# Chapter 10

# **Projections and Matrix Factorizations**

In the first part of this chapter, we say more about left- and right-inverses, and discuss the important concept of *projections*. The second part of this chapter covers *matrix factorizations*, which is about writing a matrix as a product of two or more matrices that have certain convenient structures. We cover the LU decomposition, the QR decomposition, the eigendecomposition and the singular value decomposition. The latter two decompositions are especially useful in Statistics, Econometrics and Data Science, as they reveal important insights about data matrices. We show how these factorizations can be computed in Python, and present several applications.

We begin with the idea of orthogonal matrices, which plays an important role in these topics.

# 10.1 Orthogonal Matrices, Left and Right Inverses, Projections

We refer to vectors  $x \in \mathbb{R}^n$  as *n*-vectors. Recall that the inner product of two *n*-vectors x and y is defined as

$$x \cdot y = \sum_{i=1}^{n} x_i y_i \,.$$

If x and y are column vectors, then  $x \cdot y = x^{\mathrm{T}}y$ . The norm of a vector x, ||x||, is the square root of the inner product  $x \cdot x = x^{\mathrm{T}}x$ . The angle between two vectors x and y is that value  $\theta \in [0, \pi]$  such that  $x \cdot y = ||x|| ||y|| \cos \theta$ . If  $x \cdot y = 0$ , then x and y are orthogonal. If x and y are orthogonal *unit* vectors, we say they are **orthonormal**.

## 10.1.1 Orthogonal Matrices

Matrices with orthonormal columns are called **orthogonal matrices**.<sup>1</sup> That is, a  $n \times k$  matrix

$$A = \begin{bmatrix} a_1 & a_2 & \dots & a_k \end{bmatrix}$$

where  $a_i, i = 1, 2, ..., k$  are *n*-vectors, is an orthogonal matrix if

$$A^{\mathrm{T}}A = \begin{bmatrix} a_{1}^{\mathrm{T}}a_{1} & a_{1}^{\mathrm{T}}a_{2} & \dots & a_{1}^{\mathrm{T}}a_{k} \\ a_{2}^{\mathrm{T}}a_{1} & a_{2}^{\mathrm{T}}a_{2} & \dots & a_{2}^{\mathrm{T}}a_{k} \\ \vdots & \vdots & \ddots & \vdots \\ a_{k}^{\mathrm{T}}a_{1} & a_{k}^{\mathrm{T}}a_{2} & \dots & a_{k}^{\mathrm{T}}a_{k} \end{bmatrix} = I_{k} \,.$$

<sup>&</sup>lt;sup>1</sup>It seems we ought to call matrices with orthonormal columns as *orthonormal matrices*, but unfortunately this is not standard terminology. An orthogonal matrix is one with orthonormal columns. There is no special name for matrices whose columns are merely orthogonal, but not orthonormal.

Another way of describing an orthogonal matrix is one whose transpose is its left-inverse. Orthogonal matrices need not be square. If A is orthogonal and square, then  $A^{-1} = A^{T}$ . The inverse of a square orthogonal matrix is just its transpose. This also means that the transpose of a square orthogonal matrix is orthogonal, since  $I = AA^{-1} = AA^{T}$ . If A is orthogonal, square and symmetric, then

$$A^{-1} = A^{\mathrm{T}} = A$$

A square orthogonal symmetric matrix is its own inverse!

**Example 10.1** All of the following matrices are orthogonal:

$$\begin{split} A_1 &= \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}, \ A_2 &= \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}, \ A_3 &= \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}, \\ A_4 &= \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix}, \ A_5 &= \begin{bmatrix} \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{3}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{3}} & 0 \end{bmatrix}. \end{split}$$

The orthogonal matrices  $A_1$  and  $A_2$  are square but not symmetric (unless  $\theta$  is an integer multiple of  $\pi$ ). We have  $A_i^{-1} = A_i^T \neq A_i$ , i = 1, 2. The orthogonal matrices  $A_3$  and  $A_4$  are square and symmetric. We have  $A_i^{-1} = A_i^T = A_i$ , i = 3, 4. The orthogonal matrix  $A_5$  is not square. We have  $A_5^{T}A_5 = I_2$ , but  $A_5A_5^T \neq I_3$ .

Orthogonal matrices produce transformations that preserve the norms of vectors and the angles between vectors. Suppose A is orthogonal and y = Ax. Then

$$\|y\| = (y \cdot y)^{1/2} = (x^{\mathrm{T}} A^{\mathrm{T}} A x)^{1/2} = (x^{\mathrm{T}} x)^{1/2} = \|x\| \, .$$

Furthermore, suppose  $\theta \in [0, \pi]$  is the angle between  $x_1$  and  $x_2$ , i.e.,  $\theta$  satisfies

$$\cos \theta = \frac{x_1 \cdot x_2}{\|x_1\| \|x_2\|}.$$

Then the angle  $\alpha \in [0, \pi]$  between  $y_1 = Ax_1$  and  $y_2 = Ax_2$  satisfies

$$\cos \alpha = \frac{y_1 \cdot y_2}{\|y_1\| \|y_2\|} = \frac{x_1^{\mathrm{T}} A^{\mathrm{T}} A x_2}{\|x_1\| \|x_2\|} = \frac{x_1 \cdot x_2}{\|x_1\| \|x_2\|} = \cos \theta.$$

It follows that  $\alpha = \theta$ . Tranformations that preserve norms of vectors and angles between vectors are rotations, reflections, and permutations of axes. The matrix  $A_1$  is a permutation matrix.  $A_2$  is a rotation matrix (see Ex. 10.1). The matrices  $A_3$  and  $A_4$  are reflection matrices.  $A_5$  involves reflections and rotations from  $\mathbb{R}^2$  to the  $\mathbb{R}^3$  space.

# 10.1.2 Left- and Right-Inverses

The transpose of an orthogonal matrix, which must have full column rank, is its left-inverse. In fact, left-inverses exist for *all* matrices with full column rank, orthogonal or not. Let A be any  $n \times k$  matrix with full column rank. Full column rank ensures that  $(A^{\mathrm{T}}A)^{-1}$  exists. Then the fact that

$$(A^{\rm T}A)^{-1}A^{\rm T}A = I_k \tag{10.1}$$

shows that a  $(A^{\mathrm{T}}A)^{-1}A^{\mathrm{T}}$  is a left-inverse of A. In fact, we can construct many left-inverses of A if n > k. If n > k, then the null space  $N(A^{\mathrm{T}}, n-k > 0)$  is non-trivial, and we can construct a  $k \times n$  matrix C whose rows are vectors in  $N(A^{\mathrm{T}}, n-r > 0)$ . Such a matrix will satisfy  $CA = 0_{k \times k}$ , so then

$$((A^\mathrm{T} A)^{-1}A^\mathrm{T} + C)A = (A^\mathrm{T} A)^{-1}A^\mathrm{T} A + CA = I_k\,.$$

Example 10.2 Let

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

The null space of  $A^{\mathrm{T}}$  is the set of vectors  $y = \begin{bmatrix} y_1 & y_2 & y_3 \end{bmatrix}^{\mathrm{T}}$  such that

$$A^{\mathrm{T}}y = \begin{bmatrix} 1 & 1 & 1 \\ 2 & 1 & 0 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}.$$

This is the set of all vectors of the form  $\begin{bmatrix} s & -2s & s \end{bmatrix}^{\mathrm{T}}$ . Let

$$C = \begin{bmatrix} s & -2s & s \\ t & -2t & t \end{bmatrix} \text{ for any } s, t \in \mathbb{R}.$$

Then all matrices of the form  $(A^{\mathrm{T}}A)^{-1}A^{\mathrm{T}} + C$  are left-inverses of A.

