

Thus A has two distinct eigenvalues, $\lambda_1 = 2$ and $\lambda_2 = 3$, with multiplicities 3 and 1, respectively. Let T_1 and T_2 be the restrictions of L_A to the generalized eigenspaces K_{λ_1} and K_{λ_2} , respectively.

Suppose that β_1 is a Jordan canonical basis for T_1 . Since λ_1 has multiplicity 3, it follows that $\dim(K_{\lambda_1}) = 3$ by Theorem 7.4(c) (p. 480); hence the dot diagram of T_1 has three dots. As we did earlier, let r_j denote the number of dots in the j th row of this dot diagram. Then, by Theorem 7.10,

$$r_1 = 4 - \text{rank}(A - 2I) = 4 - \text{rank} \begin{pmatrix} 0 & -1 & 0 & 1 \\ 0 & 1 & -1 & 0 \\ 0 & 1 & -1 & 0 \\ 0 & -1 & 0 & 1 \end{pmatrix} = 4 - 2 = 2,$$

and

$$r_2 = \text{rank}(A - 2I) - \text{rank}((A - 2I)^2) = 2 - 1 = 1.$$

(Actually, the computation of r_2 is unnecessary in this case because $r_1 = 2$ and the dot diagram only contains three dots.) Hence the dot diagram associated with β_1 is



So

$$A_1 = [\mathbf{T}_1]_{\beta_1} = \begin{pmatrix} 2 & 1 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 2 \end{pmatrix}.$$

Since $\lambda_2 = 3$ has multiplicity 1, it follows that $\dim(\mathbf{K}_{\lambda_2}) = 1$, and consequently any basis β_2 for \mathbf{K}_{λ_2} consists of a single eigenvector corresponding to $\lambda_2 = 3$. Therefore

$$A_2 = [\mathbf{T}_2]_{\beta_2} = (3).$$

Setting $\beta = \beta_1 \cup \beta_2$, we have

$$J = [\mathbf{L}_A]_{\beta} = \begin{pmatrix} 2 & 1 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 3 \end{pmatrix},$$

and so J is the Jordan canonical form of A .

We now find a Jordan canonical basis for $\mathbf{T} = \mathbf{L}_A$. We begin by determining a Jordan canonical basis β_1 for \mathbf{T}_1 . Since the dot diagram of \mathbf{T}_1 has two columns, each corresponding to a cycle of generalized eigenvectors, there are two such cycles. Let v_1 and v_2 denote the end vectors of the first and second cycles, respectively. We reprint below the dot diagram with the dots labeled with the names of the vectors to which they correspond.

$$\begin{array}{cc} \bullet (\mathbf{T} - 2I)(v_1) & \bullet v_2 \\ \bullet v_1 \end{array}$$

From this diagram we see that $v_1 \in \mathbf{N}((\mathbf{T} - 2I)^2)$ but $v_1 \notin \mathbf{N}(\mathbf{T} - 2I)$. Now

$$A - 2I = \begin{pmatrix} 0 & -1 & 0 & 1 \\ 0 & 1 & -1 & 0 \\ 0 & 1 & -1 & 0 \\ 0 & -1 & 0 & 1 \end{pmatrix} \quad \text{and} \quad (A - 2I)^2 = \begin{pmatrix} 0 & -2 & 1 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & -2 & 1 & 1 \end{pmatrix}.$$

It is easily seen that

$$\left\{ \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 2 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 0 \\ 2 \end{pmatrix} \right\}$$

is a basis for $\mathbf{N}((\mathbf{T} - 2I)^2) = \mathbf{K}_{\lambda_1}$. Of these three basis vectors, the last two do not belong to $\mathbf{N}(\mathbf{T} - 2I)$, and hence we select one of these for v_1 . Suppose that we choose

$$v_1 = \begin{pmatrix} 0 \\ 1 \\ 2 \\ 0 \end{pmatrix}.$$

Then

$$(\mathbf{T} - 2\mathbf{I})(v_1) = (A - 2I)(v_1) = \begin{pmatrix} 0 & -1 & 0 & 1 \\ 0 & 1 & -1 & 0 \\ 0 & 1 & -1 & 0 \\ 0 & -1 & 0 & 1 \end{pmatrix} \begin{pmatrix} 0 \\ 1 \\ 2 \\ 0 \end{pmatrix} = \begin{pmatrix} -1 \\ -1 \\ -1 \\ -1 \end{pmatrix}.$$

Now simply choose v_2 to be a vector in E_{λ_1} that is linearly independent of $(\mathbf{T} - 2\mathbf{I})(v_1)$; for example, select

$$v_2 = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}.$$

Thus we have associated the Jordan canonical basis

$$\beta_1 = \left\{ \begin{pmatrix} -1 \\ -1 \\ -1 \\ -1 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 2 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} \right\}$$

with the dot diagram in the following manner.

$$\begin{array}{c} \bullet \begin{pmatrix} -1 \\ -1 \\ -1 \\ -1 \end{pmatrix} \quad \bullet \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} \\ \bullet \begin{pmatrix} 0 \\ 1 \\ 2 \\ 0 \end{pmatrix} \end{array}$$

By Theorem 7.6 (p. 483), the linear independence of β_1 is guaranteed since v_2 was chosen to be linearly independent of $(\mathbf{T} - 2\mathbf{I})(v_1)$.

Since $\lambda_2 = 3$ has multiplicity 1, $\dim(K_{\lambda_2}) = \dim(E_{\lambda_2}) = 1$. Hence any eigenvector of L_A corresponding to $\lambda_2 = 3$ constitutes an appropriate basis β_2 . For example,

$$\beta_2 = \left\{ \begin{pmatrix} 1 \\ 0 \\ 0 \\ 1 \end{pmatrix} \right\}.$$

Thus

$$\beta = \beta_1 \cup \beta_2 = \left\{ \begin{pmatrix} -1 \\ -1 \\ -1 \\ -1 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 2 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \\ 0 \\ 1 \end{pmatrix} \right\}$$

is a Jordan canonical basis for \mathbb{L}_A .

Notice that if

$$Q = \begin{pmatrix} -1 & 0 & 1 & 1 \\ -1 & 1 & 0 & 0 \\ -1 & 2 & 0 & 0 \\ -1 & 0 & 0 & 1 \end{pmatrix},$$

then $J = Q^{-1}AQ$. \blacklozenge

Example 3

Let

$$A = \begin{pmatrix} 2 & -4 & 2 & 2 \\ -2 & 0 & 1 & 3 \\ -2 & -2 & 3 & 3 \\ -2 & -6 & 3 & 7 \end{pmatrix}.$$

We find the Jordan canonical form J of A , a Jordan canonical basis for \mathbb{L}_A , and a matrix Q such that $J = Q^{-1}AQ$.

The characteristic polynomial of A is $\det(A - tI) = (t - 2)^2(t - 4)^2$. Let $T = \mathbb{L}_A$, $\lambda_1 = 2$, and $\lambda_2 = 4$, and let T_i be the restriction of \mathbb{L}_A to K_{λ_i} for $i = 1, 2$.

We begin by computing the dot diagram of T_1 . Let r_1 denote the number of dots in the first row of this diagram. Then

$$r_1 = 4 - \text{rank}(A - 2I) = 4 - 2 = 2;$$

hence the dot diagram of T_1 is as follows.

• •

Therefore

$$A_1 = [T_1]_{\beta_1} = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix},$$

where β_1 is any basis corresponding to the dots. In this case, β_1 is an arbitrary basis for $E_{\lambda_1} = N(T - 2I)$, for example,

$$\beta_1 = \left\{ \begin{pmatrix} 2 \\ 1 \\ 0 \\ 2 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 2 \\ 0 \end{pmatrix} \right\}.$$

Next we compute the dot diagram of T_2 . Since $\text{rank}(A - 4I) = 3$, there is only $4 - 3 = 1$ dot in the first row of the diagram. Since $\lambda_2 = 4$ has

multiplicity 2, we have $\dim(K_{\lambda_2}) = 2$, and hence this dot diagram has the following form:

•
•

Thus

$$A_2 = [T_2]_{\beta_2} = \begin{pmatrix} 4 & 1 \\ 0 & 4 \end{pmatrix},$$

where β_2 is any basis for K_{λ_2} corresponding to the dots. In this case, β_2 is a cycle of length 2. The end vector of this cycle is a vector $v \in K_{\lambda_2} = N((T - 4I)^2)$ such that $v \notin N(T - 4I)$. One way of finding such a vector was used to select the vector v_1 in Example 2. In this example, we illustrate another method. A simple calculation shows that a basis for the null space of $L_A - 4I$ is

$$\left\{ \begin{pmatrix} 0 \\ 1 \\ 1 \\ 1 \end{pmatrix} \right\}.$$

Choose v to be any solution to the system of linear equations

$$(A - 4I)x = \begin{pmatrix} 0 \\ 1 \\ 1 \\ 1 \end{pmatrix},$$

for example,

$$v = \begin{pmatrix} 1 \\ -1 \\ -1 \\ 0 \end{pmatrix}.$$

Thus

$$\beta_2 = \{(L_A - 4I)(v), v\} = \left\{ \begin{pmatrix} 0 \\ 1 \\ 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ -1 \\ -1 \\ 0 \end{pmatrix} \right\}.$$

Therefore

$$\beta = \beta_1 \cup \beta_2 = \left\{ \begin{pmatrix} 2 \\ 1 \\ 0 \\ 2 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 2 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ -1 \\ -1 \\ 0 \end{pmatrix} \right\}$$

is a Jordan canonical basis for L_A . The corresponding Jordan canonical form is given by

$$J = [L_A]_\beta = \begin{pmatrix} A_1 & O \\ O & A_2 \end{pmatrix} = \left(\begin{array}{cc|cc} 2 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ \hline 0 & 0 & 4 & 1 \\ 0 & 0 & 0 & 4 \end{array} \right).$$

Finally, we define Q to be the matrix whose columns are the vectors of β listed in the same order, namely,

$$Q = \begin{pmatrix} 2 & 0 & 0 & 1 \\ 1 & 1 & 1 & -1 \\ 0 & 2 & 1 & -1 \\ 2 & 0 & 1 & 0 \end{pmatrix}.$$

Then $J = Q^{-1}AQ$. \blacklozenge

Example 4

Let V be the vector space of polynomial functions in two real variables x and y of degree at most 2. Then V is a vector space over R and $\alpha = \{1, x, y, x^2, y^2, xy\}$ is an ordered basis for V . Let T be the linear operator on V defined by

$$T(f(x, y)) = \frac{\partial}{\partial x} f(x, y).$$

For example, if $f(x, y) = x + 2x^2 - 3xy + y$, then

$$T(f(x, y)) = \frac{\partial}{\partial x}(x + 2x^2 - 3xy + y) = 1 + 4x - 3y.$$

We find the Jordan canonical form and a Jordan canonical basis for T .

Let $A = [T]_\alpha$. Then

$$A = \begin{pmatrix} 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix},$$

and hence the characteristic polynomial of T is

$$\det(A - tI) = \det \begin{pmatrix} -t & 1 & 0 & 0 & 0 & 0 \\ 0 & -t & 0 & 2 & 0 & 0 \\ 0 & 0 & -t & 0 & 0 & 1 \\ 0 & 0 & 0 & -t & 0 & 0 \\ 0 & 0 & 0 & 0 & -t & 0 \\ 0 & 0 & 0 & 0 & 0 & -t \end{pmatrix} = t^6.$$

Thus $\lambda = 0$ is the only eigenvalue of T , and $K_\lambda = V$. For each j , let r_j denote the number of dots in the j th row of the dot diagram of T . By Theorem 7.10,

$$r_1 = 6 - \text{rank}(A) = 6 - 3 = 3,$$

and since

$$A^2 = \begin{pmatrix} 0 & 0 & 0 & 2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix},$$

$$r_2 = \text{rank}(A) - \text{rank}(A^2) = 3 - 1 = 2.$$

Because there are a total of six dots in the dot diagram and $r_1 = 3$ and $r_2 = 2$, it follows that $r_3 = 1$. So the dot diagram of T is



We conclude that the Jordan canonical form of T is

$$J = \left(\begin{array}{ccc|ccc} 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{array} \right).$$

We now find a Jordan canonical basis for T . Since the first column of the dot diagram of T consists of three dots, we must find a polynomial $f_1(x, y)$ such that $\frac{\partial^2}{\partial x^2} f_1(x, y) \neq 0$. Examining the basis $\alpha = \{1, x, y, x^2, y^2, xy\}$ for $K_\lambda = V$, we see that x^2 is a suitable candidate. Setting $f_1(x, y) = x^2$, we see that

$$(T - \lambda I)(f_1(x, y)) = T(f_1(x, y)) = \frac{\partial}{\partial x}(x^2) = 2x$$

and

$$(T - \lambda I)^2(f_1(x, y)) = T^2(f_1(x, y)) = \frac{\partial^2}{\partial x^2}(x^2) = 2.$$

Likewise, since the second column of the dot diagram consists of two dots, we must find a polynomial $f_2(x, y)$ such that

$$\frac{\partial}{\partial x}(f_2(x, y)) \neq 0, \quad \text{but} \quad \frac{\partial^2}{\partial x^2}(f_2(x, y)) = 0.$$

Since our choice must be linearly independent of the polynomials already chosen for the first cycle, the only choice in α that satisfies these constraints is xy . So we set $f_2(x, y) = xy$. Thus

$$(\mathbf{T} - \lambda\mathbf{I})(f_2(x, y)) = \mathbf{T}(f_2(x, y)) = \frac{\partial}{\partial x}(xy) = y.$$

Finally, the third column of the dot diagram consists of a single polynomial that lies in the null space of \mathbf{T} . The only remaining polynomial in α is y^2 , and it is suitable here. So set $f_3(x, y) = y^2$. Therefore we have identified polynomials with the dots in the dot diagram as follows.

$$\begin{array}{ccc} \bullet 2 & \bullet y & \bullet y^2 \\ \bullet 2x & \bullet xy & \\ \bullet x^2 & & \end{array}$$

Thus $\beta = \{2, 2x, x^2, y, xy, y^2\}$ is a Jordan canonical basis for \mathbf{T} . \spadesuit

In the three preceding examples, we relied on our ingenuity and the context of the problem to find Jordan canonical bases. The reader can do the same in the exercises. We are successful in these cases because the dimensions of the generalized eigenspaces under consideration are small. We do not attempt, however, to develop a general algorithm for computing Jordan canonical bases, although one could be devised by following the steps in the proof of the existence of such a basis (Theorem 7.7 p. 484).

The following result may be thought of as a corollary to Theorem 7.10.

Theorem 7.11. *Let A and B be $n \times n$ matrices, each having Jordan canonical forms computed according to the conventions of this section. Then A and B are similar if and only if they have (up to an ordering of their eigenvalues) the same Jordan canonical form.*

Proof. If A and B have the same Jordan canonical form J , then A and B are each similar to J and hence are similar to each other.

Conversely, suppose that A and B are similar. Then A and B have the same eigenvalues. Let J_A and J_B denote the Jordan canonical forms of A and B , respectively, with the same ordering of their eigenvalues. Then A is similar to both J_A and J_B , and therefore, by the corollary to Theorem 2.23 (p. 115), J_A and J_B are matrix representations of \mathbf{L}_A . Hence J_A and J_B are Jordan canonical forms of \mathbf{L}_A . Thus $J_A = J_B$ by the corollary to Theorem 7.10. \blacksquare

Example 5

We determine which of the matrices

$$A = \begin{pmatrix} -3 & 3 & -2 \\ -7 & 6 & -3 \\ 1 & -1 & 2 \end{pmatrix}, \quad B = \begin{pmatrix} 0 & 1 & -1 \\ -4 & 4 & -2 \\ -2 & 1 & 1 \end{pmatrix},$$

$$C = \begin{pmatrix} 0 & -1 & -1 \\ -3 & -1 & -2 \\ 7 & 5 & 6 \end{pmatrix}, \quad \text{and} \quad D = \begin{pmatrix} 0 & 1 & 2 \\ 0 & 1 & 1 \\ 0 & 0 & 2 \end{pmatrix}$$

are similar. Observe that A , B , and C have the same characteristic polynomial $-(t-1)(t-2)^2$, whereas D has $-t(t-1)(t-2)$ as its characteristic polynomial. Because similar matrices have the same characteristic polynomials, D cannot be similar to A , B , or C . Let J_A , J_B , and J_C be the Jordan canonical forms of A , B , and C , respectively, using the ordering 1, 2 for their common eigenvalues. Then (see Exercise 4)

$$J_A = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 2 & 1 \\ 0 & 0 & 2 \end{pmatrix}, \quad J_B = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 2 \end{pmatrix}, \quad \text{and} \quad J_C = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 2 & 1 \\ 0 & 0 & 2 \end{pmatrix}.$$

Since $J_A = J_C$, A is similar to C . Since J_B is different from J_A and J_C , B is similar to neither A nor C . ◆

The reader should observe that any diagonal matrix is a Jordan canonical form. Thus a linear operator T on a finite-dimensional vector space V is *diagonalizable if and only if its Jordan canonical form is a diagonal matrix*. Hence T is diagonalizable if and only if the Jordan canonical basis for T consists of eigenvectors of T . Similar statements can be made about matrices. Thus, of the matrices A , B , and C in Example 5, A and C are not diagonalizable because their Jordan canonical forms are not diagonal matrices.

EXERCISES

1. Label the following statements as true or false. Assume that the characteristic polynomial of the matrix or linear operator splits.
 - (a) The Jordan canonical form of a diagonal matrix is the matrix itself.
 - (b) Let T be a linear operator on a finite-dimensional vector space V that has a Jordan canonical form J . If β is any basis for V , then the Jordan canonical form of $[T]_\beta$ is J .
 - (c) Linear operators having the same characteristic polynomial are similar.
 - (d) Matrices having the same Jordan canonical form are similar.
 - (e) Every matrix is similar to its Jordan canonical form.
 - (f) Every linear operator with the characteristic polynomial $(-1)^n(t-\lambda)^n$ has the same Jordan canonical form.
 - (g) Every linear operator on a finite-dimensional vector space has a unique Jordan canonical basis.
 - (h) The dot diagrams of a linear operator on a finite-dimensional vector space are unique.

2. Let T be a linear operator on a finite-dimensional vector space V such that the characteristic polynomial of T splits. Suppose that $\lambda_1 = 2$, $\lambda_2 = 4$, and $\lambda_3 = -3$ are the distinct eigenvalues of T and that the dot diagrams for the restriction of T to K_{λ_i} ($i = 1, 2, 3$) are as follows:

$$\begin{array}{ccc} \lambda_1 = 2 & \lambda_2 = 4 & \lambda_3 = -3 \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & & \bullet \\ & & \bullet \end{array}$$

Find the Jordan canonical form J of T .