As we will see in the next section,  $(A^TA)^{-1}A^T$  turns out to be the most important left-inverse of A, so we will use the notation

$$A_{left}^{-1} = (A^{\mathrm{T}}A)^{-1}A^{\mathrm{T}}.$$
(10.2)

A few remarks:

(a) If the  $n \times k$  matrix A has full column rank, with n > k, then the  $n \times n$  matrix  $AA_{left}^{-1} = A(A^{T}A)^{-1}A^{T}$  cannot be the identity matrix  $I_n$ , since  $AA_{left}^{-1}$  will have rank less than n, whereas the identity matrix  $I_n$  has rank n. This emphasizes that  $A_{left}^{-1}$  is not a two-sided inverse. In fact, A will not have a right-inverse, since the rank of the  $n \times n$  matrix AB will be less than n for any  $k \times n$  matrix B.

Mathematics and Programming for the Quantitative Economist

(b) If A is orthogonal, then (10.2) reduces to  $A^{\mathrm{T}}$  since  $A^{\mathrm{T}}A = I_k$ . If A is square and full rank, then (10.2) reduces to the usual two-sided inverse:

$$(A^{\mathrm{T}}A)^{-1}A^{\mathrm{T}} = A^{-1}(A^{\mathrm{T}})^{-1}A^{\mathrm{T}} = A^{-1}.$$

(c) If A is  $n \times k$ , with n < k and full row rank, then  $AA^{T}$  will have an inverse. We can then define a **right-inverse** 

$$A_{right}^{-1} = A^{\mathrm{T}} (AA^{\mathrm{T}})^{-1}$$
.

You can easily see that  $AA_{right}^{-1} = I_n$ . The right-inverse is not unique, and the left-inverse does not exist.

## 10.1.3 Projections

Suppose A is a  $3 \times 2$  matrix with independent columns (so it has full column rank = 2) and y is a  $3 \times 1$  vector. Suppose y is not in the column space of A, i.e.,  $y \notin C(A, 2)$ . Then there are no vectors  $x \in \mathbb{R}^2$  such that Ax = y. This is illustrated in Fig. 10.1 (both panels). Since A has rank 2, the column space C(A, 2) is a plane. We represent this plane in Fig. 10.1 with a parallelogram. The parallelogram is depicted as though the plane is horizontal, but this need not be the case. The figure should be viewed as though the entire  $\mathbb{R}^3$  space has been rotated so that the plane appears horizontal. All vectors of the form Ax,  $x \in \mathbb{R}^2$ , lie on the plane. Since the vector y is not in C(A, 2), it does not lie on the plane. The null space  $N(A^{\mathrm{T}}, 1)$  has dimension 1 and is orthogonal to C(A, 2). It is represented by the perpendicular line passing through the origin O. All vectors in  $N(A^{\mathrm{T}}, 1)$ can be represented by arrows on or parallel to this line.



Fig. 10.1. Non-orthogonal and orthogonal projections.

In Fig. 10.1(a) we draw two non-orthogonal vectors, one  $A\tilde{x} \in C(A, 2)$ , and the other  $\tilde{e}$  such that  $A\tilde{x} + \tilde{e} = y$ . In Fig. 10.1(b) we draw two orthogonal

vectors,  $A\hat{x} \in C(A,2)$  and  $\hat{e} \in N(A^{\mathrm{T}},1)$  such that  $A\hat{x} + \hat{e} = y$ . We call  $A\hat{x}$  the **orthogonal projection** of y onto the column space of X. Which vector  $\tilde{e}$  or  $\hat{e}$  has the smaller norm?

Your intuition will suggest to you that  $\hat{e}$  has a smaller norm than  $\tilde{e}$ . In fact  $\hat{e}$  will have the *smallest* norm among the vectors taking Ax to y. Furthermore, this is also true in the general case: If A is  $n \times k$ , with n > kand full column rank r = k, then

$$\hat{e} = y - A\hat{x}$$
,

where  $A\hat{x}$  is the orthogonal projection of y onto C(A, 2), will have the smallest norm among all vectors e such that Ax + e = y. We show this by showing that  $\hat{x}$  that minimizes  $||e||^2 = (y - Ax)^{\mathrm{T}}(y - Ax)$ , i.e.,

$$\begin{split} \hat{x} &= \mathop{\arg\min}_{x} \ (y - Ax)^{\mathrm{T}} (y - Ax) \\ &= \mathop{\arg\min}_{x} \ (y^{\mathrm{T}}y - x^{\mathrm{T}}A^{\mathrm{T}}y - y^{\mathrm{T}}Ax + x^{\mathrm{T}}A^{\mathrm{T}}Ax) \\ &= \mathop{\arg\min}_{x} \ (y^{\mathrm{T}}y - 2x^{\mathrm{T}}A^{\mathrm{T}}y + x^{\mathrm{T}}A^{\mathrm{T}}Ax) \,. \end{split}$$

The first order condition of this minimization problem is

$$\frac{d}{dx} \|e\|^2 = -2A^{\mathrm{T}}y + 2A^{\mathrm{T}}A\hat{x} = 0_k.$$
(10.3)

This condition can be rewritten as

$$A^{\rm T}(y - A \hat{x}) = A^{\rm T} \hat{e} = 0_k \,. \tag{10.4}$$

The second derivative of  $||e||^2$  is

$$\frac{d^2}{dx \, dx^{\rm T}} \|e\|^2 = 2A^{\rm T}A \,. \tag{10.5}$$

Since A is full column rank, we have  $Ac \neq 0_n$  for any  $c \neq 0_k$ , and therefore

$$c^{\mathrm{T}}(2A^{\mathrm{T}}A)c = 2(Ac)^{\mathrm{T}}(Ac) > 0 \text{ for all } c \neq 0_k$$
.

That is, the second derivative (10.5) is positive definite, and the solution to the FOC (10.3) solves the minimization problem. This condition, of course, just says that  $\hat{e}$  is orthogonal to the column space of A.

Solving (10.4) gives

$$A^{\mathrm{T}}A\,\hat{x} = A^{\mathrm{T}}y$$
  
$$\hat{x} = (A^{\mathrm{T}}A)^{-1}A^{\mathrm{T}}y.$$
 (10.6)

The orthogonal projection of y onto the column space of A is therefore

$$A\hat{x} = A(A^{\mathrm{T}}A)^{-1}A^{\mathrm{T}}y.$$
(10.7)

Mathematics and Programming for the Quantitative Economist

Incidentally, the second line is (10.6) is called the **normal equation** of the projection problem. The  $n \times n$  vector  $A(A^{\mathrm{T}}A)^{-1}A^{\mathrm{T}}$  is called a **projection** matrix. The vector  $\hat{e}$  is

$$\begin{split} \hat{e} &= y - A (A^{\mathrm{T}}A)^{-1}A^{\mathrm{T}}y \\ &= (I_n - A (A^{\mathrm{T}}A)^{-1}A^{\mathrm{T}})y \end{split}$$

Finally, notice that the expression  $(A^{\mathrm{T}}A)^{-1}A^{\mathrm{T}}$  is just  $A^+_{left}$ . Another way to think of the computations above is that y can be written the sum of  $A\hat{x} \in C(A, k)$  and  $\hat{e} \in N(A^{\mathrm{T}}, n-k)$ :

$$A\hat{x} + \hat{e} = y. \tag{10.8}$$

Pre-multiplying both sides by  $A^+_{left} = (A^{\rm T}A)^{-1}A^{\rm T}$  gives

$$\begin{split} A_{left}^{-1}A\hat{x} + A_{left}^{-1}\hat{e} &= A_{left}^{-1}y\\ \hat{x} + (A^{\mathrm{T}}A)^{-1}A^{\mathrm{T}}\hat{e} &= (A^{\mathrm{T}}A)^{-1}A^{\mathrm{T}}y\\ \hat{x} &= (A^{\mathrm{T}}A)^{-1}A^{\mathrm{T}}y \quad \text{since} \ A^{\mathrm{T}}\hat{e} = \mathbf{0}_{k}\,. \end{split}$$

The following example illustrates a very important application of the projection idea.

**Example 10.3** Suppose you have a data set containing 8 observations of two variables x and y as shown in Table 10.1, and shown as points labelled  $(x_i, y_i)$  in Fig. 10.2. The variables are believed to be related according to

$$y = \beta_0 + \beta_1 x + \varepsilon$$
, where  $E(\varepsilon \mid x) = 0$  (10.9)

so that  $E(y \mid x) = \beta_0 + \beta_1 x$ . The objective is to estimate this conditional expectation using the given data set. We do this by finding the best fitting line through these data points.

Table 10.1. A small data set

Obs. no. $i$	$x_i$	$y_i$
1	7.29	31.52
2	5.61	35.89
3	1.25	12.15
4	3.98	20.06
5	7.81	42.32
6	5.11	18.57
7	1.83	26.44
8	8.54	45.25

First we need to define what "best fitting" means. If the data all fall on a straight line, then we can find  $\hat{\beta}_0$  and  $\hat{\beta}_1$  such that  $y_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$ , but since the data do not (and in general it will not), we write

$$y_i = \hat{\beta}_0 + \hat{\beta}_1 x_i + \hat{\varepsilon}_i, \ i = 1, 2, ..., 8.$$
(10.10)

We call  $\hat{\varepsilon}_i$ , i = 1, 2, ..., 8 the "residuals". We can also write (10.10) as

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_8 \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_8 \end{bmatrix} \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix} + \begin{bmatrix} \hat{\varepsilon}_1 \\ \hat{\varepsilon}_2 \\ \vdots \\ \hat{\varepsilon}_8 \end{bmatrix}$$
(10.11)

or simply

$$y = X\hat{\beta} + \hat{\varepsilon} \tag{10.12}$$

where we have overloaded the symbol "y" to refer to both the variable y as well as the vector of observations of y. Comparing with our earlier discussion of projections, y, X,  $\hat{\beta}$  and  $\hat{\varepsilon}$  correspond with b, A,  $\hat{x}$  and  $\hat{e}$  respectively. Given the particular  $x_i$  observations that we have, the matrix X is full column rank. The  $8 \times 1$  vector  $X\hat{\beta}$  represents vectors in the 2-dimensional column space of X. Because the data do not all fall exactly in a straight line, the vector y does not lie in the column space of X.



Fig. 10.2. Fitting a straight line using OLS.

Suppose we define a best fitting line as  $\hat{\beta}_0 + \hat{\beta}_1 x_i$  such that the sum of the square of the vertical distances from the data points to the line (i.e.,

#### Mathematics and Programming for the Quantitative Economist

the "sum of squared residuals") is as small as possible. This means seeking  $\hat\beta_0$  and  $\hat\beta_1$  such that

$$\sum_{i=1}^8 \hat{\varepsilon}_i^2 = \hat{\varepsilon}^{\mathrm{T}} \hat{\varepsilon} = \|\hat{\varepsilon}\|^2.$$

Our discussion of projections tells us that we can do this by choosing  $\hat{\beta}$  so that  $X\hat{\beta}$  is the orthogonal projection of y onto the column space of X, i.e.,

$$\hat{\beta} = (X^{\mathrm{T}}X)^{-1}X^{\mathrm{T}}y.$$
 (10.13)

For our data set, we can calculate  $\hat{\beta}$  to be

$$\hat{\beta} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix} = \begin{bmatrix} 3.614 \\ 10.312 \end{bmatrix}.$$

Our fitted line is  $\hat{y} = 3.614 + 10.312x$ , which is an estimate of  $E(y \mid x)$ . The fitted values  $\hat{y}_i$  are those values such that  $(x_i, \hat{y}_i)$  lie on the fitted line, i.e.,

$$\hat{y}_i = 3.614 + 10.312x_i, \ i = 1, 2, ..., 8.$$
(10.14)

In terms of matrices, the fitted values can be calculated as

$$\hat{y} = X\hat{\beta} = X(X^{\mathrm{T}}X)^{-1}X^{\mathrm{T}}y.$$
 (10.15)

In statistical terms, the fitted values are the "predicted" values of  $y_i$  at  $x = x_i$ ,  $\hat{E}(y \mid x = x_i)$ . In the "geometry" of linear algebra, the vector  $\hat{y}$  is the orthogonal projection of y onto the column space of X.

The model (10.9) is called **simple linear regression** models, and the method we have described for estimating it is **ordinary least squares**. It extends readily to **multiple linear regression** models

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_K x_K + \varepsilon \,, \ \text{where} \ E(\varepsilon \mid x_1, x_2, \dots, x_K) = 0$$

where  $x_k$  now refers to a variable, not the kth observation of x. Given a data set  $\{y_i, x_{1,i}, x_{2,i}, \dots, x_{K,i}\}_{i=1}^n$  we would write

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_{1,1} & x_{2,1} & \dots & x_{K,1} \\ 1 & x_{1,2} & x_{2,2} & \dots & x_{K,2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{1,n} & x_{2,n} & \dots & x_{K,n} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \\ \vdots \\ \hat{\beta}_K \end{bmatrix} + \begin{bmatrix} \hat{\varepsilon}_1 \\ \hat{\varepsilon}_2 \\ \vdots \\ \hat{\varepsilon}_n \end{bmatrix}$$
(10.16)

or simply  $y = X\hat{\beta} + \hat{\varepsilon}$ . As long as X has full column rank, the  $\hat{\beta}$  that minimizes the sum of squared residuals is still given by (10.13). The fitted regression "line" is

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \dots + \hat{\beta}_K x_K.$$

Linear regression is a fundamental technique in statistics, econometrics, and data science. Even methods for working with highly non-linear models often use the linear regression methodology in some way. In our discussion here, we have merely been concerned with the geometric aspects of the "line fitting" problem. In applications, we will be concerned with the statistical properties and interpretation of the estimators, and how to use the model for prediction, causal inference, and testing theories.

#### 10.1.4 Exercises

Ex. 10.1 Show that

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

is an orthogonal matrix, and that the transformation Ax rotates the vector x by an angle of  $\theta$  without changing its norm.

Ex. 10.2 The orthogonal projection of

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \text{ onto the column space of } X = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$$

is  $\hat{y}=X\hat{\beta}$  where  $\hat{\beta}=\begin{bmatrix}\hat{\beta}_0&\hat{\beta}_1\end{bmatrix}^{\rm T}=(X^{\rm T}X)^{-1}X^{\rm T}y.$  Show that

$$\hat{\beta}_0 = \overline{y} - \hat{\beta}_1 \overline{x} \quad \text{and} \quad \hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \overline{x})(y_i - \overline{y})}{\sum_{i=1}^n (x_i - \overline{x})^2}$$

**Ex. 10.3** Suppose  $\hat{y}$  is an orthogonal projection of y onto the column space of X. Find the projection of  $\hat{y}$  onto the column space of X?

**Ex. 10.4** (Projection onto a vector) (a) Show that the formula for the orthogonal projection of a (non-zero) vector y onto another (non-zero) vector x is

$$\hat{y} = \frac{x^{\mathrm{T}}y}{x^{\mathrm{T}}x}x$$

Find  $\hat{e}$  orthogonal to y such that  $y = \hat{y} + \hat{e}$ .