3. Let T be a linear operator on a finite-dimensional vector space V with Jordan canonical form

$$\left(\begin{array}{ccc|ccc} 2 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 2 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 2 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 3 \end{array} \right).$$

- (a) Find the characteristic polynomial of T .
 - (b) Find the dot diagram corresponding to each eigenvalue of T .
 - (c) For which eigenvalues λ_i , if any, does $E_{\lambda_i} = K_{\lambda_i}$?
 - (d) For each eigenvalue λ_i , find the smallest positive integer p_i for which $K_{\lambda_i} = N((T - \lambda_i I)^{p_i})$.
 - (e) Compute the following numbers for each i , where U_i denotes the restriction of $T - \lambda_i I$ to K_{λ_i} .
 - (i) $\text{rank}(U_i)$
 - (ii) $\text{rank}(U_i^2)$
 - (iii) $\text{nullity}(U_i)$
 - (iv) $\text{nullity}(U_i^2)$
4. For each of the matrices A that follow, find a Jordan canonical form J and an invertible matrix Q such that $J = Q^{-1}AQ$. Notice that the matrices in (a), (b), and (c) are those used in Example 5.

(a) $A = \begin{pmatrix} -3 & 3 & -2 \\ -7 & 6 & -3 \\ 1 & -1 & 2 \end{pmatrix}$ (b) $A = \begin{pmatrix} 0 & 1 & -1 \\ -4 & 4 & -2 \\ -2 & 1 & 1 \end{pmatrix}$

(c) $A = \begin{pmatrix} 0 & -1 & -1 \\ -3 & -1 & -2 \\ 7 & 5 & 6 \end{pmatrix}$ (d) $A = \begin{pmatrix} 0 & -3 & 1 & 2 \\ -2 & 1 & -1 & 2 \\ -2 & 1 & -1 & 2 \\ -2 & -3 & 1 & 4 \end{pmatrix}$

5. For each linear operator T , find a Jordan canonical form J of T and a Jordan canonical basis β for T .
- V is the real vector space of functions spanned by the set of real-valued functions $\{e^t, te^t, t^2e^t, e^{2t}\}$, and T is the linear operator on V defined by $T(f) = f'$.
 - T is the linear operator on $P_3(R)$ defined by $T(f(x)) = xf''(x)$.
 - T is the linear operator on $P_3(R)$ defined by $T(f(x)) = f''(x) + 2f(x)$.
 - T is the linear operator on $M_{2 \times 2}(R)$ defined by

$$T(A) = \begin{pmatrix} 3 & 1 \\ 0 & 3 \end{pmatrix} \cdot A - A^t.$$

- (e) T is the linear operator on $M_{2 \times 2}(R)$ defined by

$$T(A) = \begin{pmatrix} 3 & 1 \\ 0 & 3 \end{pmatrix} \cdot (A - A^t).$$

- (f) V is the vector space of polynomial functions in two real variables x and y of degree at most 2, as defined in Example 4, and T is the linear operator on V defined by

$$T(f(x, y)) = \frac{\partial}{\partial x} f(x, y) + \frac{\partial}{\partial y} f(x, y).$$

6. Let A be an $n \times n$ matrix whose characteristic polynomial splits. Prove that A and A^t have the same Jordan canonical form, and conclude that A and A^t are similar. *Hint:* For any eigenvalue λ of A and A^t and any positive integer r , show that $\text{rank}((A - \lambda I)^r) = \text{rank}((A^t - \lambda I)^r)$.
7. Let A be an $n \times n$ matrix whose characteristic polynomial splits, γ be a cycle of generalized eigenvectors corresponding to an eigenvalue λ , and W be the subspace spanned by γ . Define γ' to be the ordered set obtained from γ by reversing the order of the vectors in γ .
- Prove that $[T_W]_{\gamma'} = ([T_W]_\gamma)^t$.
 - Let J be the Jordan canonical form of A . Use (a) to prove that J and J^t are similar.
 - Use (b) to prove that A and A^t are similar.
8. Let T be a linear operator on a finite-dimensional vector space, and suppose that the characteristic polynomial of T splits. Let β be a Jordan canonical basis for T .
- Prove that for any nonzero scalar c , $\{cx : x \in \beta\}$ is a Jordan canonical basis for T .

- (b) Suppose that γ is one of the cycles of generalized eigenvectors that forms β , and suppose that γ corresponds to the eigenvalue λ and has length greater than 1. Let x be the end vector of γ , and let y be a nonzero vector in E_λ . Let γ' be the ordered set obtained from γ by replacing x by $x + y$. Prove that γ' is a cycle of generalized eigenvectors corresponding to λ , and that if γ' replaces γ in the union that defines β , then the new union is also a Jordan canonical basis for T .
- (c) Apply (b) to obtain a Jordan canonical basis for L_A , where A is the matrix given in Example 2, that is different from the basis given in the example.
9. Suppose that a dot diagram has k columns and m rows with p_j dots in column j and r_i dots in row i . Prove the following results.
- $m = p_1$ and $k = r_1$.
 - $p_j = \max\{i : r_i \geq j\}$ for $1 \leq j \leq k$ and $r_i = \max\{j : p_j \geq i\}$ for $1 \leq i \leq m$. Hint: Use mathematical induction on m .
 - $r_1 \geq r_2 \geq \dots \geq r_m$.
 - Deduce that the number of dots in each column of a dot diagram is completely determined by the number of dots in the rows.
10. Let T be a linear operator whose characteristic polynomial splits, and let λ be an eigenvalue of T .
- Prove that $\dim(K_\lambda)$ is the sum of the lengths of all the blocks corresponding to λ in the Jordan canonical form of T .
 - Deduce that $E_\lambda = K_\lambda$ if and only if all the Jordan blocks corresponding to λ are 1×1 matrices.

The following definitions are used in Exercises 11–19.

Definitions. A linear operator T on a vector space V is called **nilpotent** if $T^p = T_0$ for some positive integer p . An $n \times n$ matrix A is called **nilpotent** if $A^p = O$ for some positive integer p .

11. Let T be a linear operator on a finite-dimensional vector space V , and let β be an ordered basis for V . Prove that T is nilpotent if and only if $[T]_\beta$ is nilpotent.
12. Prove that any square upper triangular matrix with each diagonal entry equal to zero is nilpotent.
13. Let T be a nilpotent operator on an n -dimensional vector space V , and suppose that p is the smallest positive integer for which $T^p = T_0$. Prove the following results.
- $N(T^i) \subseteq N(T^{i+1})$ for every positive integer i .

- (b) There is a sequence of ordered bases $\beta_1, \beta_2, \dots, \beta_p$ such that β_i is a basis for $N(T^i)$ and β_{i+1} contains β_i for $1 \leq i \leq p-1$.
- (c) Let $\beta = \beta_p$ be the ordered basis for $N(T^p) = V$ in (b). Then $[T]_\beta$ is an upper triangular matrix with each diagonal entry equal to zero.
- (d) The characteristic polynomial of T is $(-1)^n t^n$. Hence the characteristic polynomial of T splits, and 0 is the only eigenvalue of T .
14. Prove the converse of Exercise 13(d): If T is a linear operator on an n -dimensional vector space V and $(-1)^n t^n$ is the characteristic polynomial of T , then T is nilpotent.
15. Give an example of a linear operator T on a finite-dimensional vector space over the field of real numbers such that T is not nilpotent, but zero is the only eigenvalue of T . Characterize all such operators. Visit goo.gl/nDjsWm for a solution.
16. Let T be a nilpotent linear operator on a finite-dimensional vector space V . Recall from Exercise 13 that $\lambda = 0$ is the only eigenvalue of T , and hence $V = K_\lambda$. Let β be a Jordan canonical basis for T . Prove that for any positive integer i , if we delete from β the vectors corresponding to the last i dots in each column of a dot diagram of β , the resulting set is a basis for $R(T^i)$. (If a column of the dot diagram contains fewer than i dots, all the vectors associated with that column are removed from β .)
17. Let T be a linear operator on a finite-dimensional vector space V such that the characteristic polynomial of T splits, and let $\lambda_1, \lambda_2, \dots, \lambda_k$ be the distinct eigenvalues of T . For each i , let v_i denote the unique vector in K_{λ_i} such that $x = v_1 + v_2 + \dots + v_k$. (This unique representation is guaranteed by Theorem 7.3 (p. 479).) Define a mapping $S: V \rightarrow V$ by
- $$S(x) = \lambda_1 v_1 + \lambda_2 v_2 + \dots + \lambda_k v_k.$$
- (a) Prove that S is a diagonalizable linear operator on V .
- (b) Let $U = T - S$. Prove that U is nilpotent and commutes with S , that is, $SU = US$.
18. Let T be a linear operator on a finite-dimensional vector space V over C , and let J be the Jordan canonical form of T . Let D be the diagonal matrix whose diagonal entries are the diagonal entries of J , and let $M = J - D$. Prove the following results.
- (a) M is nilpotent.
- (b) $MD = DM$.
- (c) If p is the smallest positive integer for which $M^p = O$, then, for any positive integer $r < p$,

$$J^r = D^r + rD^{r-1}M + \frac{r(r-1)}{2!}D^{r-2}M^2 + \dots + rDM^{r-1} + M^r,$$

and, for any positive integer $r \geq p$,

$$\begin{aligned} J^r &= D^r + rD^{r-1}M + \frac{r(r-1)}{2!}D^{r-2}M^2 + \dots \\ &\quad + \frac{r!}{(r-p+1)!(p-1)!}D^{r-p+1}M^{p-1}. \end{aligned}$$

19. Let $F = C$ and

$$J = \begin{pmatrix} \lambda & 1 & 0 & \cdots & 0 \\ 0 & \lambda & 1 & \cdots & 0 \\ 0 & 0 & \lambda & \cdots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & 0 & 0 & \cdots & 1 \\ 0 & 0 & 0 & \cdots & \lambda \end{pmatrix}$$

be the $m \times m$ Jordan block corresponding to λ , and let $N = J - \lambda I_m$. Prove the following results:

(a) $N^m = O$, and for $1 \leq r < m$,

$$N_{ij}^r = \begin{cases} 1 & \text{if } j = i + r \\ 0 & \text{otherwise.} \end{cases}$$

(b) For any integer $r \geq m$,

$$J^r = \begin{pmatrix} \lambda^r & r\lambda^{r-1} & \frac{r(r-1)}{2!}\lambda^{r-2} & \cdots & \frac{r(r-1)\cdots(r-m+2)}{(m-1)!}\lambda^{r-m+1} \\ 0 & \lambda^r & r\lambda^{r-1} & \cdots & \frac{r(r-1)\cdots(r-m+3)}{(m-2)!}\lambda^{r-m+2} \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & 0 & 0 & \cdots & \lambda^r \end{pmatrix}.$$

(c) $\lim_{r \rightarrow \infty} J^r$ exists if and only if one of the following holds:

- (i) $|\lambda| < 1$.
- (ii) $\lambda = 1$ and $m = 1$.

(Note that $\lim_{r \rightarrow \infty} \lambda^r$ exists under these conditions. See the discussion preceding Theorem 5.12 on page 284.) Furthermore, $\lim_{r \rightarrow \infty} J^r$ is the zero matrix if condition (i) holds and is the 1×1 matrix (1) if condition (ii) holds.

(d) Prove Theorem 5.12 on page 284.

The following definition is used in Exercises 20 and 21.

Definition. For any $A \in M_{n \times n}(C)$, define the norm of A by

$$\|A\|_m = \max \{|A_{ij}| : 1 \leq i, j \leq n\}.$$

20. Let $A, B \in M_{n \times n}(C)$. Prove the following results.
- (a) $\|A\|_m \geq 0$.
 - (b) $\|A\|_m = 0$ if and only if $A = O$.
 - (c) $\|cA\|_m = |c| \cdot \|A\|_m$ for any scalar c .
 - (d) $\|A + B\|_m \leq \|A\|_m + \|B\|_m$.
 - (e) $\|AB\|_m \leq n\|A\|_m\|B\|_m$.
21. Let $A \in M_{n \times n}(C)$ be a transition matrix. (See Section 5.3.) Since C is an algebraically closed field, A has a Jordan canonical form J to which A is similar. Let P be an invertible matrix such that $P^{-1}AP = J$. Prove the following results.
- (a) $\|A^k\|_m \leq 1$ for every positive integer k .
 - (b) There exists a positive number c such that $\|J^k\|_m \leq c$ for every positive integer k .
 - (c) Each Jordan block of J corresponding to the eigenvalue $\lambda = 1$ is a 1×1 matrix.
 - (d) $\lim_{k \rightarrow \infty} A^k$ exists if and only if 1 is the only eigenvalue of A with absolute value 1.
 - (e) Theorem 5.19(a), using (c) and Theorem 5.18.

The next exercise requires knowledge of absolutely convergent series as well as the definition of e^A for a matrix A . (See page 310.)

22. Use Exercise 20(d) to prove that e^A exists for every $A \in M_{n \times n}(C)$.
23. Let $x' = Ax$ be a system of n linear differential equations, where x is an n -tuple of differentiable functions $x_1(t), x_2(t), \dots, x_n(t)$ of the real variable t , and A is an $n \times n$ coefficient matrix as in Exercise 16 of Section 5.2. In contrast to that exercise, however, do not assume that A is diagonalizable, but assume that the characteristic polynomial of A splits. Let $\lambda_1, \lambda_2, \dots, \lambda_k$ be the distinct eigenvalues of A .
- (a) Prove that if u is the end vector of a cycle of generalized eigenvectors of L_A of length p and u corresponds to the eigenvalue λ_i , then for any polynomial $f(t)$ of degree less than p , the function

$$e^{\lambda_i t}[f(t)(A - \lambda_i I)^{p-1} + f'(t)(A - \lambda_i I)^{p-2} + \dots + f^{(p-1)}(t)]u$$

is a solution to the system $x' = Ax$.

- (b) Prove that the general solution to $x' = Ax$ is a sum of the functions of the form given in (a), where the vectors u are the end vectors of the distinct cycles that constitute a fixed Jordan canonical basis for L_A .
24. Use Exercise 23 to find the general solution to each of the following systems of linear equations, where x , y , and z are real-valued differentiable functions of the real variable t .

$$(a) \begin{array}{rcl} x' = 2x + y \\ y' = 2y - z \\ z' = 3z \end{array} \quad (b) \begin{array}{rcl} x' = 2x + y \\ y' = 2y + z \\ z' = 2z \end{array}$$

7.3 THE MINIMAL POLYNOMIAL

The Cayley–Hamilton theorem (Theorem 5.22 p. 315) tells us that for any linear operator T on an n -dimensional vector space, there is a polynomial $f(t)$ of degree n such that $f(T) = T_0$, namely, the characteristic polynomial of T . Hence there is a polynomial of least degree with this property, and this degree is at most n . If $g(t)$ is such a polynomial, we can divide $g(t)$ by its leading coefficient to obtain another polynomial $p(t)$ of the same degree with leading coefficient 1, that is, $p(t)$ is a *monic* polynomial. (See Appendix E.)

Definition. Let T be a linear operator on a finite-dimensional vector space. A polynomial $p(t)$ is called a **minimal polynomial** of T if $p(t)$ is a monic polynomial of least positive degree for which $p(T) = T_0$.

The preceding discussion shows that every linear operator on a finite-dimensional vector space has a minimal polynomial. The next result shows that it is unique, and hence we can speak of *the* minimal polynomial of T .

Theorem 7.12. Let $p(t)$ be a minimal polynomial of a linear operator T on a finite-dimensional vector space V .

- (a) For any polynomial $g(t)$, if $g(T) = T_0$, then $p(t)$ divides $g(t)$. In particular, $p(t)$ divides the characteristic polynomial of T .
- (b) The minimal polynomial of T is unique.

Proof. (a) Let $g(t)$ be a polynomial for which $g(T) = T_0$. By the division algorithm for polynomials (Theorem E.1 of Appendix E, p. 555), there exist polynomials $q(t)$ and $r(t)$ such that

$$g(t) = q(t)p(t) + r(t), \tag{1}$$

where $r(t)$ has degree less than the degree of $p(t)$. Substituting T into (1) and using that $g(T) = p(T) = T_0$, we have $r(T) = T_0$. Since $r(t)$ has degree less than $p(t)$ and $p(t)$ is the minimal polynomial of T , $r(t)$ must be the zero polynomial. Thus (1) simplifies to $g(t) = q(t)p(t)$, proving (a).

(b) Suppose that $p_1(t)$ and $p_2(t)$ are each minimal polynomials of \mathbf{T} . Then $p_1(t)$ divides $p_2(t)$ by (a). Since $p_1(t)$ and $p_2(t)$ have the same degree, we have that $p_2(t) = cp_1(t)$ for some nonzero scalar c . Because $p_1(t)$ and $p_2(t)$ are monic, $c = 1$; hence $p_1(t) = p_2(t)$. ■

The minimal polynomial of a linear operator has an obvious analog for a matrix.

Definition. Let $A \in M_{n \times n}(F)$. The **minimal polynomial** $p(t)$ of A is the monic polynomial of least positive degree for which $p(A) = O$.

The following results are now immediate.

Theorem 7.13. Let \mathbf{T} be a linear operator on a finite-dimensional vector space V , and let β be an ordered basis for V . Then the minimal polynomial of \mathbf{T} is the same as the minimal polynomial of $[\mathbf{T}]_\beta$.

Proof. Exercise. ■

Corollary. For any $A \in M_{n \times n}(F)$, the minimal polynomial of A is the same as the minimal polynomial of L_A .

Proof. Exercise. ■

In view of the preceding theorem and corollary, Theorem 7.12 and all subsequent theorems in this section that are stated for operators are also valid for matrices.

For the remainder of this section, we study primarily minimal polynomials of operators (and hence matrices) whose characteristic polynomials split. A more general treatment of minimal polynomials is given in Section 7.4.

Theorem 7.14. Let \mathbf{T} be a linear operator on a finite-dimensional vector space V , and let $p(t)$ be the minimal polynomial of \mathbf{T} . A scalar λ is an eigenvalue of \mathbf{T} if and only if $p(\lambda) = 0$. Hence the characteristic polynomial and the minimal polynomial of \mathbf{T} have the same zeros.

Proof. Let $f(t)$ be the characteristic polynomial of \mathbf{T} . Since $p(t)$ divides $f(t)$, there exists a polynomial $q(t)$ such that $f(t) = q(t)p(t)$. If λ is a zero of $p(t)$, then

$$f(\lambda) = q(\lambda)p(\lambda) = q(\lambda) \cdot 0 = 0.$$

So λ is a zero of $f(t)$; that is, λ is an eigenvalue of \mathbf{T} .

Conversely, suppose that λ is an eigenvalue of \mathbf{T} , and let $x \in V$ be an eigenvector corresponding to λ . By Exercise 22 of Section 5.1, we have

$$0 = \mathbf{T}_0(x) = p(\mathbf{T})(x) = p(\lambda)x.$$

Since $x \neq 0$, it follows that $p(\lambda) = 0$, and so λ is a zero of $p(t)$. ■

The following corollary is immediate.