(b) Describe the orthogonal projection of y onto x if  $x = i_n$ , the n-vector of ones.

**Ex. 10.5** Suppose the columns of the  $n \times k$  matrix X are orthogonal. Find the formula for  $\hat{y}$ , the orthogonal projection of an *n*-vector y onto the column space of X. How does the formula simplify if X is an orthogonal matrix?

## 10.2 The LU and QR Decompositions

We begin our discussion of matrix factorizations with the LU and QR decompositions. To develop a good understanding of these factorizations (and the other ones we will discuss shortly), we have to talk about how they are constructed. However, our focus in on understanding how, why and when

395

# Mathematics and Programming for the Quantitative Economist

they work, rather than on actually computing them. Packages exists in Python for computing these factorizations.

We begin with a few additional facts regarding inverses. You should take a few moments to see if you can verify these statements. First, the inverse of a diagonal matrix

$$A = \text{diag}(a_{11}, a_{22}, \dots, a_{kk})$$

is the diagonal matrix

$$A^{-1} = \operatorname{diag}(a_{11}^{-1}, a_{22}^{-1}, \dots, a_{kk}^{-1})$$

This inverse exists only if none of the diagonal elements are zero. If A is upper (lower) triangular with no zeros along the diagonal, then  $A^{-1}$  is also upper (lower) triangular with no zeros along the diagonal. If A is upper (lower) triangular with ones along the diagonal, then  $A^{-1}$  is also upper (lower) triangular with ones along the diagonal. Finally, if A is symmetric, then  $A^{-1}$  is symmetric.

# 10.2.1 The LU Decomposition

Let A be a full rank square matrix. Recall from Chapter 4 and Chapter 8 that we can apply row operations to A to reduce it to an upper triangular matrix (see for instance Example 4.3). These operations can be achieved by pre-multiplying A with elementary row operation matrices that adds a multiple of an upper row to a lower row. These are matrices that have ones along the diagonal, one non-zero element below the diagonal, and all other elements zero. For instance, to add  $\alpha$  times the first row of a 3 × 3 matrix to the third row, we can multiply that matrix with

$$E_{[3]\leftarrow[3]+\alpha[1]} = \begin{bmatrix} 1 & 0 & 0\\ 0 & 1 & 0\\ \alpha & 0 & 1 \end{bmatrix}.$$

Such "elimination" matrices are lower triangular, with ones down the diagonal, and so its inverse will have the same structure.

Occasionally, we will have to permutate the order of the rows of the matrix, as in Example 8.15. Such permutations can be done at the start, before starting the elimination process, by pre-multiplying A by a permutation matrix P.

In other words, starting with the appropriately permutated matrix PA, we apply a series of elimination matrices  $E_1, E_2, \ldots, E_p$  until PA becomes an upper triangular matrix U, i.e.,

$$E_p E_{p-1} \dots E_2 E_1 P A = U \,.$$

This gives the decomposition

$$PA = E_p^{-1} E_{p-1}^{-1} \dots E_2^{-1} E_1^{-1} U = LU \,.$$

In this decomposition, the diagonals elements of U will in general not be ones. We can alternatively write U as DU where the diagonal elements of Uare replaced with ones, and D is a diagonal matrix containing the diagonal elements of the previous U. We have

$$PA = LDU$$
.

This is the **LU decomposition** of a full rank square matrix. If permutations are not required, we have A = LU or A = LDU. In the latter form, the diagonals elements of L and U are all ones.

### Example 10.4 We have

$$A = \begin{bmatrix} 1 & 2 & 0 \\ 3 & 5 & 1 \\ 7 & 4 & 2 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 3 & 1 & 0 \\ 7 & 10 & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 & 0 \\ 0 & -1 & 1 \\ 0 & 0 & -8 \end{bmatrix} \qquad LU$$
$$= \begin{bmatrix} 1 & 0 & 0 \\ 3 & 1 & 0 \\ 7 & 10 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & -8 \end{bmatrix} \begin{bmatrix} 1 & 2 & 0 \\ 0 & 1 & -1 \\ 0 & 0 & 1 \end{bmatrix} \qquad LDU$$

The LU decomposition is merely an expression of Gaussian elimination in matrix form. It can be used to solve full-rank systems of n equations in n unknowns by breaking them into triangular systems. Triangular systems of equations are easy to solve. If

$$a_{11}x_1 + a_{12}x_2 + a_{13}x_3 = b_1$$
  
$$a_{22}x_2 + a_{23}x_3 = b_2$$
  
$$a_{33}x_3 = b_3$$

then  $x_3 = b_3/a_{33}$ . Substituting  $x_3$  into the second equation gives  $x_2$ , and then substituting  $x_3$  and  $x_2$  into the first equation gives  $x_1$ . "Lower triangular systems" are just as easy to solve. Then given a *n*-equation *n*-unknown system Ax = b, write Ax = b as LUx = b, or Lc = b, where Ux = c. First solve Lc = b for c, then solve Ux = c for x.

10.2.1.1 The Cholesky Decomposition When the LDU decomposition is applied to symmetric positive definite matrices, we get the Cholesky Decomposition. A square matrix A is positive definite if

$$x^{\mathrm{T}}Ax > 0$$
 for all  $x \neq 0$ .

For such matrices, the LDU decomposition becomes

$$A = LDL^{\mathrm{T}}$$

where L is lower triangular with ones along the diagonal and  $D = \text{diag}(d_{11}, d_{22}, \dots, d_{nn})$  with  $d_{kk} > 0$  for all  $k = 1, 2, \dots, n$ . We can

Mathematics and Programming for the Quantitative Economist

write further write  $D = D^{1/2}D^{1/2}$  where  $D^{1/2} = \text{diag}(d_{11}^{1/2}, d_{22}^{1/2}, \dots, d_{nn}^{1/2})$ . Absorbing the  $D^{1/2}$  matrix into the L matrix gives the decomposition

$$A = LD^{1/2}D^{1/2}L^{\mathrm{T}} = LD^{1/2}(D^{1/2})^{\mathrm{T}}L^{\mathrm{T}} = L^{*}L^{*^{\mathrm{T}}}$$

where  $L^* = LD^{1/2}$ .

Example 10.5 The Cholesky decomposition of

$$A = \begin{bmatrix} 4 & -2 & 4 \\ -2 & 5 & 8 \\ 4 & 8 & 14 \end{bmatrix}$$

is

$$\begin{bmatrix}
4 & -2 & 4 \\
-2 & 5 & 8 \\
4 & 8 & 14
\end{bmatrix} = \underbrace{\begin{bmatrix}
1 & 0 & 0 \\
-0.5 & 1 & 0 \\
-1 & 1.5 & 1
\end{bmatrix}}_{L} \underbrace{\begin{bmatrix}
4 & 0 & 0 \\
0 & 4 & 0 \\
0 & 0 & 1
\end{bmatrix}}_{D} \underbrace{\begin{bmatrix}
1 & -0.5 & -1 \\
0 & 1 & 1.5 \\
0 & 0 & 1
\end{bmatrix}}_{L^{\mathrm{T}}}$$

$$= \underbrace{\begin{bmatrix}
2 & 0 & 0 \\
-0.5 & 2 & 0 \\
-1 & 1.5 & 1
\end{bmatrix}}_{L^{*}} \underbrace{\begin{bmatrix}
2 & -0.5 & -1 \\
0 & 2 & 1.5 \\
0 & 0 & 1
\end{bmatrix}}_{L^{*\mathrm{T}}}$$

## 10.2.2 The QR Decompositions

Suppose the  $n \times k$  matrix A has full column rank, with column space C(A, r). We can convert A into an orthonormal matrix Q that spans the same space C(A, r). There are two ways to do this. The first uses the **Gram-Schmidt** (GS) procedure: let

$$A = \begin{bmatrix} a_1 & a_2 & \dots & a_k \end{bmatrix}$$

where  $a_i$ , i = 1, 2, ..., k are the k columns of A. The GS procedure goes through each column from left to right. At each step, it replaces the column by a vector that is orthogonal to all the previous columns, and normalizes the new vector to one.