Corollary. Let T be a linear operator on a finite-dimensional vector space V with minimal polynomial $p(t)$ and characteristic polynomial $f(t)$. Suppose that $f(t)$ factors as

$$f(t) = (\lambda_1 - t)^{n_1}(\lambda_2 - t)^{n_2} \cdots (\lambda_k - t)^{n_k},$$

where $\lambda_1, \lambda_2, \dots, \lambda_k$ are the distinct eigenvalues of T . Then there exist integers m_1, m_2, \dots, m_k such that $1 \leq m_i \leq n_i$ for all i and

$$p(t) = (t - \lambda_1)^{m_1}(t - \lambda_2)^{m_2} \cdots (t - \lambda_k)^{m_k}.$$

Example 1

We compute the minimal polynomial of the matrix

$$A = \begin{pmatrix} 3 & -1 & 0 \\ 0 & 2 & 0 \\ 1 & -1 & 2 \end{pmatrix}.$$

Since A has the characteristic polynomial

$$f(t) = \det \begin{pmatrix} 3-t & -1 & 0 \\ 0 & 2-t & 0 \\ 1 & -1 & 2-t \end{pmatrix} = -(t-2)^2(t-3),$$

the minimal polynomial of A must be either $(t-2)(t-3)$ or $(t-2)^2(t-3)$ by the corollary to Theorem 7.14. Substituting A into $p(t) = (t-2)(t-3)$, we find that $p(A) = O$; hence $p(t)$ is the minimal polynomial of A . ♦

Example 2

Let T be the linear operator on \mathbb{R}^2 defined by

$$T(a, b) = (2a + 5b, 6a + b)$$

and β be the standard ordered basis for \mathbb{R}^2 . Then

$$[T]_\beta = \begin{pmatrix} 2 & 5 \\ 6 & 1 \end{pmatrix},$$

and hence the characteristic polynomial of T is

$$f(t) = \det \begin{pmatrix} 2-t & 5 \\ 6 & 1-t \end{pmatrix} = (t-7)(t+4).$$

Thus the minimal polynomial of T is also $(t-7)(t+4)$. ♦

Example 3

Let D be the linear operator on $P_2(R)$ defined by $D(g(x)) = g'(x)$, the derivative of $g(x)$. We compute the minimal polynomial of T . Let β be the standard ordered basis for $P_2(R)$. Then

$$[D]_{\beta} = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 2 \\ 0 & 0 & 0 \end{pmatrix},$$

and it follows that the characteristic polynomial of D is $-t^3$. So by the corollary to Theorem 7.14, the minimal polynomial of D is t , t^2 , or t^3 . Since $D^2(x^2) = 2 \neq 0$, it follows that $D^2 \neq T_0$; hence the minimal polynomial of D must be t^3 . ♦

In Example 3, it is easily verified that $P_2(R)$ is a D -cyclic subspace (of itself). Here the minimal and characteristic polynomials are of the same degree. This is no coincidence.

Theorem 7.15. *Let T be a linear operator on an n -dimensional vector space V such that V is a T -cyclic subspace of itself. Then the characteristic polynomial $f(t)$ and the minimal polynomial $p(t)$ have the same degree, and hence $f(t) = (-1)^n p(t)$.*

Proof. Since V is a T -cyclic space, there exists an $x \in V$ such that

$$\beta = \{x, T(x), \dots, T^{n-1}(x)\}$$

is a basis for V (Theorem 5.21 p. 314). Let

$$g(t) = a_0 + a_1 t + \cdots + a_k t^k,$$

be a polynomial of degree $k < n$. Then $a_k \neq 0$ and

$$g(T)(x) = a_0 x + a_1 T(x) + \cdots + a_k T^k(x),$$

and so $g(T)(x)$ is a linear combination of the vectors of β having at least one nonzero coefficient, namely, a_k . Since β is linearly independent, it follows that $g(T)(x) \neq 0$; hence $g(T) \neq T_0$. Therefore the minimal polynomial of T has degree n , which is also the degree of the characteristic polynomial of T . ■

Theorem 7.15 gives a condition under which the degree of the minimal polynomial of an operator is as large as possible. We now investigate the other extreme. By Theorem 7.14, the degree of the minimal polynomial of an operator must be greater than or equal to the number of distinct eigenvalues of the operator. The next result shows that the operators for which the degree of the minimal polynomial is as small as possible are precisely the diagonalizable operators.

Theorem 7.16. Let T be a linear operator on a finite-dimensional vector space V . Then T is diagonalizable if and only if the minimal polynomial of T is of the form

$$p(t) = (t - \lambda_1)(t - \lambda_2) \cdots (t - \lambda_k),$$

where $\lambda_1, \lambda_2, \dots, \lambda_k$ are the distinct eigenvalues of T .

Proof. Suppose that T is diagonalizable. Let $\lambda_1, \lambda_2, \dots, \lambda_k$ be the distinct eigenvalues of T , and define

$$p(t) = (t - \lambda_1)(t - \lambda_2) \cdots (t - \lambda_k).$$

By Theorem 7.14, $p(t)$ divides the minimal polynomial of T . Let $\beta = \{v_1, v_2, \dots, v_n\}$ be a basis for V consisting of eigenvectors of T , and consider any $v_i \in \beta$. Then $(T - \lambda_j I)(v_i) = 0$ for some eigenvalue λ_j . Since $t - \lambda_j$ divides $p(t)$, there is a polynomial $q_j(t)$ such that $p(t) = q_j(t)(t - \lambda_j)$. Hence

$$p(T)(v_i) = q_j(T)(T - \lambda_j I)(v_i) = 0.$$

It follows that $p(T) = T_0$, since $p(T)$ takes each vector in a basis for V into 0. Therefore $p(t)$ is the minimal polynomial of T .

Conversely, suppose that there are distinct scalars $\lambda_1, \lambda_2, \dots, \lambda_k$ such that the minimal polynomial $p(t)$ of T factors as

$$p(t) = (t - \lambda_1)(t - \lambda_2) \cdots (t - \lambda_k).$$

By Theorem 7.14, the λ_i 's are eigenvalues of T . We apply mathematical induction on $n = \dim(V)$. Clearly T is diagonalizable for $n = 1$. Now assume that T is diagonalizable whenever $\dim(V) < n$ for some $n > 1$, and let $\dim(V) = n$ and $W = R(T - \lambda_k I)$. Obviously $W \neq V$, because λ_k is an eigenvalue of T . If $W = \{0\}$, then $T = \lambda_k I$, which is clearly diagonalizable. So suppose that $0 < \dim(W) < n$. Then W is T -invariant, and for any $x \in W$,

$$(T - \lambda_1 I)(T - \lambda_2 I) \cdots (T - \lambda_{k-1} I)(x) = 0.$$

It follows that the minimal polynomial of T_W divides the polynomial $(t - \lambda_1)(t - \lambda_2) \cdots (t - \lambda_{k-1})$. Hence by the induction hypothesis, T_W is diagonalizable. Furthermore, λ_k is not an eigenvalue of T_W by Theorem 7.14. Therefore $W \cap N(T - \lambda_k I) = \{0\}$. Now let $\beta_1 = \{v_1, v_2, \dots, v_m\}$ be a basis for W consisting of eigenvectors of T_W (and hence of T), and let $\beta_2 = \{w_1, w_2, \dots, w_p\}$ be a basis for $N(T - \lambda_k I)$, the eigenspace of T corresponding to λ_k . Then β_1 and β_2 are disjoint by the previous comment. Moreover, $m + p = n$ by the dimension theorem applied to $T - \lambda_k I$. We show that $\beta = \beta_1 \cup \beta_2$ is linearly independent. Consider scalars a_1, a_2, \dots, a_m and b_1, b_2, \dots, b_p such that

$$a_1 v_1 + a_2 v_2 + \cdots + a_m v_m + b_1 w_1 + b_2 w_2 + \cdots + b_p w_p = 0.$$

Let

$$x = \sum_{i=1}^m a_i v_i \quad \text{and} \quad y = \sum_{i=1}^p b_i w_i.$$

Then $x \in W$, $y \in N(T - \lambda_k I)$, and $x + y = 0$. It follows that $x = -y \in W \cap N(T - \lambda_k I)$, and therefore $x = 0$. Since β_1 is linearly independent, we have that $a_1 = a_2 = \dots = a_m = 0$. Similarly, $b_1 = b_2 = \dots = b_p = 0$, and we conclude that β is a linearly independent subset of V consisting of n eigenvectors. It follows that β is a basis for V consisting of eigenvectors of T , and consequently T is diagonalizable. ■

In addition to diagonalizable operators, there are methods for determining the minimal polynomial of any linear operator on a finite-dimensional vector space. In the case that the characteristic polynomial of the operator splits, the minimal polynomial can be described using the Jordan canonical form of the operator. (See Exercise 13.) In the case that the characteristic polynomial does not split, the minimal polynomial can be described using the *rational canonical form*, which we study in the next section. (See Exercise 7 of Section 7.4.)

Example 4

We determine all matrices $A \in M_{2 \times 2}(R)$ for which $A^2 - 3A + 2I = O$. Let $g(t) = t^2 - 3t + 2 = (t - 1)(t - 2)$. Since $g(A) = O$, the minimal polynomial $p(t)$ of A divides $g(t)$. Hence the only possible candidates for $p(t)$ are $t - 1$, $t - 2$, and $(t - 1)(t - 2)$. If $p(t) = t - 1$ or $p(t) = t - 2$, then $A = I$ or $A = 2I$, respectively. If $p(t) = (t - 1)(t - 2)$, then A is diagonalizable with eigenvalues 1 and 2, and hence A is similar to

$$D = \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix},$$

that is, $A = QDQ^{-1}$ for some invertible matrix Q . ♦

Example 5

Let $A \in M_{n \times n}(R)$ satisfy $A^3 = A$. We show that A is diagonalizable. Let $g(t) = t^3 - t = t(t + 1)(t - 1)$. Then $g(A) = O$, and hence the minimal polynomial $p(t)$ of A divides $g(t)$. Since $g(t)$ has no repeated factors, neither does $p(t)$. Thus A is diagonalizable by Theorem 7.16. ♦

Example 6

In Example 3, we saw that the minimal polynomial of the differential operator D on $P_2(R)$ is t^3 . Hence, by Theorem 7.16, D is not diagonalizable. ♦

EXERCISES

1. Label the following statements as true or false. Assume that all vector spaces are finite-dimensional.
 - (a) Every linear operator T has a polynomial $p(t)$ of largest degree for which $p(T) = T_0$.
 - (b) Every linear operator has a unique minimal polynomial.
 - (c) The characteristic polynomial of a linear operator divides the minimal polynomial of that operator.
 - (d) The minimal and the characteristic polynomials of any diagonalizable operator are equal.
 - (e) Let T be a linear operator on an n -dimensional vector space V , $p(t)$ be the minimal polynomial of T , and $f(t)$ be the characteristic polynomial of T . Suppose that $f(t)$ splits. Then $f(t)$ divides $[p(t)]^n$.
 - (f) The minimal polynomial of a linear operator always has the same degree as the characteristic polynomial of the operator.
 - (g) A linear operator is diagonalizable if its minimal polynomial splits.
 - (h) Let T be a linear operator on a vector space V such that V is a T -cyclic subspace of itself. Then the degree of the minimal polynomial of T equals $\dim(V)$.
 - (i) Let T be a linear operator on a vector space V such that T has n distinct eigenvalues, where $n = \dim(V)$. Then the degree of the minimal polynomial of T equals n .
2. Find the minimal polynomial of each of the following matrices.

(a) $\begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}$	(b) $\begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}$
(c) $\begin{pmatrix} 4 & -14 & 5 \\ 1 & -4 & 2 \\ 1 & -6 & 4 \end{pmatrix}$	(d) $\begin{pmatrix} 3 & 0 & 1 \\ 2 & 2 & 2 \\ -1 & 0 & 1 \end{pmatrix}$
3. For each linear operator T on V , find the minimal polynomial of T .
 - (a) $V = \mathbb{R}^2$ and $T(a, b) = (a + b, a - b)$
 - (b) $V = P_2(R)$ and $T(g(x)) = g'(x) + 2g(x)$
 - (c) $V = P_2(R)$ and $T(f(x)) = -xf''(x) + f'(x) + 2f(x)$
 - (d) $V = M_{n \times n}(R)$ and $T(A) = A^t$. Hint: Note that $T^2 = I$.
4. Determine which of the matrices and operators in Exercises 2 and 3 are diagonalizable.
5. Describe all linear operators T on \mathbb{R}^2 such that T is diagonalizable and $T^3 - 2T^2 + T = T_0$.

6. Prove Theorem 7.13 and its corollary.
7. Prove the corollary to Theorem 7.14.
8. Let T be a linear operator on a finite-dimensional vector space, and let $p(t)$ be the minimal polynomial of T . Prove the following results.
 - (a) T is invertible if and only if $p(0) \neq 0$.
 - (b) If T is invertible and $p(t) = t^m + a_{m-1}t^{m-1} + \cdots + a_1t + a_0$, then

$$T^{-1} = -\frac{1}{a_0} (T^{m-1} + a_{m-1}T^{m-2} + \cdots + a_2T + a_1I).$$

9. Let T be a diagonalizable linear operator on a finite-dimensional vector space V . Prove that V is a T -cyclic subspace if and only if each of the eigenspaces of T is one-dimensional.
10. Let T be a linear operator on a finite-dimensional vector space V , and suppose that W is a T -invariant subspace of V . Prove that the minimal polynomial of T_W divides the minimal polynomial of T .
11. Let $g(t)$ be the auxiliary polynomial associated with a homogeneous linear differential equation with constant coefficients (as defined in Section 2.7), and let V denote the solution space of this differential equation. Prove the following results.
 - (a) V is a D -invariant subspace, where D is the differentiation operator on C^∞ .
 - (b) The minimal polynomial of D_V (the restriction of D to V) is $g(t)$.
 - (c) If the degree of $g(t)$ is n , then the characteristic polynomial of D_V is $(-1)^n g(t)$.

Hint: Use Theorem 2.32 (p. 135) for (b) and (c).

12. Let D be the differentiation operator on $P(R)$, the space of polynomials over R . Prove that there exists no polynomial $g(t)$ for which $g(D) = T_0$. Hence D has no minimal polynomial.
13. Let T be a linear operator on a finite-dimensional vector space, and suppose that the characteristic polynomial of T splits. Let $\lambda_1, \lambda_2, \dots, \lambda_k$ be the distinct eigenvalues of T , and for each i let p_i be the number of rows in the largest Jordan block corresponding to λ_i in a Jordan canonical form of T . Prove that the minimal polynomial of T is

$$(t - \lambda_1)^{p_1}(t - \lambda_2)^{p_2} \cdots (t - \lambda_k)^{p_k}.$$

The following exercise requires knowledge of direct sums (see Section 5.2).

14. Let T be a linear operator on a finite-dimensional vector space V , and let W_1 and W_2 be T -invariant subspaces of V such that $V = W_1 \oplus W_2$. Suppose that $p_1(t)$ and $p_2(t)$ are the minimal polynomials of T_{W_1} and T_{W_2} , respectively. Either prove that the minimal polynomial $f(t)$ of T always equals $p_1(t)p_2(t)$ or give an example in which $f(t) \neq p_1(t)p_2(t)$.

Exercise 15 uses the following definition.

Definition. Let T be a linear operator on a finite-dimensional vector space V , and let x be a nonzero vector in V . The polynomial $p(t)$ is called a **T -annihilator** of x if $p(t)$ is a monic polynomial of least degree for which $p(T)(x) = 0$.

- 15.[†] Let T be a linear operator on a finite-dimensional vector space V , and let x be a nonzero vector in V . Prove the following results.

- (a) The vector x has a unique T -annihilator.
- (b) The T -annihilator of x divides any polynomial $g(t)$ for which $g(T) = T_0$.
- (c) If $p(t)$ is the T -annihilator of x and W is the T -cyclic subspace generated by x , then $p(t)$ is the minimal polynomial of T_W , and $\dim(W)$ equals the degree of $p(t)$.
- (d) The degree of the T -annihilator of x is 1 if and only if x is an eigenvector of T .

Visit goo.gl/8KD6Gw for a solution.

16. Let T be a linear operator on a finite-dimensional vector space V , and let W_1 be a T -invariant subspace of V . Let $x \in V$ such that $x \notin W_1$. Prove the following results.

- (a) There exists a unique monic polynomial $g_1(t)$ of least positive degree such that $g_1(T)(x) \in W_1$.
- (b) If $h(t)$ is a polynomial for which $h(T)(x) \in W_1$, then $g_1(t)$ divides $h(t)$.
- (c) $g_1(t)$ divides the minimal and the characteristic polynomials of T .
- (d) Let W_2 be a T -invariant subspace of V such that $W_2 \subseteq W_1$, and let $g_2(t)$ be the unique monic polynomial of least degree such that $g_2(T)(x) \in W_2$. Then $g_1(t)$ divides $g_2(t)$.

7.4* THE RATIONAL CANONICAL FORM

Until now we have used eigenvalues, eigenvectors, and generalized eigenvectors in our analysis of linear operators with characteristic polynomials that

split. In general, characteristic polynomials need not split, and indeed, operators need not have eigenvalues! However, the unique factorization theorem for polynomials (see page 562) guarantees that the characteristic polynomial $f(t)$ of any linear operator T on an n -dimensional vector space factors uniquely as

$$f(t) = (-1)^n(\phi_1(t))^{n_1}(\phi_2(t))^{n_2} \cdots (\phi_k(t))^{n_k},$$

where the $\phi_i(t)$'s ($1 \leq i \leq k$) are distinct irreducible monic polynomials and the n_i 's are positive integers. In the case that $f(t)$ splits, each irreducible monic polynomial factor is of the form $\phi_i(t) = t - \lambda_i$, where λ_i is an eigenvalue of T , and there is a one-to-one correspondence between eigenvalues of T and the irreducible monic factors of the characteristic polynomial. In general, eigenvalues need not exist, but the irreducible monic factors always exist. In this section, we establish structure theorems based on the irreducible monic factors of the characteristic polynomial instead of eigenvalues.

In this context, the following definition is the appropriate replacement for eigenspace and generalized eigenspace.

Definition. Let T be a linear operator on a finite-dimensional vector space V with characteristic polynomial

$$f(t) = (-1)^n(\phi_1(t))^{n_1}(\phi_2(t))^{n_2} \cdots (\phi_k(t))^{n_k},$$

where the $\phi_i(t)$'s ($1 \leq i \leq k$) are distinct irreducible monic polynomials and the n_i 's are positive integers. For $1 \leq i \leq k$, we define the subset K_{ϕ_i} of V by

$$K_{\phi_i} = \{x \in V : (\phi_i(T))^p(x) = 0 \text{ for some positive integer } p\}.$$

We show that each K_{ϕ_i} is a nonzero T -invariant subspace of V . Note that if $\phi_i(t) = t - \lambda$ is of degree one, then K_{ϕ_i} is the generalized eigenspace of T corresponding to the eigenvalue λ .

Having obtained suitable generalizations of the related concepts of eigenvalue and eigenspace, our next task is to describe a canonical form of a linear operator suitable to this context. The one that we study is called the *rational canonical form*. Since a canonical form is a description of a matrix representation of a linear operator, it can be defined by specifying the form of the ordered bases allowed for these representations.

Here the bases of interest naturally arise from the generators of certain cyclic subspaces. For this reason, the reader should recall the definition of a T -cyclic subspace generated by a vector and Theorem 5.21 (p. 314). We briefly review this concept and introduce some new notation and terminology.