First define

$$\begin{split} \tilde{A}_1 &= \begin{bmatrix} \tilde{a}_1 & a_2 & \dots & a_k \end{bmatrix} & \text{where} \quad \tilde{a}_1 &= a_1 \ ( \text{ i.e., do nothing } ) \\ A_1 &= \begin{bmatrix} q_1 & a_2 & \dots & a_k \end{bmatrix} & \text{where} \quad q_1 &= a_1 / \|a_1\| \end{split}$$

This is the same as *post*-multiplying A by  $D_1 = \text{diag}(||a_1||^{-1}, 1, \dots, 1)$ . Next, project  $a_2$  onto  $q_1$  to get the linear projection  $(q_1^{\mathrm{T}}a_2)q_1$  (remember  $q^{\mathrm{T}}q = 1$ ) and define

$$\tilde{a}_2 = a_2 - (q_1^{\rm T} a_2) q_1 \, .$$

Set

$$\begin{split} \dot{A}_2 &= \begin{bmatrix} q_1 & \tilde{a}_2 & \dots & a_k \end{bmatrix} \\ A_2 &= \begin{bmatrix} q_1 & q_2 & \dots & a_k \end{bmatrix} \quad \text{where} \quad q_2 &= \tilde{a}_2 / \| \tilde{a}_2 \| \,. \end{split}$$

Notice that to get  $\tilde{A}_2$  we subtract  $q_1^{\mathrm{T}}a_2$  times the first column of  $A_1$  from the second column of  $A_1$ . The is an "elementary column operation" and is equivalent to *post*-multiplying  $A_1$  by an appropriate *upper*-triangular matrix " $E_2$ ". Then to normalize the second column, we again post-multiply with an appropriate diagonal matrix  $D_2$ , i.e.,  $A_2 = AD_1E_2D_2$ .

Next, project  $a_3$  onto the column space spanned by the first two columns of  $A_2$ . Since the first two columns of  $A_2$  make up an orthogonal matrix, the linear project is  $(q_1 \cdot a_3)q_1 + (q_2 \cdot a_3)q_2$  (see Ex. 10.5). Define

$$\tilde{a}_3 = a_3 - (q_1 \cdot a_3)q_1 - (q_2 \cdot a_3)q_2$$

and normalize it to unit length. Set

$$\begin{split} A_3 &= \begin{bmatrix} q_1 & \tilde{a}_2 & \tilde{a}_3 & \dots & a_k \end{bmatrix} \\ A_3 &= \begin{bmatrix} q_1 & q_2 & q_3 & \dots & a_k \end{bmatrix} \quad \text{where} \quad q_3 &= \tilde{a}_3 / \|\tilde{a}_3\| \,. \end{split}$$

This is equivalent to post-multiplying  $A_2$  by an appropriate upper triangular matrix  $E_3$  followed by an appropriate diagonal matrix  $D_3$ , i.e.,  $A_3 = AD_1E_2D_2E_3D_3$ .

Repeat this process until you have orthogonalized every row of A:

$$AD_1E_2D_2E_3D_3\ldots E_pD_p=Q=\begin{bmatrix}q_1&q_2&\ldots&q_k\end{bmatrix}\,.$$

The matrix A and Q are  $n \times k$  whereas the D and E matrices are full rank diagonal or upper triangular matrices. We have

$$A = QD_p^{-1}E_p^{-1} \dots D_2^{-1}E_2^{-1}D_1^{-1} = QR$$

where R is a  $k \times k$  upper triangular matrix.

In summary, the GS procedure converts the columns of A into orthonormal columns in Q. The R matrix is a summary of all of the steps needed to do this. The columns of Q forms an orthonormal basis for C(A, k). Since Q is orthogonal,  $Q^{-1} = Q^{T}$ .

The QR decomposition provides an efficient way to do least squares computations. Recall the formula for orthogonally projecting an *n*-vector y onto the column space of an  $n \times k$  matrix X with full column rank. The orthogonal projection is  $\hat{y} = X\hat{\beta}$ , where  $\hat{\beta} = (X^{\mathrm{T}}X)^{-1}X^{\mathrm{T}}y$ . For large matrices, the inverse of  $X^{\mathrm{T}}X$  can be computationally expensive. However, if we have the QR decomposition

$$X = QR$$
 where  $Q^{T}Q = I$  and R is upper triangular,

399

then from the normal equations  $X^{\mathrm{T}} X \hat{\beta} = X^{\mathrm{T}} y$ , we have

$$\begin{split} X^{\mathrm{T}} X \hat{\beta} &= X^{\mathrm{T}} y \\ R^{\mathrm{T}} Q^{\mathrm{T}} Q R \hat{\beta} &= R^{\mathrm{T}} Q^{\mathrm{T}} y \\ R^{\mathrm{T}} R \hat{\beta} &= R^{\mathrm{T}} Q^{\mathrm{T}} y \\ R \hat{\beta} &= Q^{\mathrm{T}} y \end{split}$$

where the last line holds because R is a full rank square matrix. We have turned the problem of solving the normal equations for  $\hat{\beta}$  into solving a triangular system of equations.

**Example 10.6** The following is a small example of the QR decomposition obtained via the GS procedure:

$$A = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 2 \\ 1 & 0 & 3 \\ 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} \frac{1}{2} & 0 & -\frac{1}{2} \\ \frac{1}{2} & \frac{1}{\sqrt{2}} & \frac{1}{2} \\ \frac{1}{2} & -\frac{1}{\sqrt{2}} & \frac{1}{2} \\ \frac{1}{2} & 0 & -\frac{1}{2} \end{bmatrix} \begin{bmatrix} 2 & 2 & \frac{7}{2} \\ 0 & \sqrt{2} & -\frac{1}{\sqrt{2}} \\ 0 & 0 & \frac{3}{2} \end{bmatrix} = QR.$$

10.2.2.1 Householder Reflectors An approach to generating QR decompositions that is numerically more stable than the GS procedure uses what are known as **Householder matrices** or **Householder reflectors**.<sup>2</sup> For any *n*-vector v, define the corresponding Householder reflector to be the matrix

$$H = I_n - \frac{2}{\|v\|^2} v v^{\mathrm{T}}$$
 .

It can be shown that this matrix, when pre-multiplied into an *n*-vector x, reflects x about a plane whose normal vector is v/||v||.

**Example 10.7** Consider  $\mathbb{R}^2$  and the line  $x_2 = -x_1$ . A normal vector to this line is

$$\frac{v}{\|v\|} = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix} .$$

The Householder matrix is

$$H = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - 2 \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix} = \begin{bmatrix} 0 & -1 \\ -1 & 0 \end{bmatrix}$$

which transforms any vector  $\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$  into the vector  $\begin{bmatrix} -x_2 \\ -x_1 \end{bmatrix}$ , which is a reflection across the line  $x_2 = -x_1$ .

400

<sup>&</sup>lt;sup>2</sup>Alston Scott Householder (1904-1993).

Householder reflectors are symmetric and orthogonal (see exercises), so we have  $H H = I_n$ . They can be used to eliminate entries below pivots, providing an alternative to Gaussian elimination. In particular, if v is defined as

$$v = x \pm \|x\|e_1 \text{ where } x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \text{ and } e_1 = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix},$$

then  $H x = \begin{bmatrix} \alpha_1 & 0 & \dots & 0 \end{bmatrix}^{\mathrm{T}}$ .