Let T be a linear operator on a finite-dimensional vector space V , and let x be a nonzero vector in V . We use the notation C_x for the T -cyclic subspace generated by x . Recall (Theorem 5.21) that if $\dim(C_x) = k$, then the set

$$\{x, T(x), T^2(x), \dots, T^{k-1}(x)\}$$

is an ordered basis for C_x . To distinguish this basis from all other ordered bases for C_x , we call it the **T-cyclic basis generated by x** and denote it by β_x . Let A be the matrix representation of the restriction of T to C_x in the ordered basis β_x . Recall from the proof of Theorem 5.21 that

$$A = \begin{pmatrix} 0 & 0 & \cdots & 0 & -a_0 \\ 1 & 0 & \cdots & 0 & -a_1 \\ 0 & 1 & \cdots & 0 & -a_2 \\ \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & -a_{k-1} \end{pmatrix},$$

where

$$a_0x + a_1T(x) + \cdots + a_{k-1}T^{k-1}(x) + T^k(x) = 0.$$

Furthermore, the characteristic polynomial of A is given by

$$\det(A - tI) = (-1)^k(a_0 + a_1t + \cdots + a_{k-1}t^{k-1} + t^k).$$

The matrix A is called the **companion matrix** of the monic polynomial $h(t) = a_0 + a_1t + \cdots + a_{k-1}t^{k-1} + t^k$. Every monic polynomial has a companion matrix, and the characteristic polynomial of the companion matrix of a monic polynomial $g(t)$ of degree k is equal to $(-1)^k g(t)$. (See Exercise 19 of Section 5.4.) By Theorem 7.15 (p. 512), the monic polynomial $h(t)$ is also the minimal polynomial of A . Since A is the matrix representation of the restriction of T to C_x , $h(t)$ is also the minimal polynomial of this restriction. By Exercise 15 of Section 7.3, $h(t)$ is also the **T-annihilator** of x .

It is the object of this section to prove that for every linear operator T on a finite-dimensional vector space V , there exists an ordered basis β for V such that the matrix representation $[T]_\beta$ is of the form

$$\begin{pmatrix} C_1 & O & \cdots & O \\ O & C_2 & \cdots & O \\ \vdots & \vdots & & \vdots \\ O & O & \cdots & C_r \end{pmatrix},$$

where each C_i is the companion matrix of a polynomial $(\phi(t))^m$ such that $\phi(t)$ is a monic irreducible divisor of the characteristic polynomial of T and m is a positive integer. A matrix representation of this kind is called a **rational canonical form** of T . We call the accompanying basis a **rational canonical basis** for T .

The next theorem is a simple consequence of the following lemma, which relies on the concept of **T-annihilator**, introduced in the Exercises of Section 7.3.

Lemma. Let T be a linear operator on a finite-dimensional vector space V , let x be a nonzero vector in V , and suppose that the T -annihilator of x is of the form $(\phi(t))^p$ for some irreducible monic polynomial $\phi(t)$. Then $\phi(t)$ divides the minimal polynomial of T , and $x \in K_\phi$.

Proof. By Exercise 15(b) of Section 7.3, $(\phi(t))^p$ divides the minimal polynomial of T . Therefore $\phi(t)$ divides the minimal polynomial of T . Furthermore, $x \in K_\phi$ by the definition of K_ϕ . ■

Theorem 7.17. Let T be a linear operator on a finite-dimensional vector space V , and let β be an ordered basis for V . Then β is a rational canonical basis for T if and only if β is the disjoint union of T -cyclic bases β_{v_i} , where each v_i lies in K_ϕ for some irreducible monic divisor $\phi(t)$ of the characteristic polynomial of T .

Proof. Exercise. ■

Example 1

Suppose that T is a linear operator on \mathbb{R}^8 and

$$\beta = \{v_1, v_2, v_3, v_4, v_5, v_6, v_7, v_8\}$$

is a rational canonical basis for T such that

$$C = [T]_\beta = \left(\begin{array}{cc|cccccc} 0 & -3 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 & 0 & -1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & -2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & -1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \end{array} \right)$$

is a rational canonical form of T . In this case, the submatrices C_1 , C_2 , and C_3 are the companion matrices of the polynomials $\phi_1(t)$, $(\phi_2(t))^2$, and $\phi_2(t)$, respectively, where

$$\phi_1(t) = t^2 - t + 3 \quad \text{and} \quad \phi_2(t) = t^2 + 1.$$

In the context of Theorem 7.17, β is the disjoint union of the T -cyclic bases; that is,

$$\begin{aligned} \beta &= \beta_{v_1} \cup \beta_{v_3} \cup \beta_{v_7} \\ &= \{v_1, v_2\} \cup \{v_3, v_4, v_5, v_6\} \cup \{v_7, v_8\}. \end{aligned}$$

By Exercise 39 of Section 5.4, the characteristic polynomial $f(t)$ of T is the product of the characteristic polynomials of the companion matrices:

$$f(t) = \phi_1(t)(\phi_2(t))^2\phi_2(t) = \phi_1(t)(\phi_2(t))^3. \quad \blacklozenge$$

The rational canonical form C of the operator T in Example 1 is constructed from matrices of the form C_i , each of which is the companion matrix of some power of a monic irreducible divisor of the characteristic polynomial of T . Furthermore, each such divisor is used in this way at least once.

In the course of showing that every linear operator T on a finite dimensional vector space has a rational canonical form C , we show that the companion matrices C_i that constitute C are always constructed from powers of the monic irreducible divisors of the characteristic polynomial of T . A key role in our analysis is played by the subspaces K_ϕ , where $\phi(t)$ is an irreducible monic divisor of the minimal polynomial of T . Since the minimal polynomial of an operator divides the characteristic polynomial of the operator, every irreducible divisor of the former is also an irreducible divisor of the latter. We eventually show that the converse is also true; that is, the minimal polynomial and the characteristic polynomial have the same irreducible divisors.

We begin with a result that lists several properties of irreducible divisors of the minimal polynomial. The reader is advised to review the definition of T -annihilator and the accompanying Exercise 15 of Section 7.3.

Theorem 7.18. *Let T be a linear operator on a finite-dimensional vector space V , and suppose that*

$$p(t) = (\phi_1(t))^{m_1}(\phi_2(t))^{m_2} \cdots (\phi_k(t))^{m_k}$$

is the minimal polynomial of T , where the $\phi_i(t)$'s ($1 \leq i \leq k$) are the distinct irreducible monic factors of $p(t)$ and the m_i 's are positive integers. Then the following statements are true.

- (a) K_{ϕ_i} is a nonzero T -invariant subspace of V for each i .
- (b) If x is a nonzero vector in some K_{ϕ_i} , then the T -annihilator of x is of the form $(\phi_i(t))^p$ for some integer p .
- (c) $K_{\phi_i} \cap K_{\phi_j} = \{0\}$ for $i \neq j$.
- (d) K_{ϕ_i} is invariant under $\phi_j(T)$ for $i \neq j$, and the restriction of $\phi_j(T)$ to K_{ϕ_i} is one-to-one and onto.
- (e) $K_{\phi_i} = N((\phi_i(T))^{m_i})$ for each i .

Proof. If $k = 1$, then (a), (b), and (e) are obvious, while (c) and (d) are vacuously true. Now suppose that $k > 1$.

(a) The proof that K_{ϕ_i} is a T -invariant subspace of V is left as an exercise. Let $f_i(t)$ be the polynomial obtained from $p(t)$ by omitting the factor $(\phi_i(t))^{m_i}$. To prove that K_{ϕ_i} is nonzero, first observe that $f_i(t)$ is a proper divisor of $p(t)$; therefore there exists a vector $z \in V$ such that $x = f_i(T)(z) \neq 0$. Then $x \in K_{\phi_i}$ because

$$(\phi_i(T))^{m_i}(x) = (\phi_i(T))^{m_i}f_i(T)(z) = p(T)(z) = 0.$$

(b) Assume the hypothesis. Then $(\phi_i(T))^q(x) = 0$ for some positive integer q . Hence the T -annihilator of x divides $(\phi_i(t))^q$ by Exercise 15(b) of Section 7.3, and the result follows.

(c) Assume $i \neq j$. Let $x \in K_{\phi_i} \cap K_{\phi_j}$, and suppose that $x \neq 0$. By (b), the T -annihilator of x is a power of both $\phi_i(t)$ and $\phi_j(t)$. But this is impossible because $\phi_i(t)$ and $\phi_j(t)$ are relatively prime (see Appendix E). We conclude that $x = 0$.

(d) Assume $i \neq j$. Since K_{ϕ_i} is T -invariant, it is also $\phi_j(T)$ -invariant. Suppose that $\phi_j(T)(x) = 0$ for some $x \in K_{\phi_i}$. Then $x \in K_{\phi_i} \cap K_{\phi_j} = \{0\}$ by (c). Therefore the restriction of $\phi_j(T)$ to K_{ϕ_i} is one-to-one. Since V is finite-dimensional, this restriction is also onto.

(e) Suppose that $1 \leq i \leq k$. Clearly, $N((\phi_i(T))^{m_i}) \subseteq K_{\phi_i}$. Let $f_i(t)$ be the polynomial defined in (a). Since $f_i(t)$ is a product of polynomials of the form $\phi_j(t)$ for $j \neq i$, we have by (d) that the restriction of $f_i(T)$ to K_{ϕ_i} is onto. Let $x \in K_{\phi_i}$. Then there exists $y \in K_{\phi_i}$ such that $f_i(T)(y) = x$. Therefore

$$((\phi_i(T))^{m_i})(x) = ((\phi_i(T))^{m_i})f_i(T)(y) = p(T)(y) = 0,$$

and hence $x \in N((\phi_i(T))^{m_i})$. Thus $K_{\phi_i} = N((\phi_i(T))^{m_i})$. ■

Since a rational canonical basis for an operator T is obtained from a union of T -cyclic bases, we need to know when such a union is linearly independent. The next major result, Theorem 7.19, reduces this problem to the study of T -cyclic bases within K_ϕ , where $\phi(t)$ is an irreducible monic divisor of the minimal polynomial of T . We begin with the following lemma.

Lemma. *Let T be a linear operator on a finite-dimensional vector space V , and suppose that*

$$p(t) = (\phi_1(t))^{m_1}(\phi_2(t))^{m_2} \cdots (\phi_k(t))^{m_k}$$

is the minimal polynomial of T , where the ϕ_i 's ($1 \leq i \leq k$) are the distinct irreducible monic factors of $p(t)$ and the m_i 's are positive integers. For $1 \leq i \leq k$, let $v_i \in K_{\phi_i}$ be such that

$$v_1 + v_2 + \cdots + v_k = 0. \tag{2}$$

Then $v_i = 0$ for all i .

Proof. The result is trivial if $k = 1$, so suppose that $k > 1$. Consider any i . Let $f_i(t)$ be the polynomial obtained from $p(t)$ by omitting the factor $(\phi_i(t))^{m_i}$. As a consequence of Theorem 7.18, $f_i(T)$ is one-to-one on K_{ϕ_i} , and $f_i(T)(v_j) = 0$ for $i \neq j$. Thus, applying $f_i(T)$ to (2), we obtain $f_i(T)(v_i) = 0$, from which it follows that $v_i = 0$. ■

Theorem 7.19. *Let T be a linear operator on a finite-dimensional vector space V , and suppose that*

$$p(t) = (\phi_1(t))^{m_1}(\phi_2(t))^{m_2} \cdots (\phi_k(t))^{m_k}$$

is the minimal polynomial of T , where the ϕ_i 's ($1 \leq i \leq k$) are the distinct irreducible monic factors of $p(t)$ and the m_i 's are positive integers. For $1 \leq i \leq k$, let S_i be a linearly independent subset of K_{ϕ_i} . Then

- (a) $S_i \cap S_j = \emptyset$ for $i \neq j$
- (b) $S_1 \cup S_2 \cup \dots \cup S_k$ is linearly independent.

Proof. If $k = 1$, then (a) is vacuously true and (b) is obvious. Now suppose that $k > 1$. Then (a) follows immediately from Theorem 7.18(c). Furthermore, the proof of (b) is identical to the proof of Theorem 5.5 (p. 261) with the eigenvectors replaced by the generalized eigenvectors. ■

In view of Theorem 7.19, we can focus on bases of individual spaces of the form K_ϕ , where $\phi(t)$ is an irreducible monic divisor of the minimal polynomial of T . The next several results give us ways to construct bases for these spaces that are unions of T -cyclic bases. These results serve the dual purposes of leading to the existence theorem for the rational canonical form and of providing methods for constructing rational canonical bases.

For Theorems 7.20 and 7.21 and the latter's corollary, we fix a linear operator T on a finite-dimensional vector space V and an irreducible monic divisor $\phi(t)$ of the minimal polynomial of T .

Theorem 7.20. Let v_1, v_2, \dots, v_k be distinct vectors in K_ϕ such that

$$S_1 = \beta_{v_1} \cup \beta_{v_2} \cup \dots \cup \beta_{v_k}$$

is linearly independent. For each i , suppose there exists $w_i \in V$ such that $\phi(T)(w_i) = v_i$. Then

$$S_2 = \beta_{w_1} \cup \beta_{w_2} \cup \dots \cup \beta_{w_k}$$

is also linearly independent.

Proof. Consider any linear combination of vectors in S_2 that sums to zero, say,

$$\sum_{i=1}^k \sum_{j=0}^{n_i} a_{ij} T^j(w_i) = 0. \quad (3)$$

For each i , let $f_i(t)$ be the polynomial defined by

$$f_i(t) = \sum_{j=0}^{n_i} a_{ij} t^j.$$

Then (3) can be rewritten as

$$\sum_{i=1}^k f_i(T)(w_i) = 0. \quad (4)$$

Apply $\phi(T)$ to both sides of (4) to obtain

$$\sum_{i=1}^k \phi(T)f_i(T)(w_i) = \sum_{i=1}^k f_i(T)\phi(T)(w_i) = \sum_{i=1}^k f_i(T)(v_i) = 0.$$

This last sum can be rewritten as a linear combination of the vectors in S_1 so that each $f_i(T)(v_i)$ is a linear combination of the vectors in β_{v_i} . Since S_1 is linearly independent, it follows that

$$f_i(T)(v_i) = 0 \quad \text{for all } i.$$

Therefore the T -annihilator of v_i divides $f_i(t)$ for all i . (See Exercise 15 of Section 7.3.) By Theorem 7.18(b), $\phi(t)$ divides the T -annihilator of v_i , and hence $\phi(t)$ divides $f_i(t)$ for all i . Thus, for each i , there exists a polynomial $g_i(t)$ such that $f_i(t) = g_i(t)\phi(t)$. So (4) becomes

$$\sum_{i=1}^k g_i(T)\phi(T)(w_i) = \sum_{i=1}^k g_i(T)(v_i) = 0.$$

Again, linear independence of S_1 requires that

$$f_i(T)(w_i) = g_i(T)(v_i) = 0 \quad \text{for all } i.$$

But $f_i(T)(w_i)$ is the result of grouping the terms of the linear combination in (3) that arise from the linearly independent set β_{w_i} . We conclude that for each i , $a_{ij} = 0$ for all j . Therefore S_2 is linearly independent. ■

We now show that K_ϕ has a basis consisting of a union of T -cycles.

Lemma. Let W be a T -invariant subspace of K_ϕ , and let β be a basis for W . Then the following statements are true.

- (a) Suppose that $x \in N(\phi(T))$, but $x \notin W$. Then $\beta \cup \beta_x$ is linearly independent.
- (b) For some w_1, w_2, \dots, w_s in $N(\phi(T))$, β can be extended to the linearly independent set

$$\beta' = \beta \cup \beta_{w_1} \cup \beta_{w_2} \cup \dots \cup \beta_{w_s},$$

whose span contains $N(\phi(T))$.

Proof. (a) Let $\beta = \{v_1, v_2, \dots, v_k\}$, and suppose that

$$\sum_{i=1}^k a_i v_i + z = 0 \quad \text{and} \quad z = \sum_{j=0}^{d-1} b_j T^j(x),$$

where d is the degree of $\phi(t)$. Then $z \in C_x \cap W$, and hence $C_z \subseteq C_x \cap W$. Suppose that $z \neq 0$. Then z has $\phi(t)$ as its T -annihilator, and therefore

$$d = \dim(C_z) \leq \dim(C_x \cap W) \leq \dim(C_x) = d.$$

It follows that $C_x \cap W = C_x$, and consequently $x \in W$, contrary to hypothesis. Therefore $z = 0$, from which it follows that $b_j = 0$ for all j . Since β is linearly independent, it follows that $a_i = 0$ for all i . Thus $\beta \cup \beta_x$ is linearly independent.

(b) Suppose that W does not contain $N(\phi(T))$. Choose a vector $w_1 \in N(\phi(T))$ that is not in W . By (a), $\beta_1 = \beta \cup \beta_{w_1}$ is linearly independent. Let $W_1 = \text{span}(\beta_1)$. If W_1 does not contain $N(\phi(T))$, choose a vector w_2 in $N(\phi(T))$, but not in W_1 , so that $\beta_2 = \beta_1 \cup \beta_{w_2} = \beta \cup \beta_{w_1} \cup \beta_{w_2}$ is linearly independent. Continuing this process, we eventually obtain vectors w_1, w_2, \dots, w_s in $N(\phi(T))$ such that the union

$$\beta' = \beta \cup \beta_{w_1} \cup \beta_{w_2} \cup \dots \cup \beta_{w_s}$$

is a linearly independent set whose span contains $N(\phi(T))$. ■

Theorem 7.21. *If the minimal polynomial of T is of the form $p(t) = (\phi(t))^m$, then there exists a rational canonical basis for T .*

Proof. The proof is by mathematical induction on m . Suppose that $m = 1$. Apply (b) of the lemma to $W = \{0\}$ to obtain a linearly independent subset of V of the form $\beta_{v_1} \cup \beta_{v_2} \cup \dots \cup \beta_{v_k}$, whose span contains $N(\phi(T))$. Since $V = N(\phi(T))$, this set is a rational canonical basis for V .

Now suppose that, for some integer $m > 1$, the result is valid whenever the minimal polynomial of T is of the form $(\phi(t))^k$, where $k < m$, and assume that the minimal polynomial of T is $p(t) = (\phi(t))^m$. Let $r = \text{rank}(\phi(T))$. Then $R(\phi(T))$ is a T -invariant subspace of V , and the restriction of T to this subspace has $(\phi(t))^{m-1}$ as its minimal polynomial. Therefore we may apply the induction hypothesis to obtain a rational canonical basis for the restriction of T to $R(T)$. Suppose that v_1, v_2, \dots, v_k are the generating vectors of the T -cyclic bases that constitute this rational canonical basis. For each i , choose w_i in V such that $v_i = \phi(T)(w_i)$. By Theorem 7.20, the union β of the sets β_{w_i} is linearly independent. Let $W = \text{span}(\beta)$. Then W contains $R(\phi(T))$. Apply (b) of the lemma and adjoin additional T -cyclic bases $\beta_{w_{k+1}}, \beta_{w_{k+2}}, \dots, \beta_{w_s}$ to β , if necessary, where w_i is in $N(\phi(T))$ for $i \geq k$, to obtain a linearly independent set

$$\beta' = \beta_{w_1} \cup \beta_{w_2} \cup \dots \cup \beta_{w_k} \cup \dots \cup \beta_{w_s}$$

whose span W' contains both W and $N(\phi(T))$.

We show that $W' = V$. Let U denote the restriction of $\phi(T)$ to W' , which is $\phi(T)$ -invariant. By the way in which W' was obtained from $R(\phi(T))$, it follows that $R(U) = R(\phi(T))$ and $N(U) = N(\phi(T))$. Therefore

$$\begin{aligned}\dim(W') &= \text{rank}(U) + \text{nullity}(U) \\ &= \text{rank}(\phi(T)) + \text{nullity}(\phi(T)) \\ &= \dim(V).\end{aligned}$$

Thus $W' = V$, and β' is a rational canonical basis for T . ■

Corollary. K_ϕ has a basis consisting of the union of T -cyclic bases.

Proof. Apply Theorem 7.21 to the restriction of T to K_ϕ . ■

We are now ready to study the general case.

Theorem 7.22. Every linear operator on a finite-dimensional vector space has a rational canonical basis and, hence, a rational canonical form.

Proof. Let T be a linear operator on a finite-dimensional vector space V , and let $p(t) = (\phi_1(t))^{m_1}(\phi_2(t))^{m_2} \cdots (\phi_k(t))^{m_k}$ be the minimal polynomial of T , where the $\phi_i(t)$'s are the distinct irreducible monic factors of $p(t)$ and $m_i > 0$ for all i . The proof is by mathematical induction on k . The case $k = 1$ is proved in Theorem 7.21.

Suppose that the result is valid whenever the minimal polynomial contains fewer than k distinct irreducible factors for some $k > 1$, and suppose that $p(t)$ contains k distinct factors. Let U be the restriction of T to the T -invariant subspace $W = R((\phi_k(T))^{m_k})$, and let $q(t)$ be the minimal polynomial of U . Then $q(t)$ divides $p(t)$ by Exercise 10 of Section 7.3. Furthermore, $\phi_k(t)$ does not divide $q(t)$. For otherwise, there would exist a nonzero vector $x \in W$ such that $\phi_k(U)(x) = 0$ and a vector $y \in V$ such that $x = (\phi_k(T))^{m_k}(y)$. It follows that $(\phi_k(T))^{m_k+1}(y) = 0$, and hence $y \in K_{\phi_k}$ and $x = (\phi_k(T))^{m_k}(y) = 0$ by Theorem 7.18(e), a contradiction. Thus $q(t)$ contains fewer than k distinct irreducible divisors. So by the induction hypothesis, U has a rational canonical basis β_1 consisting of a union of U -cyclic bases (and hence T -cyclic bases) of vectors from some of the subspaces K_{ϕ_i} , $1 \leq i \leq k-1$. By the corollary to Theorem 7.21, K_{ϕ_k} has a basis β_2 consisting of a union of T -cyclic bases. By Theorem 7.19, β_1 and β_2 are disjoint, and $\beta = \beta_1 \cup \beta_2$ is linearly independent. Let s denote the number of vectors in β . Then

$$\begin{aligned}s &= \dim(R((\phi_k(T))^{m_k})) + \dim(K_{\phi_k}) \\ &= \text{rank}((\phi_k(T))^{m_k}) + \text{nullity}((\phi_k(T))^{m_k}) \\ &= n.\end{aligned}$$

We conclude that β is a basis for V . Therefore β is a rational canonical basis, and T has a rational canonical form. ■

In our study of the rational canonical form, we relied on the minimal polynomial. We are now able to relate the rational canonical form to the characteristic polynomial.

Theorem 7.23. *Let T be a linear operator on an n -dimensional vector space V with characteristic polynomial*

$$f(t) = (-1)^n(\phi_1(t))^{n_1}(\phi_2(t))^{n_2} \cdots (\phi_k(t))^{n_k},$$

where the $\phi_i(t)$'s ($1 \leq i \leq k$) are distinct irreducible monic polynomials and the n_i 's are positive integers. Then the following statements are true.

- (a) $\phi_1(t), \phi_2(t), \dots, \phi_k(t)$ are the irreducible monic factors of the minimal polynomial.
- (b) For each i , $\dim(K_{\phi_i}) = d_i n_i$, where d_i is the degree of $\phi_i(t)$.
- (c) If β is a rational canonical basis for T , then $\beta_i = \beta \cap K_{\phi_i}$ is a basis for K_{ϕ_i} for each i .
- (d) If γ_i is a basis for K_{ϕ_i} for each i , then $\gamma = \gamma_1 \cup \gamma_2 \cup \cdots \cup \gamma_k$ is a basis for V . In particular, if each γ_i is a disjoint union of T -cyclic bases, then γ is a rational canonical basis for T .

Proof. (a) By Theorem 7.22, T has a rational canonical form C . By Exercise 39 of Section 5.4, the characteristic polynomial of C , and hence of T , is the product of the characteristic polynomials of the companion matrices that compose C . Therefore each irreducible monic divisor $\phi_i(t)$ of $f(t)$ divides the characteristic polynomial of at least one of the companion matrices, and hence for some integer p , $(\phi_i(t))^p$ is the T -annihilator of a nonzero vector of V . We conclude that $(\phi_i(t))^p$, and so $\phi_i(t)$, divides the minimal polynomial of T . Conversely, if $\phi(t)$ is an irreducible monic polynomial that divides the minimal polynomial of T , then $\phi(t)$ divides the characteristic polynomial of T because the minimal polynomial divides the characteristic polynomial.

(b), (c), and (d) Let $C = [T]_{\beta}$, which is a rational canonical form of T . Consider any i ($1 \leq i \leq k$). Since $f(t)$ is the product of the characteristic polynomials of the companion matrices that compose C , we may multiply those characteristic polynomials that arise from the T -cyclic bases in β_i to obtain the factor $(\phi_i(t))^{n_i}$ of $f(t)$. Since this polynomial has degree $n_i d_i$, and the union of these bases is a linearly independent subset β_i of K_{ϕ_i} , we have

$$n_i d_i \leq \dim(K_{\phi_i}).$$

Furthermore, $n = \sum_{i=1}^k d_i n_i$, because this sum is equal to the degree of $f(t)$.

Now let s denote the number of vectors in γ . By Theorem 7.19, γ is linearly independent, and therefore

$$n = \sum_{i=1}^k d_i n_i \leq \sum_{i=1}^k \dim(K_{\phi_i}) = s \leq n.$$

Hence $n = s$, and $d_i n_i = \dim(K_{\phi_i})$ for all i . It follows that γ is a basis for V and β_i is a basis for K_{ϕ_i} for each i . ■

Uniqueness of the Rational Canonical Form

Having shown that a rational canonical form exists, we are now in a position to ask about the extent to which it is unique. Certainly, the rational canonical form of a linear operator T can be modified by permuting the T -cyclic bases that constitute the corresponding rational canonical basis. This has the effect of permuting the companion matrices that make up the rational canonical form. As in the case of the Jordan canonical form, we show that except for these permutations, the rational canonical form is unique, although the rational canonical bases are not.

To simplify this task, we adopt the convention of ordering every rational canonical basis so that all the T -cyclic bases associated with the same irreducible monic divisor of the characteristic polynomial are grouped together. Furthermore, within each such grouping, we arrange the T -cyclic bases in decreasing order of size. Our task is to show that, subject to this order, the rational canonical form of a linear operator is unique up to the arrangement of the irreducible monic divisors.

As in the case of the Jordan canonical form, we introduce arrays of dots from which we can reconstruct the rational canonical form. For the Jordan canonical form, we devised a dot diagram for each eigenvalue of the given operator. In the case of the rational canonical form, we define a dot diagram for each irreducible monic divisor of the characteristic polynomial of the given operator. A proof that the resulting dot diagrams are completely determined by the operator is also a proof that the rational canonical form is unique.

In what follows, T is a linear operator on a finite-dimensional vector space with rational canonical basis β ; $\phi(t)$ is an irreducible monic divisor of the characteristic polynomial of T ; $\beta_{v_1}, \beta_{v_2}, \dots, \beta_{v_k}$ are the T -cyclic bases of β that are contained in K_ϕ ; and d is the degree of $\phi(t)$. For each j , let $(\phi(t))^{p_j}$ be the annihilator of v_j . This polynomial has degree $d p_j$; therefore, by Exercise 15 of Section 7.3, β_{v_j} contains $d p_j$ vectors. Furthermore, $p_1 \geq p_2 \geq \dots \geq p_k$ since the T -cyclic bases are arranged in decreasing order of size. We define the **dot diagram** of $\phi(t)$ to be the array consisting of k columns of dots with p_j dots in the j th column, arranged so that the j th column begins at the top and terminates after p_j dots. For example, if $k = 3$, $p_1 = 4$, $p_2 = 2$, and $p_3 = 2$, then the dot diagram is

$$\begin{array}{ccc} \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \vdots & & \vdots \\ \bullet & & \end{array}$$

Although each column of a dot diagram corresponds to a T -cyclic basis

β_{v_i} in K_ϕ , there are fewer dots in the column than there are vectors in the basis.

Example 2

Recall the linear operator T of Example 1 with the rational canonical basis β and the rational canonical form $C = [T]_\beta$. Since there are two irreducible monic divisors of the characteristic polynomial of T , $\phi_1(t) = t^2 - t + 3$ and $\phi_2(t) = t^2 + 1$, there are two dot diagrams to consider. Because $\phi_1(t)$ is the T -annihilator of v_1 and β_{v_1} is a basis for K_{ϕ_1} , the dot diagram for $\phi_1(t)$ consists of a single dot. The other two T -cyclic bases, β_{v_3} and β_{v_7} , lie in K_{ϕ_2} . Since v_3 has T -annihilator $(\phi_2(t))^2$ and v_7 has T -annihilator $\phi_2(t)$, in the dot diagram of $\phi_2(t)$ we have $p_1 = 2$ and $p_2 = 1$. These diagrams are as follows:

Dot diagram for $\phi_1(t)$ Dot diagram for $\phi_2(t)$ 

In practice, we obtain the rational canonical form of a linear operator from the information provided by dot diagrams. This is illustrated in the next example.

Example 3

Let T be a linear operator on a finite-dimensional vector space over R , and suppose that the irreducible monic divisors of the characteristic polynomial of T are

$$\phi_1(t) = t - 1, \quad \phi_2(t) = t^2 + 2, \quad \text{and} \quad \phi_3(t) = t^2 + t + 1.$$

Suppose, furthermore, that the dot diagrams associated with these divisors are as follows:

Diagram for $\phi_1(t)$ Diagram for $\phi_2(t)$ Diagram for $\phi_3(t)$

Since the dot diagram for $\phi_1(t)$ has two columns, it contributes two companion matrices to the rational canonical form. The first column has two dots, and therefore corresponds to the 2×2 companion matrix of $(\phi_1(t))^2 = (t - 1)^2$. The second column, with only one dot, corresponds to the 1×1 companion matrix of $\phi_1(t) = t - 1$. These two companion matrices are given by

$$C_1 = \begin{pmatrix} 0 & -1 \\ 1 & 2 \end{pmatrix} \quad \text{and} \quad C_2 = (1).$$

The dot diagram for $\phi_2(t) = t^2 + 2$ consists of two columns, each containing a single dot; hence this diagram contributes two copies of the 2×2 companion

matrix for $\phi_2(t)$, namely,

$$C_3 = C_4 = \begin{pmatrix} 0 & -2 \\ 1 & 0 \end{pmatrix}.$$

The dot diagram for $\phi_3(t) = t^2 + t + 1$ consists of a single column with a single dot contributing the single 2×2 companion matrix

$$C_5 = \begin{pmatrix} 0 & -1 \\ 1 & -1 \end{pmatrix}.$$

Therefore the rational canonical form of T is the 9×9 matrix

$$C = \begin{pmatrix} C_1 & O & O & O & O \\ O & C_2 & O & O & O \\ O & O & C_3 & O & O \\ O & O & O & C_4 & O \\ O & O & O & O & C_5 \end{pmatrix}$$

$$= \left(\begin{array}{cc|ccccc|c} 0 & -1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 2 & 0 & 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & -2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & -2 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & -1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{array} \right). \quad \blacklozenge$$

We return to the general problem of finding dot diagrams. As we did before, we fix a linear operator T on a finite-dimensional vector space and an irreducible monic divisor $\phi(t)$ of the characteristic polynomial of T . Let U denote the restriction of the linear operator $\phi(T)$ to K_ϕ . By Theorem 7.18(d), $U^q = T_0$ for some positive integer q . Consequently, by Exercise 12 of Section 7.2, the characteristic polynomial of U is $(-1)^m t^m$, where $m = \dim(K_\phi)$. Therefore K_ϕ is the generalized eigenspace of U corresponding to $\lambda = 0$, and U has a Jordan canonical form. The dot diagram associated with the Jordan canonical form of U gives us a key to understanding the dot diagram of T that is associated with $\phi(t)$. We now relate the two diagrams.

Let β be a rational canonical basis for T , and $\beta_{v_1}, \beta_{v_2}, \dots, \beta_{v_k}$ be the T -cyclic bases of β that are contained in K_ϕ . Consider one of these T -cyclic bases β_{v_j} , and suppose again that the T -annihilator of v_j is $(\phi(t))^{p_j}$. Then β_{v_j} consists of $d p_j$ vectors in β . For $0 \leq i < d$, let γ_i be the cycle of generalized eigenvectors of U corresponding to $\lambda = 0$ with end vector $T^i(v_j)$,

where $T^0(v_j) = b_j$. Then

$$\gamma_i = \{(\phi(T))^{p_j-1}T^i(v_j), (\phi(T))^{p_j-2}T^i(v_j), \dots, (\phi(T))T^i(v_j), T^i(v_j)\}.$$

By Theorem 7.1 (p. 478), γ_i is a linearly independent subset of C_{v_i} . Now let

$$\alpha_j = \gamma_0 \cup \gamma_1 \cup \dots \cup \gamma_{d-1}.$$

Notice that α_j contains dp_j vectors.

Lemma 1. α_j is an ordered basis for C_{v_j} .

Proof. The key to this proof is Theorem 7.4 (p. 480). Since α_j is the union of cycles of generalized eigenvectors of U corresponding to $\lambda = 0$, it suffices to show that the set of initial vectors of these cycles

$$\{(\phi(T))^{p_j-1}(v_j), (\phi(T))^{p_j-2}T(v_j), \dots, (\phi(T))^{p_j-1}T^{d-1}(v_j)\}$$

is linearly independent. Consider any linear combination of these vectors

$$a_0(\phi(T))^{p_j-1}(v_j) + a_1(\phi(T))^{p_j-2}T(v_j) + \dots + a_{d-1}(\phi(T))^{p_j-1}T^{d-1}(v_j),$$

where not all of the coefficients are zero. Let $g(t)$ be the polynomial defined by $g(t) = a_0 + a_1t + \dots + a_{d-1}t^{d-1}$. Then $g(t)$ is a nonzero polynomial of degree less than d , and hence $(\phi(t))^{p_j-1}g(t)$ is a nonzero polynomial with degree less than dp_j . Since $(\phi(t))^{p_j}$ is the T -annihilator of v_j , it follows that $(\phi(t))^{p_j-1}g(T)(v_j) \neq 0$. Therefore the set of initial vectors is linearly independent. So by Theorem 7.4, α_j is linearly independent, and the γ_i 's are disjoint. Consequently, α_j consists of dp_j linearly independent vectors in C_{v_j} , which has dimension dp_j . We conclude that α_j is a basis for C_{v_j} . ■

Thus we may replace β_{v_j} by α_j as a basis for C_{v_j} . We do this for each j to obtain a subset $\alpha = \alpha_1 \cup \alpha_2 \cup \dots \cup \alpha_k$ of K_ϕ .

Lemma 2. α is a Jordan canonical basis for K_ϕ .

Proof. Since $\beta_{v_1} \cup \beta_{v_2} \cup \dots \cup \beta_{v_k}$ is a basis for K_ϕ , and since $\text{span}(\alpha_i) = \text{span}(\beta_{v_i}) = C_{v_i}$, Exercise 9 implies that α is a basis for K_ϕ . Because α is a union of cycles of generalized eigenvectors of U , we conclude that α is a Jordan canonical basis. ■

We are now in a position to relate the dot diagram of T corresponding to $\phi(t)$ to the dot diagram of U , bearing in mind that in the first case we are considering a rational canonical form and in the second case we are considering a Jordan canonical form. For convenience, we designate the first diagram D_1 , and the second diagram D_2 . For each j , the presence of the T -cyclic basis β_{x_j} results in a column of p_j dots in D_1 . By Lemma 1, this basis is

replaced by the union α_j of d cycles of generalized eigenvectors of U , each of length p_j , which becomes part of the Jordan canonical basis for U . In effect, α_j determines d columns each containing p_j dots in D_2 . So each column in D_1 determines d columns in D_2 of the same length, and all columns in D_2 are obtained in this way. Alternatively, each row in D_2 has d times as many dots as the corresponding row in D_1 . Since Theorem 7.10 (p. 493) gives us the number of dots in any row of D_2 , we may divide the appropriate expression in this theorem by d to obtain the number of dots in the corresponding row of D_1 . Thus we have the following result.

Theorem 7.24. *Let T be a linear operator on a finite-dimensional vector space V , let $\phi(t)$ be an irreducible monic divisor of the characteristic polynomial of T of degree d , and let r_i denote the number of dots in the i th row of the dot diagram for $\phi(t)$ with respect to a rational canonical basis for T . Then*

- (a) $r_1 = \frac{1}{d}[\dim(V) - \text{rank}(\phi(T))];$
- (b) $r_i = \frac{1}{d}[\text{rank}((\phi(T))^{i-1}) - \text{rank}((\phi(T))^i)] \quad \text{for } i > 1.$

Thus the dot diagrams associated with a rational canonical form of an operator are completely determined by the operator. Since the rational canonical form is completely determined by its dot diagrams, we have the following uniqueness condition.

Corollary. *Under the conventions described earlier, the rational canonical form of a linear operator is unique up to the arrangement of the irreducible monic divisors of the characteristic polynomial.*

Since the rational canonical form of a linear operator is unique, the polynomials corresponding to the companion matrices that determine this form are also unique. These polynomials, which are powers of the irreducible monic divisors, are called the **elementary divisors** of the linear operator. Since a companion matrix may occur more than once in a rational canonical form, the same is true for the elementary divisors. We call the number of such occurrences the **multiplicity** of the elementary divisor.

Conversely, the elementary divisors and their multiplicities determine the companion matrices and, therefore, the rational canonical form of a linear operator.