We show this for  $v = x + ||x||e_1$ . As an exercise, you should repeat the argument for  $v = x - ||x||e_1$ . In practice, we choose "+" or "-" to make ||v|| as large as possible, to avoid dividing by small numbers when forming H.

We first note that

$$\begin{aligned} 2v^{\mathrm{T}}x &= 2(x^{\mathrm{T}}x + \|x\|e^{\mathrm{T}}x) = 2x^{\mathrm{T}}x + 2\|x\|x_{1} \text{ and} \\ v^{\mathrm{T}}v &= (x^{\mathrm{T}} + \|x\|e_{1}^{\mathrm{T}})(x + \|x\|e_{1}) \\ &= x^{\mathrm{T}}x + 2\|x\|x_{1} + \|x\|^{2} \\ &= 2x^{\mathrm{T}}x + 2\|x\|x_{1}, \end{aligned}$$

therefore

$$\begin{split} H \, x \, &= \, \left( I_n - \frac{2}{\|v\|^2} v v^{\mathrm{T}} \right) x \, = \, x - \frac{2}{\|v\|^2} v v^{\mathrm{T}} x \\ &= \, x - \frac{2 v^{\mathrm{T}} x}{v^{\mathrm{T}} v} v \, = \, x - v \, = \, -\|x\| e_1 = \begin{bmatrix} -\|x\| \\ 0 \\ \vdots \\ 0 \end{bmatrix} \, . \end{split}$$

The idea is this: given a full column  $n \times k$  matrix  $A = \begin{bmatrix} a_1 & a_2 & \dots & a_k \end{bmatrix}$ , pre-multiply A by  $H_1$  with  $v = a_1 \pm ||a_1||e_1$ . This places a non-zero term in the (1, 1)th position of A and makes everything below it 0. Next define

$$\widetilde{H}_2 = \begin{bmatrix} I_1 & 0 \\ 0 & H_2 \end{bmatrix}$$

where  $H_2$  is the  $(n-1) \times (n-1)$  Householder reflector that turns the  $(n-1) \times 1$  vector  $[y_2 \ y_3 \ \dots \ y_n]^{\mathrm{T}}$  into a  $(n-1) \times 1$  vector of the form  $[\alpha_2 \ 0 \ \dots \ 0]^{\mathrm{T}}$ . The structure of  $\widetilde{H}_2$  means that the first column of the previously transformed A matrix is not affected. Pre-multiplying  $\widetilde{H}_2$  into  $H_1A$ , we get a matrix with non-zero terms in the (1,1)th and (2,2)th position, and all terms below them equal to zero. Repeat in similar fashion

until we get a matrix of the form

$$\widetilde{H_p} \dots \widetilde{H_2} H_1 A = \begin{bmatrix} \alpha_1 & * & * & \dots & * \\ 0 & \alpha_2 & * & \dots & * \\ 0 & 0 & \alpha_3 & \dots & * \\ 0 & 0 & 0 & \dots & \alpha_k \\ 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 0 \end{bmatrix}$$

Define the matrix on the right to be the R matrix. All of the Householder and modified Householder matrices on the left are symmetric and orthogonal. Therefore

$$A = H_1 \widetilde{H}_2 \dots \widetilde{H}_p R = QR \,.$$

Since the product of orthogonal matrices are orthogonal<sup>3</sup>, Q is orthogonal.

The QR decomposition produced here is slightly different from the one produced by the GS procedure. For one thing, the GS procedure produced A = QR where Q and R are  $n \times k$  and  $k \times k$  respectively, with the Q matrix containing an orthogonal basis for C(A, k). In the Householder approach, Q and R are  $n \times n$  and  $n \times k$  respectively. The Householder Q contains an orthogonal basis for the entire  $\mathbb{R}^n$  space, with the first k columns serving as an orthogonal basis for C(A, k). Since the bottom n - k rows of the Householder R are zeros, we can discard the right-most n - k columns of Q and bottom n - k rows of R to get a QR decomposition akin to the one produced by GS.

**Example 10.8** For the matrix A in Example 10.6, we have the following QR decomposition obtained via Householder reflectors:

$$\mathbf{A} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 2 \\ 1 & 0 & 3 \\ 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} -\frac{1}{2} & 0 & \frac{1}{2} & -\frac{1}{\sqrt{2}} \\ -\frac{1}{2} & -\frac{1}{\sqrt{2}} & -\frac{1}{2} & 0 \\ -\frac{1}{2} & \frac{1}{\sqrt{2}} & -\frac{1}{2} & 0 \\ -\frac{1}{2} & 0 & \frac{1}{2} & \frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} -2 & -2 & -\frac{7}{2} \\ 0 & -\sqrt{2} & \frac{1}{\sqrt{2}} \\ 0 & 0 & -\frac{3}{2} \\ 0 & 0 & 0 \end{bmatrix}$$
$$= \begin{bmatrix} -\frac{1}{2} & 0 & \frac{1}{2} \\ -\frac{1}{2} & -\frac{1}{\sqrt{2}} & -\frac{1}{2} \\ -\frac{1}{2} & \frac{1}{\sqrt{2}} & -\frac{1}{2} \\ -\frac{1}{2} & 0 & \frac{1}{2} \end{bmatrix} \begin{bmatrix} -2 & -2 & -\frac{7}{2} \\ 0 & 0 & -\frac{3}{2} \\ 0 & 0 & 0 \end{bmatrix}.$$

In the last line, we dropped the fourth column and last row of the previous line's "expanded" Q and R matrices, respectively.

Most computer packages now use the Householder approach to generate QR decompositions.

402

A

 $<sup>{}^{3}\</sup>mathrm{If}\;A^{\mathrm{T}}A=I\;\mathrm{and}\;B^{\mathrm{T}}B=I,\;\mathrm{then}\;(AB)^{\mathrm{T}}AB=B^{\mathrm{T}}A^{\mathrm{T}}AB=I.$ 

## 10.2.3 Exercises

**Ex. 10.6** (a) Show that Householder reflectors are symmetric and orthogonal. (b) Show that matrices of the form

$$\widetilde{H} = \begin{bmatrix} I_p & 0_{p \times (n-p)} \\ 0_{(n-p) \times p} & H_{n-p} \end{bmatrix}$$

are orthogonal and symmetric, where  $H_{n-p}$  is a  $n-p \times n-p$  Householder reflector.

**Ex.** 10.7 (a) Show that the columns of the matrix A in Example 10.6 and Example 10.8 are not orthogonal.

(b) Show that the Q matrices in both of those examples are orthogonal matrices.

## 10.3 The Eigendecomposition and the SVD

The eigendecomposition of a square matrix A refers to the factorization  $A = S\Lambda S^{-1}$  where the columns of S are eigenvectors of A and  $\Lambda$  is a diagonal matrix whose diagonal elements are the *n* eigenvalues of A, both terms to be defined shortly.<sup>4</sup> Many square matrices can be factorized in this manner. If so, we say that A is diagonalizable, since it follows that  $S^{-1}AS = \Lambda$ . All symmetric matrices are diagonalizable, and their eigendecompositions have additional useful structure. When applied to symmetric matrices, the eigendecomposition is called the spectral decomposition of A.

The eigenvalues of a square matrix A contain information regarding many of its key properties. For instance, the product of the eigenvalues of A gives its determinant and the sum of the eigenvalues gives its trace. The rank of A is equal to the number of non-zero eigenvalues in  $A^{T}A$ . A symmetric matrix A is positive semi-definite  $(c^{T}Ac \ge 0 \text{ for all } c \ne 0_{n})$  if and only if all of its eigenvalues are non-negative. It is positive definite  $(c^{T}Ac > 0 \text{ for all } c \ne 0)$  if and only if all of its eigenvalues are positive.