Example 4

Let

$$\beta = \{e^x \cos 2x, e^x \sin 2x, xe^x \cos 2x, xe^x \sin 2x\}$$

be viewed as a subset of $\mathcal{F}(R, R)$, the space of all real-valued functions defined on R , and let $V = \text{span}(\beta)$. Then V is a four-dimensional subspace of $\mathcal{F}(R, R)$,

and β is an ordered basis for V . Let D be the linear operator on V defined by $D(y) = y'$, the derivative of y , and let $A = [D]_\beta$. Then

$$A = \begin{pmatrix} 1 & 2 & 1 & 0 \\ -2 & 1 & 0 & 1 \\ 0 & 0 & 1 & 2 \\ 0 & 0 & -2 & 1 \end{pmatrix},$$

and the characteristic polynomial of D , and hence of A , is

$$f(t) = (t^2 - 2t + 5)^2.$$

Thus $\phi(t) = t^2 - 2t + 5$ is the only irreducible monic divisor of $f(t)$. Since $\phi(t)$ has degree 2 and V is four-dimensional, the dot diagram for $\phi(t)$ contains only two dots. Therefore the dot diagram is determined by r_1 , the number of dots in the first row. Because ranks are preserved under matrix representations, we can use A in place of D in the formula given in Theorem 7.24. Now

$$\phi(A) = \begin{pmatrix} 0 & 0 & 0 & 4 \\ 0 & 0 & -4 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix},$$

and so

$$r_1 = \frac{1}{2}[4 - \text{rank}(\phi(A))] = \frac{1}{2}[4 - 2] = 1.$$

It follows that the second dot lies in the second row, and the dot diagram is as follows:

•
•

Hence V is a D -cyclic space generated by a single function with D -annihilator $(\phi(t))^2$. Furthermore, its rational canonical form is given by the companion matrix of $(\phi(t))^2 = t^4 - 4t^3 + 14t^2 - 20t + 25$, which is

$$\begin{pmatrix} 0 & 0 & 0 & -25 \\ 1 & 0 & 0 & 20 \\ 0 & 1 & 0 & -14 \\ 0 & 0 & 1 & 4 \end{pmatrix}.$$

Thus $(\phi(t))^2$ is the only elementary divisor of D , and it has multiplicity 1. For the cyclic generator, it suffices to find a function g in V for which $\phi(D)(g) \neq 0$. Since $\phi(A)(e_3) \neq 0$, it follows that $\phi(D)(xe^x \cos 2x) \neq 0$; therefore $g(x) = xe^x \cos 2x$ can be chosen as the cyclic generator. Hence

$$\beta_g = \{xe^x \cos 2x, D(xe^x \cos 2x), D^2(xe^x \cos 2x), D^3(xe^x \cos 2x)\}$$

is a rational canonical basis for D . Notice that the function h defined by $h(x) = xe^x \sin 2x$ can be chosen in place of g . This shows that the rational canonical basis is not unique. ◆

It is convenient to refer to the rational canonical form and elementary divisors of a matrix, which are defined in the obvious way.

Definitions. Let $A \in M_{n \times n}(F)$. The **rational canonical form** of A is defined to be the rational canonical form of L_A . Likewise, for A , the **elementary divisors** and their **multiplicities** are the same as those of L_A .

Let A be an $n \times n$ matrix, let C be a rational canonical form of A , and let β be the appropriate rational canonical basis for L_A . Then $C = [L_A]_\beta$, and therefore A is similar to C . In fact, if Q is the matrix whose columns are the vectors of β in the same order, then $Q^{-1}AQ = C$.

Example 5

For the following real matrix A , we find the rational canonical form C of A and a matrix Q such that $Q^{-1}AQ = C$.

$$A = \begin{pmatrix} 0 & 2 & 0 & -6 & 2 \\ 1 & -2 & 0 & 0 & 2 \\ 1 & 0 & 1 & -3 & 2 \\ 1 & -2 & 1 & -1 & 2 \\ 1 & -4 & 3 & -3 & 4 \end{pmatrix}$$

The characteristic polynomial of A is $f(t) = -(t^2 + 2)^2(t - 2)$; therefore $\phi_1(t) = t^2 + 2$ and $\phi_2(t) = t - 2$ are the distinct irreducible monic divisors of $f(t)$. By Theorem 7.23, $\dim(K_{\phi_1}) = 4$ and $\dim(K_{\phi_2}) = 1$. Since the degree of $\phi_1(t)$ is 2, the total number of dots in the dot diagram of $\phi_1(t)$ is $4/2 = 2$, and the number of dots r_1 in the first row is given by

$$\begin{aligned} r_1 &= \frac{1}{2}[\dim(R^5) - \text{rank}(\phi_1(A))] \\ &= \frac{1}{2}[5 - \text{rank}(A^2 + 2I)] \\ &= \frac{1}{2}[5 - 1] = 2. \end{aligned}$$

Thus the dot diagram of $\phi_1(t)$ is

• •

and each column contributes the companion matrix

$$\begin{pmatrix} 0 & -2 \\ 1 & 0 \end{pmatrix}$$

for $\phi_1(t) = t^2 + 2$ to the rational canonical form C . Consequently $\phi_1(t)$ is an elementary divisor with multiplicity 2. Since $\dim(K_{\phi_2}) = 1$, the dot diagram of $\phi_2(t) = t - 2$ consists of a single dot, which contributes the 1×1 matrix

(2). Hence $\phi_2(t)$ is an elementary divisor with multiplicity 1. Therefore the rational canonical form C is

$$C = \left(\begin{array}{cc|ccc} 0 & -2 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & -2 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 2 \end{array} \right).$$

We can infer from the dot diagram of $\phi_1(t)$ that if β is a rational canonical basis for L_A , then $\beta \cap K_{\phi_1}$ is the union of two cyclic bases β_{v_1} and β_{v_2} , where v_1 and v_2 each have annihilator $\phi_1(t)$. It follows that both v_1 and v_2 lie in $N(\phi_1(L_A))$. It can be shown that

$$\left\{ \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ 2 \\ 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ -1 \\ 0 \\ 1 \end{pmatrix} \right\}$$

is a basis for $N(\phi_1(L_A))$. Setting $v_1 = e_1$, we see that

$$Av_1 = \begin{pmatrix} 0 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}.$$

Next choose v_2 in $K_{\phi_1} = N(\phi(L_A))$, but not in the span of $\beta_{v_1} = \{v_1, Av_1\}$. For example, $v_2 = e_2$. Then it can be seen that

$$Av_2 = \begin{pmatrix} 2 \\ -2 \\ 0 \\ -2 \\ -4 \end{pmatrix},$$

and $\beta_{v_1} \cup \beta_{v_2}$ is a basis for K_{ϕ_1} .

Since the dot diagram of $\phi_2(t) = t - 2$ consists of a single dot, any nonzero vector in K_{ϕ_2} is an eigenvector of A corresponding to the eigenvalue $\lambda = 2$. For example, choose

$$v_3 = \begin{pmatrix} 0 \\ 1 \\ 1 \\ 1 \\ 2 \end{pmatrix}.$$

By Theorem 7.23, $\beta = \{v_1, Av_1, v_2, Av_2, v_3\}$ is a rational canonical basis for \mathbb{L}_A . So setting

$$Q = \begin{pmatrix} 1 & 0 & 0 & 2 & 0 \\ 0 & 1 & 1 & -2 & 1 \\ 0 & 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & -2 & 1 \\ 0 & 1 & 0 & -4 & 2 \end{pmatrix},$$

we have $Q^{-1}AQ = C$. \blacklozenge

Example 6

For the following matrix A , we find the rational canonical form C and a matrix Q such that $Q^{-1}AQ = C$.

$$A = \begin{pmatrix} 2 & 1 & 0 & 0 \\ 0 & 2 & 1 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 2 \end{pmatrix}$$

Since the characteristic polynomial of A is $f(t) = (t-2)^4$, the only irreducible monic divisor of $f(t)$ is $\phi(t) = t-2$, and so $K_\phi = \mathbb{R}^4$. In this case, $\phi(t)$ has degree 1; hence in applying Theorem 7.24 to compute the dot diagram for $\phi(t)$, we obtain

$$\begin{aligned} r_1 &= 4 - \text{rank}(\phi(A)) = 4 - 2 = 2, \\ r_2 &= \text{rank}(\phi(A)) - \text{rank}((\phi(A))^2) = 2 - 1 = 1, \end{aligned}$$

and

$$r_3 = \text{rank}((\phi(A))^2) - \text{rank}((\phi(A))^3) = 1 - 0 = 1,$$

where r_i is the number of dots in the i th row of the dot diagram. Since there are $\dim(\mathbb{R}^4) = 4$ dots in the diagram, we may terminate these computations with r_3 . Thus the dot diagram for A is

$$\begin{array}{c} \bullet \quad \bullet \\ \vdots \\ \bullet \\ \vdots \end{array}$$

Since $(t-2)^3$ has the companion matrix

$$\begin{pmatrix} 0 & 0 & 8 \\ 1 & 0 & -12 \\ 0 & 1 & 6 \end{pmatrix}$$

and $(t - 2)$ has the companion matrix (2), the rational canonical form of A is given by

$$C = \left(\begin{array}{ccc|c} 0 & 0 & 8 & 0 \\ 1 & 0 & -12 & 0 \\ 0 & 1 & 6 & 0 \\ \hline 0 & 0 & 0 & 2 \end{array} \right).$$

Next we find a rational canonical basis for L_A . The preceding dot diagram indicates that there are two vectors v_1 and v_2 in \mathbb{R}^4 with annihilators $(\phi(t))^3$ and $\phi(t)$, respectively, and such that

$$\beta = \beta_{v_1} \cup \beta_{v_2} = \{v_1, Av_1, A^2v_1, v_2\}$$

is a rational canonical basis for L_A . Furthermore, $v_1 \notin N((L_A - 2I)^2)$, and $v_2 \in N(L_A - 2I)$. It can easily be shown that

$$N(L_A - 2I) = \text{span}(\{e_1, e_4\})$$

and

$$N((L_A - 2I)^2) = \text{span}(\{e_1, e_2, e_4\}).$$

The standard vector e_3 meets the criteria for v_1 ; so we set $v_1 = e_3$. It follows that

$$Av_1 = \begin{pmatrix} 0 \\ 1 \\ 2 \\ 0 \end{pmatrix} \quad \text{and} \quad A^2v_1 = \begin{pmatrix} 1 \\ 4 \\ 4 \\ 0 \end{pmatrix}.$$

Next we choose a vector $v_2 \in N(L_A - 2I)$ that is not in the span of β_{v_1} . Clearly, $v_2 = e_4$ satisfies this condition. Thus

$$\left\{ \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 2 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 4 \\ 4 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \end{pmatrix} \right\}$$

is a rational canonical basis for L_A .

Finally, let Q be the matrix whose columns are the vectors of β in the same order:

$$Q = \begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 1 & 4 & 0 \\ 1 & 2 & 4 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}.$$

Then $C = Q^{-1}AQ$. ◆

Direct Sums*

The next theorem is a simple consequence of Theorem 7.23.

Theorem 7.25 (Primary Decomposition Theorem). Let T be a linear operator on an n -dimensional vector space V with characteristic polynomial

$$f(t) = (-1)^n (\phi_1(t))^{n_1} (\phi_2(t))^{n_2} \cdots (\phi_k(t))^{n_k},$$

where the $\phi_i(t)$'s ($1 \leq i \leq k$) are distinct irreducible monic polynomials and the n_i 's are positive integers. Then the following statements are true.

- (a) $V = K_{\phi_1} \oplus K_{\phi_2} \oplus \cdots \oplus K_{\phi_k}$.
- (b) If T_i ($1 \leq i \leq k$) is the restriction of T to K_{ϕ_i} and C_i is the rational canonical form of T_i , then $C_1 \oplus C_2 \oplus \cdots \oplus C_k$ is the rational canonical form of T .

Proof. Exercise. ■

The next theorem is a simple consequence of Theorem 7.17.

Theorem 7.26. Let T be a linear operator on a finite-dimensional vector space V . Then V is a direct sum of T -cyclic subspaces C_{v_i} , where each v_i lies in K_ϕ for some irreducible monic divisor $\phi(t)$ of the characteristic polynomial of T .

Proof. Exercise. ■

EXERCISES

1. Label the following statements as true or false.

- (a) Every rational canonical basis for a linear operator T is the union of T -cyclic bases.
- (b) If a basis is the union of T -cyclic bases for a linear operator T , then it is a rational canonical basis for T .
- (c) There exist square matrices having no rational canonical form.
- (d) A square matrix is similar to its rational canonical form.
- (e) For any linear operator T on a finite-dimensional vector space, any irreducible factor of the characteristic polynomial of T divides the minimal polynomial of T .
- (f) Let $\phi(t)$ be an irreducible monic divisor of the characteristic polynomial of a linear operator T . The dots in the diagram used to compute the rational canonical form of the restriction of T to K_ϕ are in one-to-one correspondence with the vectors in a basis for K_ϕ .

- (g) If a matrix has a Jordan canonical form, then its Jordan canonical form and rational canonical form are similar.
2. For each of the following matrices $A \in M_{n \times n}(F)$, find the rational canonical form C of A and a matrix $Q \in M_{n \times n}(F)$ such that $Q^{-1}AQ = C$.
- (a) $A = \begin{pmatrix} 3 & 1 & 0 \\ 0 & 3 & 1 \\ 0 & 0 & 3 \end{pmatrix}$ $F = R$
- (b) $A = \begin{pmatrix} 0 & -1 \\ 1 & -1 \end{pmatrix}$ $F = R$
- (c) $A = \begin{pmatrix} 0 & -1 \\ 1 & -1 \end{pmatrix}$ $F = C$
- (d) $A = \begin{pmatrix} 0 & -7 & 14 & -6 \\ 1 & -4 & 6 & -3 \\ 0 & -4 & 9 & -4 \\ 0 & -4 & 11 & -5 \end{pmatrix}$ $F = R$
- (e) $A = \begin{pmatrix} 0 & -4 & 12 & -7 \\ 1 & -1 & 3 & -3 \\ 0 & -1 & 6 & -4 \\ 0 & -1 & 8 & -5 \end{pmatrix}$ $F = R$
3. For each of the following linear operators T , find the elementary divisors, the rational canonical form C , and a rational canonical basis β .
- (a) T is the linear operator on $P_3(R)$ defined by
- $$T(f(x)) = f(0)x - f'(1).$$
- (b) Let $S = \{\sin x, \cos x, x \sin x, x \cos x\}$, a subset of $\mathcal{F}(R, R)$, and let $V = \text{span}(S)$. Define T to be the linear operator on V such that
- $$T(f) = f'.$$
- (c) T is the linear operator on $M_{2 \times 2}(R)$ defined by
- $$T(A) = \begin{pmatrix} 0 & 1 \\ -1 & 1 \end{pmatrix} \cdot A.$$
- (d) Let $S = \{\sin x \sin y, \sin x \cos y, \cos x \sin y, \cos x \cos y\}$, a subset of $\mathcal{F}(R \times R, R)$, and let $V = \text{span}(S)$. Define T to be the linear operator on V such that
- $$T(f)(x, y) = \frac{\partial f(x, y)}{\partial x} + \frac{\partial f(x, y)}{\partial y}.$$
4. Let T be a linear operator on a finite-dimensional vector space V with minimal polynomial $(\phi(t))^m$ for some positive integer m .
- (a) Prove that $R(\phi(T)) \subseteq N((\phi(T))^{m-1})$.

- (b) Give an example to show that the subspaces in (a) need not be equal.
- (c) Prove that the minimal polynomial of the restriction of T to $R(\phi(T))$ equals $(\phi(t))^{m-1}$.
- 5.** Let T be a linear operator on a finite-dimensional vector space. Prove that the rational canonical form of T is a diagonal matrix if and only if T is diagonalizable. Visit [goo.gl/tK8pru](#) for a solution.
- 6.** Let T be a linear operator on a finite-dimensional vector space V with characteristic polynomial $f(t) = (-1)^n \phi_1(t)\phi_2(t)$, where $\phi_1(t)$ and $\phi_2(t)$ are distinct irreducible monic polynomials and $n = \dim(V)$.
- Prove that there exist $v_1, v_2 \in V$ such that v_1 has T -annihilator $\phi_1(t)$, v_2 has T -annihilator $\phi_2(t)$, and $\beta_{v_1} \cup \beta_{v_2}$ is a basis for V .
 - Prove that there is a vector $v_3 \in V$ with T -annihilator $\phi_1(t)\phi_2(t)$ such that β_{v_3} is a basis for V .
 - Describe the difference between the matrix representation of T with respect to $\beta_{v_1} \cup \beta_{v_2}$ and the matrix representation of T with respect to β_{v_3} .

Thus, to assure the uniqueness of the rational canonical form, we require that the generators of the T -cyclic bases that constitute a rational canonical basis have T -annihilators equal to powers of irreducible monic factors of the characteristic polynomial of T .

- 7.** Let T be a linear operator on a finite-dimensional vector space with minimal polynomial

$$f(t) = (\phi_1(t))^{m_1}(\phi_2(t))^{m_2} \cdots (\phi_k(t))^{m_k},$$

where the $\phi_i(t)$'s are distinct irreducible monic factors of $f(t)$. Prove that for each i , m_i is the number of entries in the first column of the dot diagram for $\phi_i(t)$.

- 8.** Let T be a linear operator on a finite-dimensional vector space V . Prove that for any irreducible polynomial $\phi(t)$, if $\phi(T)$ is not one-to-one, then $\phi(t)$ divides the characteristic polynomial of T . *Hint:* Apply Exercise 15 of Section 7.3.
- 9.** Let V be a vector space and $\beta_1, \beta_2, \dots, \beta_k$ be disjoint subsets of V whose union is a basis for V . Now suppose that $\gamma_1, \gamma_2, \dots, \gamma_k$ are linearly independent subsets of V such that $\text{span}(\gamma_i) = \text{span}(\beta_i)$ for all i . Prove that $\gamma_1 \cup \gamma_2 \cup \cdots \cup \gamma_k$ is also a basis for V .
- 10.** Let T be a linear operator on a finite-dimensional vector space, and suppose that $\phi(t)$ is an irreducible monic factor of the characteristic polynomial of T . Prove that if $\phi(t)$ is the T -annihilator of vectors x and y , then $x \in C_y$ if and only if $C_x = C_y$.

Exercises 11 and 12 are concerned with direct sums.

11. Prove Theorem 7.25.

12. Prove Theorem 7.26.

INDEX OF DEFINITIONS FOR CHAPTER 7

- | | |
|---|---|
| Companion matrix 519 | Jordan canonical form of a linear operator 476 |
| Cycle of generalized eigenvectors 482 | Jordan canonical form of a matrix 485 |
| Cyclic basis 519 | Length of a cycle 482 |
| Dot diagram for Jordan canonical form 491 | Minimal polynomial of a linear operator 509 |
| Dot diagram for rational canonical form 528 | Minimal polynomial of a matrix 510 |
| Elementary divisor of a linear operator 532 | Multiplicity of an elementary divisor 532 |
| Elementary divisor of a matrix 534 | Rational canonical basis of a linear operator 519 |
| End vector of a cycle 482 | Rational canonical form for a linear operator 519 |
| Generalized eigenspace 477 | Rational canonical form of a matrix 534 |
| Generalized eigenvector 477 | |
| Generator of a cyclic basis 519 | |
| Initial vector of a cycle 482 | |
| Jordan block 476 | |
| Jordan canonical basis 476 | |

Appendices

APPENDIX A SETS

A **set**¹ is a collection of objects, called **elements** of the set. If x is an element of the set A , then we write $x \in A$; otherwise, we write $x \notin A$. For example, if Z is the set of integers, then $3 \in Z$ and $\frac{1}{2} \notin Z$.

One set that appears frequently is the set of real numbers, which we denote by R throughout this text.

Two sets A and B are called **equal**, written $A = B$, if they contain exactly the same elements. Sets may be described in one of two ways:

1. By listing the elements of the set between set braces { }.
2. By describing the elements of the set in terms of some characteristic property.

For example, the set consisting of the elements 1, 2, 3, and 4 can be written as $\{1, 2, 3, 4\}$ or as

$$\{x: x \text{ is a positive integer less than } 5\}.$$

Note that the order in which the elements of a set are listed is immaterial, and listing an element in a set more than once does not change the set. Hence

$$\{1, 2, 3, 4\} = \{3, 1, 2, 4\} = \{1, 3, 1, 4, 2\}.$$

Example 1

Let A denote the set of real numbers strictly between 1 and 2. Then A may be written as

$$A = \{x \in R: 1 < x < 2\}. \quad \blacklozenge$$

A set B is called a **subset** of a set A , written $B \subseteq A$ or $A \supseteq B$, if every element of B is an element of A . For example, $\{1, 2, 6\} \subseteq \{2, 8, 7, 6, 1\}$. If $B \subseteq A$, and $B \neq A$, then B is called a **proper subset** of A . Observe that $A = B$ if and only if $A \subseteq B$ and $B \subseteq A$, a fact that is often used to prove that two sets are equal.

The **empty set**, denoted by \emptyset , is the set containing no elements. The empty set is a subset of every set.

¹It is impossible to give a formal definition of the word *set* in terms of simpler concepts. In this appendix we explain the conceptual ideas behind sets and the facts about them that we will use in this book.

Sets may be combined to form other sets in two basic ways. The **union** of two sets A and B , denoted $A \cup B$, is the set of elements that are in A , or B , or both; that is,

$$A \cup B = \{x: x \in A \text{ or } x \in B\}.$$

The **intersection** of two sets A and B , denoted $A \cap B$, is the set of elements that are in both A and B ; that is,

$$A \cap B = \{x: x \in A \text{ and } x \in B\}.$$

Two sets are called **disjoint** if their intersection equals the empty set.

Example 2

Let $A = \{1, 3, 5\}$ and $B = \{1, 5, 7, 8\}$. Then

$$A \cup B = \{1, 3, 5, 7, 8\} \quad \text{and} \quad A \cap B = \{1, 5\}.$$

Likewise, if $X = \{1, 2, 8\}$ and $Y = \{3, 4, 5\}$, then

$$X \cup Y = \{1, 2, 3, 4, 5, 8\} \quad \text{and} \quad X \cap Y = \emptyset.$$

Thus X and Y are disjoint sets. \blacklozenge

The union and intersection of more than two sets can be defined analogously. Specifically, if A_1, A_2, \dots, A_n are sets, then the union and intersections of these sets are defined, respectively, by

$$\bigcup_{i=1}^n A_i = \{x: x \in A_i \text{ for some } i = 1, 2, \dots, n\}$$

and

$$\bigcap_{i=1}^n A_i = \{x: x \in A_i \text{ for all } i = 1, 2, \dots, n\}.$$

In this situation, we have one set A_i for each $i \in \{1, 2, \dots, n\}$. It can be convenient to regard these sets A_i as being *indexed* by the set $I = \{1, 2, \dots, n\}$. When doing so, I is called an **index set** for the collection $\{A_i : i \in I\}$.

The use of an index set is especially useful with an infinite collection of sets, as in Example 3 below. If Λ is an index set and $\{A_\alpha : \alpha \in \Lambda\}$ is a collection of sets, the union and intersection of these sets are defined, respectively, by

$$\bigcup_{\alpha \in \Lambda} A_\alpha = \{x: x \in A_\alpha \text{ for some } \alpha \in \Lambda\}$$

and

$$\bigcap_{\alpha \in \Lambda} A_\alpha = \{x: x \in A_\alpha \text{ for all } \alpha \in \Lambda\}.$$

Example 3

Let $\Lambda = \{\alpha \in R: \alpha > 1\}$, and let

$$A_\alpha = \left\{ x \in R: \frac{-1}{\alpha} \leq x \leq 1 + \alpha \right\}$$

for each $\alpha \in \Lambda$. Then

$$\bigcup_{\alpha \in \Lambda} A_\alpha = \{x \in R: x > -1\} \quad \text{and} \quad \bigcap_{\alpha \in \Lambda} A_\alpha = \{x \in R: 0 \leq x \leq 2\}. \quad \blacklozenge$$

By a relation on a set A , we mean a rule for determining whether or not, for any elements x and y in A , x stands in a given relationship to y . To make this concept precise, we define a **relation** on A to be any set S of ordered pairs of elements of A , and say that the elements x and y of A satisfy the relation S if and only if $(x, y) \in S$. On the set of real numbers, for instance, “is equal to,” “is less than,” and “is greater than or equal to” are familiar relations. If S is a relation on a set A , we often write $x \sim y$ in place of $(x, y) \in S$.

A relation S on a set A is called an **equivalence relation** on A if the following three conditions hold:

1. For each $x \in A$, $x \sim x$ (reflexivity).
2. If $x \sim y$, then $y \sim x$ (symmetry).
3. If $x \sim y$ and $y \sim z$, then $x \sim z$ (transitivity).

For example, if we define $x \sim y$ to mean that $x - y$ is divisible by a fixed integer n , then \sim is an equivalence relation on the set of integers.

APPENDIX B FUNCTIONS

If A and B are sets, then a **function** f from A to B , written $f: A \rightarrow B$, is a rule that associates to each element x in A a unique element denoted $f(x)$ in B . The element $f(x)$ is called the **image** of x (under f), and x is called a **preimage** of $f(x)$ (under f). If $f: A \rightarrow B$, then A is called the **domain** of f , B is called the **codomain** of f , and the set $\{f(x): x \in A\}$ is called the **range** of f . Note that the range of f is a subset of B . If $S \subseteq A$, we denote by $f(S)$ the set $\{f(x): x \in S\}$ of all images of elements of S . Likewise, if $T \subseteq B$, we denote by $f^{-1}(T)$ the set $\{x \in A: f(x) \in T\}$ of all preimages of elements in T . Finally, two functions $f: A \rightarrow B$ and $g: A \rightarrow B$ are **equal**, written $f = g$, if $f(x) = g(x)$ for all $x \in A$.

Example 1

Suppose that $A = [-10, 10]$. Let $f: A \rightarrow R$ be the function that assigns to each element x in A the element $x^2 + 1$ in R ; that is, f is defined by

$f(x) = x^2 + 1$. Then A is the domain of f , R is the codomain of f , and $[1, 101]$ is the range of f . Since $f(2) = 5$, the image of 2 is 5, and 2 is a preimage of 5. Notice that -2 is another preimage of 5. Moreover, if $S = [1, 2]$ and $T = [82, 101]$, then $f(S) = [2, 5]$ and $f^{-1}(T) = [-10, -9] \cup [9, 10]$. ♦

As Example 1 shows, the preimage of an element in the range need not be unique. Functions such that each element of the range has a unique preimage are called **one-to-one**; that is $f: A \rightarrow B$ is one-to-one if $f(x) = f(y)$ implies $x = y$ or, equivalently, if $x \neq y$ implies $f(x) \neq f(y)$.

If $f: A \rightarrow B$ is a function with range B , that is, if $f(A) = B$, then f is called **onto**. So f is onto if and only if the range of f equals the codomain of f .

Let $f: A \rightarrow B$ be a function and $S \subseteq A$. Then a function $f_S: S \rightarrow B$, called the **restriction** of f to S , can be formed by defining $f_S(x) = f(x)$ for each $x \in S$.

The next example illustrates these concepts.

Example 2

Let $f: [-1, 1] \rightarrow [0, 1]$ be defined by $f(x) = x^2$. This function is onto, but not one-to-one since $f(-1) = f(1) = 1$. Note that if $S = [0, 1]$, then f_S is both onto and one-to-one. Finally, if $T = [\frac{1}{2}, 1]$, then f_T is one-to-one, but not onto. ♦

Let A , B , and C be sets and $f: A \rightarrow B$ and $g: B \rightarrow C$ be functions. By following f with g , we obtain a function $g \circ f: A \rightarrow C$ called the **composite** of g and f . Thus $(g \circ f)(x) = g(f(x))$ for all $x \in A$. For example, let $A = B = C = R$, $f(x) = \sin x$, and $g(x) = x^2 + 3$. Then $(g \circ f)(x) = (g(f(x))) = \sin^2 x + 3$, whereas $(f \circ g)(x) = f(g(x)) = \sin(x^2 + 3)$. Hence, $g \circ f \neq f \circ g$. Functional composition is associative, however; that is, if $h: C \rightarrow D$ is another function, then $h \circ (g \circ f) = (h \circ g) \circ f$.

A function $f: A \rightarrow B$ is said to be **invertible** if there exists a function $g: B \rightarrow A$ such that $(f \circ g)(y) = y$ for all $y \in B$ and $(g \circ f)(x) = x$ for all $x \in A$. If such a function g exists, then it is unique and is called the **inverse** of f . We denote the inverse of f (when it exists) by f^{-1} . It can be shown that f is invertible if and only if f is both one-to-one and onto.

Example 3

The function $f: R \rightarrow R$ defined by $f(x) = 3x + 1$ is one-to-one and onto; hence f is invertible. The inverse of f is the function $f^{-1}: R \rightarrow R$ defined by $f^{-1}(x) = (x - 1)/3$. ♦

The following facts about invertible functions are easily proved.

1. If $f: A \rightarrow B$ is invertible, then f^{-1} is invertible, and $(f^{-1})^{-1} = f$.

2. If $f: A \rightarrow B$ and $g: B \rightarrow C$ are invertible, then $g \circ f$ is invertible, and $(g \circ f)^{-1} = f^{-1} \circ g^{-1}$.

APPENDIX C FIELDS

The set of real numbers is an example of an algebraic structure called a **field**. Basically, a field is a set in which four operations (called addition, multiplication, subtraction, and division) can be defined so that, with the exception of division by zero, the sum, product, difference, and quotient of any two elements in the set is an element of the set. More precisely, a field is defined as follows.

Definitions. A field F is a set on which two operations $+$ and \cdot (called **addition** and **multiplication**, respectively) are defined so that, for each pair of elements x, y in F , there are unique elements in F , denoted $x + y$ and $x \cdot y$, and such that the following conditions hold for all elements a, b, c in F .

- (F 1) $a + b = b + a$ and $a \cdot b = b \cdot a$
(commutativity of addition and multiplication)
- (F 2) $(a + b) + c = a + (b + c)$ and $(a \cdot b) \cdot c = a \cdot (b \cdot c)$
(associativity of addition and multiplication)
- (F 3) There exist distinct elements 0 and 1 in F such that

$$0 + a = a \quad \text{and} \quad 1 \cdot a = a$$

(existence of identity elements for addition and multiplication)

- (F 4) For each element a in F and each nonzero element b in F , there exist elements c and d in F such that

$$a + c = 0 \quad \text{and} \quad b \cdot d = 1$$

(existence of inverses for addition and multiplication)

- (F 5) $a \cdot (b + c) = a \cdot b + a \cdot c$
(distributivity of multiplication over addition)

The elements $x + y$ and $x \cdot y$ are called the **sum** and **product**, respectively, of x and y . The elements 0 (read “zero”) and 1 (read “one”) mentioned in (F 3) are called **identity elements** for addition and multiplication, respectively, and the elements c and d referred to in (F 4) are called an **additive inverse** for a and a **multiplicative inverse** for b , respectively.

Example 1

The set of real numbers R with the usual definitions of addition and multiplication is a field. ♦

Example 2

The set of rational numbers with the usual definitions of addition and multiplication is a field. ♦

Example 3

The set of all real numbers of the form $a + b\sqrt{2}$, where a and b are rational numbers, with addition and multiplication as in \mathbb{R} is a field. ♦

Example 4

The field Z_2 consists of two elements 0 and 1 with the operations of addition and multiplication defined by the equations

$$\begin{aligned}0 + 0 &= 0, & 0 + 1 &= 1 + 0 = 1, & 1 + 1 &= 0, \\0 \cdot 0 &= 0, & 0 \cdot 1 &= 1 \cdot 0 = 0, & \text{and } 1 \cdot 1 &= 1.\end{aligned}\quad \blacklozenge$$

Example 5

Neither the set of positive integers nor the set of integers with the usual definitions of addition and multiplication is a field, for in either case (F 4) does not hold. ♦

The identity and inverse elements guaranteed by (F 3) and (F 4) are unique; this is a consequence of the following theorem.

Theorem C.1 (Cancellation Laws). *For arbitrary elements a , b , and c in a field, the following statements are true.*

- (a) *If $a + b = c + b$, then $a = c$.*
- (b) *If $a \cdot b = c \cdot b$ and $b \neq 0$, then $a = c$.*

Proof. (a) The proof of (a) is left as an exercise.

(b) If $b \neq 0$, then (F 4) guarantees the existence of an element d in the field such that $b \cdot d = 1$. Multiply both sides of the equality $a \cdot b = c \cdot b$ by d to obtain $(a \cdot b) \cdot d = (c \cdot b) \cdot d$. Consider the left side of this equality: By (F 2) and (F 3), we have

$$(a \cdot b) \cdot d = a \cdot (b \cdot d) = a \cdot 1 = a.$$

Similarly, the right side of the equality reduces to c . Thus $a = c$. ■

Corollary. *The elements 0 and 1 mentioned in (F 3), and the elements c and d mentioned in (F 4), are unique.*

Proof. Suppose that $0' \in F$ satisfies $0' + a = a$ for each $a \in F$. Since $0 + a = a$ for each $a \in F$, we have $0' + a = 0 + a$ for each $a \in F$. Thus $0' = 0$ by Theorem C.1.

The proofs of the remaining parts are similar. ■

Thus each element b in a field has a unique additive inverse and, if $b \neq 0$, a unique multiplicative inverse. (It is shown in the corollary to Theorem C.2 that 0 has no multiplicative inverse.) The additive inverse and the multiplicative inverse of b are denoted by $-b$ and b^{-1} , respectively. Note that $-(-b) = b$ and $(b^{-1})^{-1} = b$.

Subtraction and *division* can be defined in terms of addition and multiplication by using the additive and multiplicative inverses. Specifically, subtraction of b is defined to be addition of $-b$ and division by $b \neq 0$ is defined to be multiplication by b^{-1} ; that is,

$$a - b = a + (-b) \quad \text{and} \quad \frac{a}{b} = a \cdot b^{-1}.$$

In particular, the symbol $\frac{1}{b}$ denotes b^{-1} . Division by zero is undefined, but, with this exception, the sum, product, difference, and quotient of any two elements of a field are defined.

Many of the familiar properties of multiplication of real numbers are true in any field, as the next theorem shows.

Theorem C.2. *Let a and b be arbitrary elements of a field. Then each of the following statements is true.*

- (a) $a \cdot 0 = 0$.
- (b) $(-a) \cdot b = a \cdot (-b) = -(a \cdot b)$.
- (c) $(-a) \cdot (-b) = a \cdot b$.

Proof. (a) Since $0 + 0 = 0$, (F 5) shows that

$$0 + a \cdot 0 = a \cdot 0 = a \cdot (0 + 0) = a \cdot 0 + a \cdot 0.$$

Thus $0 = a \cdot 0$ by Theorem C.1.

(b) By definition, $-(a \cdot b)$ is the unique element of F with the property $a \cdot b + [-(a \cdot b)] = 0$. So in order to prove that $(-a) \cdot b = -(a \cdot b)$, it suffices to show that $a \cdot b + (-a) \cdot b = 0$. But $-a$ is the element of F such that $a + (-a) = 0$; so

$$a \cdot b + (-a) \cdot b = [a + (-a)] \cdot b = 0 \cdot b = b \cdot 0 = 0$$

by (F 5) and (a). Thus $(-a) \cdot b = -(a \cdot b)$. The proof that $a \cdot (-b) = -(a \cdot b)$ is similar.

(c) By applying (b) twice, we find that

$$(-a) \cdot (-b) = -[a \cdot (-b)] = -[-(a \cdot b)] = a \cdot b. \quad \blacksquare$$

Corollary. *The additive identity of a field has no multiplicative inverse.*

In an arbitrary field F , it may happen that a sum $1 + 1 + \cdots + 1$ (p summands) equals 0 for some positive integer p . For example, in the field Z_2 (defined in Example 4), $1 + 1 = 0$. In this case, the smallest positive integer p for which a sum of p 1's equals 0 is called the **characteristic** of F ; if no such positive integer exists, then F is said to have **characteristic zero**. Thus Z_2 has characteristic two, and R has characteristic zero. Observe that if F is a field of characteristic $p \neq 0$, then $x + x + \cdots + x$ (p summands) equals 0 for all $x \in F$. In a field having nonzero characteristic (especially characteristic two), many unexpected problems arise. For this reason, some of the results about vector spaces stated in this book require that the field over which the vector space is defined be of characteristic zero (or, at least, of some characteristic other than two).

Finally, note that in other sections of this book, the product of two elements a and b in a field is usually denoted ab rather than $a \cdot b$.

APPENDIX D COMPLEX NUMBERS

For the purposes of algebra, the field of real numbers is not sufficient, for there are polynomials of nonzero degree with real coefficients that have no zeros in the field of real numbers (for example, $x^2 + 1$). It is often desirable to have a field in which any polynomial of nonzero degree with coefficients from that field has a zero in that field. It is possible to “enlarge” the field of real numbers to obtain such a field.

Definitions. A **complex number** is an expression of the form $z = a + bi$, where a and b are real numbers called the **real part** and the **imaginary part** of z , respectively.

The **sum** and **product** of two complex numbers $z = a + bi$ and $w = c + di$ (where a, b, c , and d are real numbers) are defined, respectively, as follows:

$$z + w = (a + bi) + (c + di) = (a + c) + (b + d)i$$

and

$$zw = (a + bi)(c + di) = (ac - bd) + (bc + ad)i.$$

Example 1

The sum and product of $z = 3 - 5i$ and $w = 9 + 7i$ are, respectively,

$$z + w = (3 - 5i) + (9 + 7i) = (3 + 9) + [(-5) + 7]i = 12 + 2i$$

and

$$zw = (3 - 5i)(9 + 7i) = [3 \cdot 9 - (-5) \cdot 7] + [(-5) \cdot 9 + 3 \cdot 7]i = 62 - 24i. \quad \blacklozenge$$

Any real number c may be regarded as a complex number by identifying c with the complex number $c + 0i$, that is, treating them as identical. Observe that this correspondence preserves sums and products; that is,

$$(c + 0i) + (d + 0i) = (c + d) + 0i \quad \text{and} \quad (c + 0i)(d + 0i) = cd + 0i.$$

Any complex number of the form $bi = 0 + bi$, where b is a nonzero real number, is called **imaginary**. The product of two imaginary numbers is real since

$$(bi)(di) = (0 + bi)(0 + di) = (0 - bd) + (b \cdot 0 + 0 \cdot d)i = -bd.$$

In particular, for $i = 0 + 1i$, we have $i \cdot i = -1$.

The observation that $i^2 = i \cdot i = -1$ provides an easy way to remember the definition of multiplication of complex numbers: simply multiply two complex numbers as you would any two algebraic expressions, and replace i^2 by -1 . Example 2 illustrates this technique.