The eigendecomposition of A is useful in dynamic models. If

$$x_{t+1} = Ax_t, \ t = 0, 1, 2, \dots$$

then from an initial  $x_0 \in \mathbb{R}^n$ , we have

$$x_1 = Ax_0, \ x_2 = Ax_1 = A^2x_0, \ x_t = A^tx_0,$$

and so on. If  $A = S\Lambda S^{-1}$ , then

$$A^t = \underbrace{S\Lambda S^{-1} S\Lambda S^{-1} S\Lambda S^{-1} \dots S\Lambda S^{-1}}_{t \text{ terms}} = S\Lambda^t S^{-1}.$$

The *t*th power of a diagonal matrix is easy to compute — just raise each of the diagonal elements to the *t*th power. This also means that if the eigenvalues of A all have modulus less than one, then  $x_t$  will eventually converge to the zero vector, since  $\Lambda^t \to 0_{n \times n}$  as  $t \to \infty$ .

403

 $<sup>^{4}</sup>$ The term "eigen" is the German word for "own". It is often used in mathematics to construct names of characteristic features of various mathematical objects.

# 10.3.1 Eigenvalues and Eigenvectors

Let A be an  $n \times n$  matrix of *real* numbers and consider the transformation

$$y = Ax, x \in \mathbb{C}^n. \tag{10.17}$$

It will be convenient for certain statements we would like to make to consider the domain of (10.17) to be the set of all *n*-vectors of complex numbers, which includes the set of *n*-vectors of real numbers. The matrix *A*, however, will be restricted to square matrices containing only real numbers.

Given A, there will be certain vectors  $x \in \mathbb{C}^n$  for which the transformed vector y = Ax is simply a scalar multiple of x, i.e.,

$$Ax = \lambda x \,. \tag{10.18}$$

where  $\lambda$  is possibly complex-valued scalar. This will typically not be the case for most vectors  $x \in \mathbb{C}^n$ , but there will be some vectors for which (10.18) holds. Such vectors are called **eigenvectors** corresponding to the **eigenvalue**  $\lambda$ .

If  $\lambda$  is real, then for real vectors x satisfying (10.18) the transformation y = Ax merely stretches or shrinks the vector (if  $|\lambda| \neq 1$ ), flipping its direction if  $\lambda < 0$ . If  $\lambda = 0$ , then vectors for which (10.18) holds are "eliminated". We are unaware of any useful "geometric intuition" for the case where  $\lambda$  or x are complex-valued, but this does not make such eigenvalues and eigenvectors any less useful or valid. Of course, (10.18) always holds for  $x = 0_n$ . We are interested in situations where (10.18) holds for non-trivial x.

Example 10.9 The matrix

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

has an eigenvalue  $\lambda = 3$ . Any vector of the form

$$x = \begin{bmatrix} \frac{3}{2}s \\ s \end{bmatrix}, \, s \in \mathbb{C}$$

is an eigenvector corresponding to  $\lambda = 3$ . We can verify this by direct multiplication:

$$Ax = \begin{bmatrix} 1 & 3\\ 2 & 0 \end{bmatrix} \begin{bmatrix} \frac{3}{2}s\\ s \end{bmatrix} = \begin{bmatrix} \frac{3}{2}s + 3s\\ 2\left(\frac{3}{2}s\right) \end{bmatrix} = 3 \begin{bmatrix} \frac{3}{2}s\\ s \end{bmatrix} = \lambda x \,. \tag{10.19}$$

Notice that there is a whole subspace of vectors associated with the eigenvalue  $\lambda = 3$ . Furthermore, this set of eigenvectors includes both real and complex vectors.

How do we find eigenvalues and eigenvectors for a general  $n \times n$  matrix A? Rewrite (10.18) as

$$\lambda x - Ax = (\lambda I_n - A)x = 0. \tag{10.20}$$

This tell us that (a) for any eigenvalue  $\lambda$ , the eigenvectors associated with  $\lambda$  make up the nullspace of the matrix  $\lambda I_n - A$  (we call this the **eigenspace** of A), and (b) this nullspace of eigenvectors will be non-trivial only if  $\lambda I_n - A$  has determinant zero. That is, the eigenvalues must satisfy

$$\det(\lambda I_n - A) = 0 \tag{10.21}$$

or equivalently,  $\det(A - \lambda I_n) = 0$ . The left-hand-side of (10.21) is an order-*n* polynomial in  $\lambda$  called the **characteristic polynomial** of the matrix *A*. Equation (10.21) can be factorized as described in the Fundamental Theorem of Algebra, i.e.,

$$\det(\lambda I_n - A) = \lambda^n + \zeta_{n-1}\lambda^{n-1} + \dots + \zeta_1\lambda + \zeta_0$$
  
=  $(\lambda - \lambda_1)(\lambda - \lambda_2) \times \dots \times (\lambda - \lambda_n) = 0.$  (10.22)

The eigenvalues of A are the n roots  $\lambda_j$ , j = 1, 2, ..., n, of this equation. The eigenvalues may be real or complex-valued, and there may be repeated values. Complex-valued eigenvalues will appear in conjugate pairs, i.e., if a+bi is an eigenvalue, then a-bi will also be an eigenvalue, where i is the "imaginary number"  $i = \sqrt{-1}$ . For each eigenvalue  $\lambda$ , the corresponding eigenvectors are found as the solutions to (10.20). We will look at (10.22) in more detail shortly.

Example 10.10 The eigenvalues of

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

are the solutions to the equation

$$\begin{split} \det(\lambda I_3 - A) &= \begin{vmatrix} \lambda - 2 & -1 & 0 \\ 1 & \lambda - 2 & 0 \\ 0 & 0 & \lambda - 3 \end{vmatrix} \\ &= (\lambda - 3)((\lambda - 2)^2 + 1) \\ &= (\lambda - 3)(\lambda - (2 + i))(\lambda - (2 - i)) = 0 \end{split}$$

The eigenvalues are  $\lambda_1 = 3$ ,  $\lambda = 2 + i$  and  $\lambda = 2 - i$ . The eigenvectors associated with  $\lambda_1 = 3$  are the solutions to (10.20) with  $\lambda_1 = 3$ , i.e.

$$(3I_3 - A)x_1 = \begin{bmatrix} 1 & -1 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} x_{11} \\ x_{12} \\ x_{13} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}.$$

The solution to this equation are all vectors  $x_1$  such that  $x_{11} - x_{12} = 0$  and  $x_{11} + x_{12} = 0$  (which together implies  $x_{11} = x_{12} = 0$ ), with  $x_{13}$  "free". That is, the eigenvectors are all vectors

$$x_1 = \begin{bmatrix} 0 & 0 & s \end{bmatrix}^{\mathrm{T}}, s \in \mathbb{C}.$$

For  $\lambda_2 = 2 - i$ , (10.20) becomes

$$((2-i)I_3 - A)x_2 = \begin{bmatrix} -i & -1 & 0 \\ 1 & -i & 0 \\ 0 & 0 & -1 - i \end{bmatrix} \begin{bmatrix} x_{21} \\ x_{22} \\ x_{23} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}.$$

The solutions are all vectors  $x_2$  satisfying

$$\begin{aligned} -ix_{21} - x_{22} &= 0\\ x_{21} - ix_{22} &= 0\\ (-1 - i)x_{23} &= 0 \end{aligned}$$