Example 2

The product of $-5 + 2i$ and $1 - 3i$ is

$$\begin{aligned} (-5 + 2i)(1 - 3i) &= -5(1 - 3i) + 2i(1 - 3i) \\ &= -5 + 15i + 2i - 6i^2 \\ &= -5 + 15i + 2i - 6(-1) \\ &= 1 + 17i. \quad \blacklozenge \end{aligned}$$

The real number 0, regarded as a complex number, is an additive identity element for the complex numbers since

$$(a + bi) + 0 = (a + bi) + (0 + 0i) = (a + 0) + (b + 0)i = a + bi.$$

Likewise the real number 1, regarded as a complex number, is a multiplicative identity element for the set of complex numbers since

$$(a + bi) \cdot 1 = (a + bi)(1 + 0i) = (a \cdot 1 - b \cdot 0) + (b \cdot 1 + a \cdot 0)i = a + bi.$$

Every complex number $a + bi$ has an additive inverse, namely $(-a) + (-b)i$. But also each complex number except 0 has a multiplicative inverse. In fact,

$$(a + bi)^{-1} = \left(\frac{a}{a^2 + b^2} \right) - \left(\frac{b}{a^2 + b^2} \right) i.$$

In view of the preceding statements, the following result is not surprising.

Theorem D.1. *The set of complex numbers with the operations of addition and multiplication previously defined is a field.*

Proof. Exercise. ■

Definition. The (**complex**) **conjugate** of a complex number $a + bi$ is the complex number $a - bi$. We denote the conjugate of the complex number z by \bar{z} .

Example 3

The conjugates of $-3 + 2i$, $4 - 7i$, and 6 are, respectively,

$$\overline{-3 + 2i} = -3 - 2i, \quad \overline{4 - 7i} = 4 + 7i, \quad \text{and} \quad \overline{6} = \overline{6 + 0i} = 6 - 0i = 6. \quad \blacklozenge$$

The next theorem contains some important properties of the conjugate of a complex number.

Theorem D.2. Let z and w be complex numbers. Then the following statements are true.

- (a) $\overline{\bar{z}} = z$.
- (b) $\overline{(z + w)} = \bar{z} + \bar{w}$.
- (c) $\overline{zw} = \bar{z} \cdot \bar{w}$.
- (d) $\overline{\left(\frac{z}{w}\right)} = \frac{\bar{z}}{\bar{w}}$ if $w \neq 0$.
- (e) z is a real number if and only if $\bar{z} = z$.

Proof. We leave the proofs of (a), (d), and (e) to the reader.

(b) Let $z = a + bi$ and $w = c + di$, where $a, b, c, d \in R$. Then

$$\begin{aligned} \overline{(z + w)} &= \overline{(a + c) + (b + d)i} = (a + c) - (b + d)i \\ &= (a - bi) + (c - di) = \bar{z} + \bar{w}. \end{aligned}$$

(c) For z and w , we have

$$\begin{aligned} \overline{zw} &= \overline{(a + bi)(c + di)} = \overline{(ac - bd) + (ad + bc)i} \\ &= (ac - bd) - (ad + bc)i = (a - bi)(c - di) = \bar{z} \cdot \bar{w}. \end{aligned} \quad \blacksquare$$

For any complex number $z = a + bi$, $z\bar{z}$ is real and nonnegative, for

$$z\bar{z} = (a + bi)(a - bi) = a^2 + b^2.$$

This fact can be used to define the absolute value of a complex number.

Definition. Let $z = a + bi$, where $a, b \in R$. The **absolute value** (or **modulus**) of z is the real number $\sqrt{a^2 + b^2}$. We denote the absolute value of z by $|z|$.

Observe that $z\bar{z} = |z|^2$. The fact that the product of a complex number and its conjugate is real provides an easy method for determining the quotient of two complex numbers; for if $c + di \neq 0$, then

$$\frac{a+bi}{c+di} = \frac{a+bi}{c+di} \cdot \frac{c-di}{c-di} = \frac{(ac+bd)+(bc-ad)i}{c^2+d^2} = \frac{ac+bd}{c^2+d^2} + \frac{bc-ad}{c^2+d^2}i.$$

Example 4

To illustrate this procedure, we compute the quotient $(1+4i)/(3-2i)$:

$$\frac{1+4i}{3-2i} = \frac{1+4i}{3-2i} \cdot \frac{3+2i}{3+2i} = \frac{-5+14i}{9+4} = -\frac{5}{13} + \frac{14}{13}i. \quad \blacklozenge$$

The absolute value of a complex number has the familiar properties of the absolute value of a real number, as the following result shows.

Theorem D.3. *Let z and w denote any two complex numbers. Then the following statements are true.*

- (a) $|zw| = |z| \cdot |w|$.
- (b) $\left| \frac{z}{w} \right| = \frac{|z|}{|w|}$ if $w \neq 0$.
- (c) $|z+w| \leq |z| + |w|$.
- (d) $|z|-|w| \leq |z+w|$.

Proof. (a) By Theorem D.2, we have

$$|zw|^2 = (zw)\overline{(zw)} = (zw)(\bar{z} \cdot \bar{w}) = (z\bar{z})(w\bar{w}) = |z|^2|w|^2,$$

proving (a).

- (b) For the proof of (b), apply (a) to the product $\left(\frac{z}{w}\right)w$.
- (c) For any complex number $x = a+bi$, where $a, b \in R$, observe that

$$x + \bar{x} = (a+bi) + (a-bi) = 2a \leq 2\sqrt{a^2+b^2} = 2|x|.$$

Thus $x + \bar{x}$ is real and satisfies the inequality $x + \bar{x} \leq 2|x|$. Taking $x = w\bar{z}$, we have, by Theorem D.2 and (a),

$$w\bar{z} + \bar{w}z \leq 2|w\bar{z}| = 2|w||\bar{z}| = 2|z||w|.$$

Using Theorem D.2 again gives

$$\begin{aligned} |z+w|^2 &= (z+w)\overline{(z+w)} = (z+w)(\bar{z}+\bar{w}) = z\bar{z} + w\bar{z} + z\bar{w} + w\bar{w} \\ &\leq |z|^2 + 2|z||w| + |w|^2 = (|z| + |w|)^2. \end{aligned}$$

By taking square roots, we obtain (c).

(d) From (a) and (c), it follows that

$$|z| = |(z + w) - w| \leq |z + w| + |-w| = |z + w| + |w|.$$

So

$$|z| - |w| \leq |z + w|,$$

proving (d). ■

It is interesting as well as useful that complex numbers have both a geometric and an algebraic representation. Suppose that $z = a + bi$, where a and b are real numbers. We may represent z as a vector in the complex plane (see Figure D.1(a)). Notice that, as in \mathbb{R}^2 , there are two axes, the **real axis** and the **imaginary axis**. The real and imaginary parts of z are the first and second coordinates, and the absolute value of z gives the length of the vector z . It is clear that addition of complex numbers may be represented as in \mathbb{R}^2 using the parallelogram law.

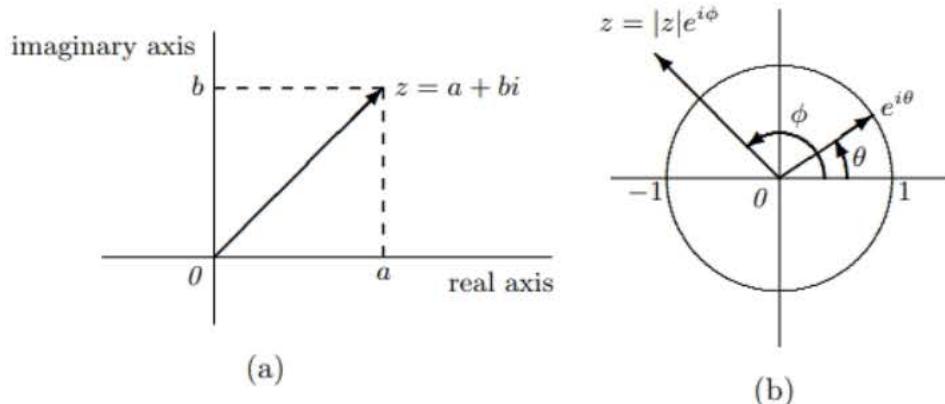


Figure D.1

In Section 2.7 (p.133), we introduce Euler's formula. The special case $e^{i\theta} = \cos \theta + i \sin \theta$ is of particular interest. Because of the geometry we have introduced, we may represent the vector $e^{i\theta}$ as in Figure D.1(b); that is, $e^{i\theta}$ is the unit vector that makes an angle θ with the positive real axis. From this figure, we see that any nonzero complex number z may be depicted as a multiple of a unit vector, namely, $z = |z|e^{i\phi}$, where ϕ is the angle that the vector z makes with the positive real axis. Thus multiplication, as well as addition, has a simple geometric interpretation: If $z = |z|e^{i\theta}$ and $w = |w|e^{i\omega}$ are two nonzero complex numbers, then from the properties established in Section 2.7 and Theorem D.3, we have

$$zw = |z|e^{i\theta} \cdot |w|e^{i\omega} = |zw|e^{i(\theta+\omega)}.$$

So zw is the vector whose length is the product of the lengths of z and w , and makes the angle $\theta + \omega$ with the positive real axis.

Our motivation for enlarging the set of real numbers to the set of complex numbers is to obtain a field such that every polynomial with nonzero degree having coefficients in that field has a zero. Our next result guarantees that the field of complex numbers has this property.

Theorem D.4 (The Fundamental Theorem of Algebra). Suppose that $p(z) = a_n z^n + a_{n-1} z^{n-1} + \cdots + a_1 z + a_0$ is a polynomial in $\mathbb{P}(C)$ of degree $n \geq 1$. Then $p(z)$ has a zero.

The following proof is based on one in the book *Principles of Mathematical Analysis* 3d., by Walter Rudin (McGraw-Hill Higher Education, New York, 1976).

Proof. We want to find z_0 in C such that $p(z_0) = 0$. Let m be the greatest lower bound of $\{|p(z)| : z \in C\}$. For $|z| = s > 0$, we have

$$\begin{aligned}|p(z)| &= |a_n z^n + a_{n-1} z^{n-1} + \cdots + a_0| \\ &\geq |a_n| |z|^n - |a_{n-1}| |z|^{n-1} - \cdots - |a_0| \\ &= |a_n| s^n - |a_{n-1}| s^{n-1} - \cdots - |a_0| \\ &= s^n [|a_n| - |a_{n-1}| s^{-1} - \cdots - |a_0| s^{-n}].\end{aligned}$$

Because the last expression approaches infinity as s approaches infinity, we may choose a closed disk D about the origin such that $|p(z)| > m + 1$ if z is not in D . It follows that m is the greatest lower bound of $\{|p(z)| : z \in D\}$. Because D is closed and bounded and $p(z)$ is continuous, there exists z_0 in D such that $|p(z_0)| = m$. We want to show that $m = 0$. We argue by contradiction.

Assume that $m \neq 0$. Let $q(z) = \frac{p(z + z_0)}{p(z_0)}$. Then $q(z)$ is a polynomial of degree n , $q(0) = 1$, and $|q(z)| \geq 1$ for all z in C . So we may write

$$q(z) = 1 + b_k z^k + b_{k+1} z^{k+1} + \cdots + b_n z^n,$$

where $b_k \neq 0$. Because $-\frac{|b_k|}{b_k}$ has modulus one, we may pick a real number θ such that $e^{ik\theta} = -\frac{|b_k|}{b_k}$, or $e^{ik\theta} b_k = -|b_k|$. For any $r > 0$, we have

$$\begin{aligned}q(re^{i\theta}) &= 1 + b_k r^k e^{ik\theta} + b_{k+1} r^{k+1} e^{i(k+1)\theta} + \cdots + b_n r^n e^{in\theta} \\ &= 1 - |b_k| r^k + b_{k+1} r^{k+1} e^{i(k+1)\theta} + \cdots + b_n r^n e^{in\theta}.\end{aligned}$$

Choose r small enough so that $1 - |b_k| r^k > 0$. Then

$$|q(re^{i\theta})| \leq 1 - |b_k| r^k + |b_{k+1}| r^{k+1} + \cdots + |b_n| r^n$$

$$= 1 - r^k[|b_k| - |b_{k+1}|r - \cdots - |b_n|r^{n-k}].$$

Now choose r even smaller, if necessary, so that the expression within the brackets is positive. We obtain that $|q(re^{i\theta})| < 1$. But this is a contradiction. ■

The following important corollary is a consequence of Theorem D.4 and the division algorithm for polynomials (Theorem E.1).

Corollary. *If $p(z) = a_n z^n + a_{n-1} z^{n-1} + \cdots + a_1 z + a_0$ is a polynomial of degree $n \geq 1$ with complex coefficients, then there exist complex numbers c_1, c_2, \dots, c_n (not necessarily distinct) such that*

$$p(z) = a_n(z - c_1)(z - c_2) \cdots (z - c_n).$$

Proof. Exercise. ■

A field is called **algebraically closed** if it has the property that every polynomial of positive degree with coefficients from that field factors as a product of polynomials of degree 1. Thus the preceding corollary asserts that the field of complex numbers is algebraically closed.

APPENDIX E POLYNOMIALS

In this appendix, we discuss some useful properties of the polynomials with coefficients from a field. For the definition of a polynomial, refer to Section 1.2. Throughout this appendix, we assume that all polynomials have coefficients from a fixed field F .

Definition. *A polynomial $f(x)$ divides a polynomial $g(x)$ if there exists a polynomial $q(x)$ such that $g(x) = f(x)q(x)$.*

Our first result shows that the familiar long division process for polynomials with real coefficients is valid for polynomials with coefficients from an arbitrary field.

Theorem E.1 (The Division Algorithm for Polynomials). *Let $f(x)$ be a polynomial of degree n , and let $g(x)$ be a polynomial of degree $m \geq 0$. Then there exist unique polynomials $q(x)$ and $r(x)$ such that*

$$f(x) = q(x)g(x) + r(x), \quad (1)$$

where the degree of $r(x)$ is less than m .

Proof. We begin by establishing the existence of $q(x)$ and $r(x)$ that satisfy (1). If $n < m$, take $q(x) = 0$ and $r(x) = f(x)$ to satisfy (1). Otherwise, if $0 \leq m \leq n$, we apply mathematical induction on n . First suppose that

$n = 0$. Then $m = 0$, and it follows that $f(x)$ and $g(x)$ are nonzero constants. Hence we may take $q(x) = f(x)/g(x)$ and $r(x) = 0$ to satisfy (1).

Now suppose that the result is valid for all polynomials with degree less than n for some fixed $n > 0$, and assume that $f(x)$ has degree n . Suppose that

$$f(x) = a_n x^n + a_{n-1} x^{n-1} + \cdots + a_1 x + a_0$$

and

$$g(x) = b_m x^m + b_{m-1} x^{m-1} + \cdots + b_1 x + b_0,$$

and let $h(x)$ be the polynomial defined by

$$h(x) = f(x) - a_n b_m^{-1} x^{n-m} g(x). \quad (2)$$

Then $h(x)$ is a polynomial of degree $k < n$, and therefore we may apply the induction hypothesis or the case where $k < m$ (whichever is relevant) to obtain polynomials $q_1(x)$ and $r(x)$ such that $r(x)$ has degree less than m and

$$h(x) = q_1(x)g(x) + r(x). \quad (3)$$

Combining (2) and (3) and solving for $f(x)$ gives us $f(x) = q(x)g(x) + r(x)$ with $q(x) = a_n b_m^{-1} x^{n-m} + q_1(x)$, which establishes the existence of $q(x)$ and $r(x)$ by mathematical induction.

We now show the uniqueness of $q(x)$ and $r(x)$. Suppose that $q_1(x)$, $q_2(x)$, $r_1(x)$, and $r_2(x)$ exist such that $r_1(x)$ and $r_2(x)$ each has degree less than m and

$$f(x) = q_1(x)g(x) + r_1(x) = q_2(x)g(x) + r_2(x).$$

Then

$$[q_1(x) - q_2(x)]g(x) = r_2(x) - r_1(x). \quad (4)$$

The right side of (4) is a polynomial of degree less than m . Since $g(x)$ has degree m , it must follow that $q_1(x) - q_2(x)$ is the zero polynomial. Hence $q_1(x) = q_2(x)$, and thus $r_1(x) = r_2(x)$ by (4). ■

In the context of Theorem E.1, we call $q(x)$ and $r(x)$ the **quotient** and **remainder**, respectively, for the division of $f(x)$ by $g(x)$. For example, suppose that F is the field of complex numbers. Then the quotient and remainder for the division of

$$f(x) = (3+i)x^5 - (1-i)x^4 + 6x^3 + (-6+2i)x^2 + (2+i)x + 1$$

by

$$g(x) = (3+i)x^2 - 2ix + 4$$

are, respectively,

$$q(x) = x^3 + ix^2 - 2 \quad \text{and} \quad r(x) = (2-3i)x + 9.$$

Corollary 1. Let $f(x)$ be a polynomial of positive degree, and let $a \in F$. Then $f(a) = 0$ if and only if $x - a$ divides $f(x)$.

Proof. Suppose that $x - a$ divides $f(x)$. Then there exists a polynomial $q(x)$ such that $f(x) = (x - a)q(x)$. Thus $f(a) = (a - a)q(a) = 0 \cdot q(a) = 0$.

Conversely, suppose that $f(a) = 0$. By the division algorithm, there exist polynomials $q(x)$ and $r(x)$ such that $r(x)$ has degree less than 1 and

$$f(x) = q(x)(x - a) + r(x).$$

Substituting a for x in the equation above, we obtain $r(a) = 0$. Since $r(x)$ has degree less than 1, it must be the constant polynomial $r(x) = 0$. Thus $f(x) = q(x)(x - a)$. ■

For any polynomial $f(x)$ with coefficients from a field F , an element $a \in F$ is called a **zero** of $f(x)$ if $f(a) = 0$. With this terminology, the preceding corollary states that a is a zero of $f(x)$ if and only if $x - a$ divides $f(x)$.

Corollary 2. Any polynomial of degree $n \geq 1$ has at most n distinct zeros.

Proof. The proof is by mathematical induction on n . The result is obvious if $n = 1$. Now suppose that the result is true for some positive integer n , and let $f(x)$ be a polynomial of degree $n + 1$. If $f(x)$ has no zeros, then there is nothing to prove. Otherwise, if a is a zero of $f(x)$, then by Corollary 1 we may write $f(x) = (x - a)q(x)$ for some polynomial $q(x)$. Note that $q(x)$ must be of degree n ; therefore, by the induction hypothesis, $q(x)$ can have at most n distinct zeros. Since any zero of $f(x)$ distinct from a is also a zero of $q(x)$, it follows that $f(x)$ can have at most $n + 1$ distinct zeros. ■

Polynomials having no common divisors arise naturally in the study of *canonical forms*. (See Chapter 7.)

Definition. Polynomials $f_1(x), f_2(x), \dots, f_n(x)$ are called **relatively prime** if no polynomial of positive degree divides all of them.

For example, the polynomials with real coefficients $f(x) = x^2(x - 1)$ and $h(x) = (x - 1)(x - 2)$ are not relatively prime because $x - 1$ divides each of them. On the other hand, consider $f(x)$ and $g(x) = (x - 2)(x - 3)$, which do not appear to have common factors. Could other factorizations of $f(x)$ and $g(x)$ reveal a hidden common factor? We will soon see (Theorem E.9) that the preceding factors are the only ones. Thus $f(x)$ and $g(x)$ are relatively prime because they have no common factors of positive degree.