The third equation says that  $x_{23}$  must be 0. Multiplying the first equation by *i* gives the second equation, so  $x_{21}$  and  $x_{22}$  must satisfy  $x_{21} - ix_{22} = 0$ . Letting  $x_{21}$  be the free parameter *s*, we see that the solutions are all vectors of the form

$$x_2 = \begin{bmatrix} s & -is & 0 \end{bmatrix}^{\mathrm{T}}, s \in \mathbb{C}.$$

This is a 1-dimensional vector subspace (it has one free parameter) of the vector space  $\mathbb{C}^3$ . Notice that there are no real eigenvectors associated with  $\lambda_2 = 2-i$ . As an exercise, you should show that the eigenvectors associated with the eigenvalue  $\lambda_3 = 2+i$  are all vectors of the form

$$x_3 = \begin{bmatrix} s & is & 0 \end{bmatrix}^{\mathrm{T}}, s \in \mathbb{C}.$$

Eigenvectors are sometimes required to have unit norm<sup>5</sup> in addition to satisfying  $Ax = \lambda x$ . We will sometimes use unit eigenvectors specifically, but our presentation so far emphasizes that the set of eigenvectors corresponding to any eigenvalue comprise entire subspaces. Statistical libraries will return eigenvectors with unit norm, and real-valued if possible. In

$$x_1 = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1n} \end{bmatrix}^{\mathrm{T}}$$
 where  $x_{1j} = a_j + b_j \, i$  ,  $j = 1, 2, \dots, n$ 

is defined as

$$\|x_1\| = (|x_{11}|^2 + |x_{12}|^2 + \dots + |x_{1n}|^2)^{\frac{1}{2}}$$

where  $|x_{1j}|^2 = a_j^2 + b_j^2$ . It can be computed as the square root of  $\overline{x}_1^{\mathrm{T}} x_1$  where  $\overline{x}_1^{\mathrm{T}}$  is the "conjugate transpose" of  $x_1$ , i.e., transpose  $x_1$  and replace each complex entry with its conjugate.

<sup>&</sup>lt;sup>5</sup>The norm of a complex-valued vector

Example 10.10, the eigenvectors corresponding to  $\lambda_1 = 3, \lambda_2 = 2 - i$  and  $\lambda_3 = 2 + i$  would be reported as, respectively,

$$x_{1} = \begin{bmatrix} 0\\ 0\\ 1 \end{bmatrix}, x_{2} = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}}i \\ 0 \end{bmatrix}, x_{3} = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}}i \\ 0 \end{bmatrix}.$$

## 10.3.2 The Eigendecomposition

The n eigenvalues  $\lambda_j$  of a  $n \times n$  matrix A and their corresponding eigenvectors  $x_j$  satisfy the equations

$$Ax_i = \lambda_i x_i, \ i = j, 2, ..., n \,.$$

These n equations can be summarized into a single equation as

$$A\begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix} = \begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_n \end{bmatrix}$$

or  $AS = S\Lambda$ . Assuming that S is invertible, we get the eigendecomposition

$$A = S\Lambda S^{-1} \,. \tag{10.23}$$

**Example 10.11** The eigendecomposition of matrix A in Example 10.10 is

$$\underbrace{\begin{bmatrix} 2 & 1 & 0 \\ -1 & 2 & 0 \\ 0 & 0 & 3 \end{bmatrix}}_{A} = \underbrace{\begin{bmatrix} u & v & 0 \\ -ui & vi & 0 \\ 0 & 0 & s \end{bmatrix}}_{S} \underbrace{\begin{bmatrix} 2-i & 0 & 0 \\ 0 & 2+i & 0 \\ 0 & 0 & 3 \end{bmatrix}}_{\Lambda} \underbrace{\begin{bmatrix} \frac{1}{2u} & -\frac{1}{2ui} & 0 \\ \frac{1}{2v} & \frac{1}{2vi} & 0 \\ 0 & 0 & \frac{1}{s} \end{bmatrix}}_{S-1}.$$

You can verify on your own that  $S^{-1}$  is indeed the inverse of S. Note that you can arrange the eigenvectors in S in any order as long as the corresponding eigenvalues are placed in  $\Lambda$  from top left to botton right in the same order.

In forming (10.23), we assumed that S is non-singular, which means that the columns of S, i.e., the eigenvectors  $x_1, x_2, \ldots, x_n$ , are independent. It turns out that for some matrices we cannot find n independent eigenvectors. Such matrices cannot be factorized as (10.23). They are not diagonalizable.

Non-diagonalizability may occur if there are repeated eigenvalues. It can be shown that if all of the eigenvalues of a matrix have distinct values, then there are n independent eigenvectors, one per eigenvalue. All matrices that have n distinct eigenvalues (no repeats) are therefore diagonalizable. If there are repeated eigenvalues, then diagonalizability depends on the

"multiplicities" of the eigenvalues. If r eigenvalues share the same value (we call r the **algebraic multiplicity** of the eigenvalue), then diagonalizability requires that the eigenspace associated with that eigenvalue have dimension r (we call this the **geometric multiplicity** of the eigenvalue). This then allows us to choose r independent eigenvectors from the eigenspace for the r eigenvalues with the same value. If the geometric multiplicity of the repeated eigenvalue is lower than its algebraic multiplicity, the matrix is not diagonalizable.

**Example 10.12** You can easily show that the matrix

	Γ3	0	[0		[3	1	0
A =	0	3	0	and $B =$	= 0	3	0
	0	0	1		0	0	1

both have eigenvalues  $\lambda_1 = 3$ ,  $\lambda_2 = 3$ , and  $\lambda_3 = 1$ . The eigenvalue  $\lambda = 3$  has algebraic multiplicity 2 for both A and B. Consider first matrix A. For  $\lambda_3 = 1$ , we have

$$(A - \lambda_3 I_3) x_3 = \begin{bmatrix} 2 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} x_{31} \\ x_{32} \\ x_{33} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

or  $x_{31} = 0$ ,  $x_{32} = 0$ , with  $x_{33}$  taking any value. That is, the eigenvectors are all vectors of the form

$$x_3 = \begin{bmatrix} x_{31} \\ x_{32} \\ x_{33} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ s \end{bmatrix}.$$

For  $\lambda_1 = \lambda_2 = 3$ , we have

$$(\lambda_i I_3 - A) x_i = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 2 \end{bmatrix} \begin{bmatrix} x_{i1} \\ x_{i2} \\ x_{i3} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \ i = 1, 2,$$

which says that  $x_{i3} = 0$ , and  $x_{i1}$  and  $x_{i2}$  can take any values. The eigenvectors corresponding to  $\lambda_1 = 3$  make up the 2-dimensional subspace comprising vectors of the form

$$x_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ x_{i3} \end{bmatrix} = \begin{bmatrix} s_1 \\ s_2 \\ 0 \end{bmatrix}, \ i = 1, 2.$$

Because the geometric multiplicity of this eigenvalue is equal to its algebraic multiplicity, we can find two independent eigenvectors for  $\lambda_1 = \lambda_2 = 3$ , e.g.,

$$x_1 = \begin{bmatrix} q \\ 0 \\ 0 \end{bmatrix}$$
 and  $x_2 = \begin{bmatrix} 0 \\ r \\ 0 \end{bmatrix}$ ,  $q, r \neq 0$ .

We have the eigendecomposition

 $A = \begin{bmatrix} 3 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} q & 0 & 0 \\ 0 & r & 0 \\ 0 & 0 & s \end{bmatrix} \begin{bmatrix} 3 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1/q & 0 & 0 \\ 0 & 1/r & 0 \\ 0 & 0 & 1/s \end{bmatrix}.$ 

For matrix B, the eigenvectors  $x_3$  corresponding to  $\lambda_3 = 1$  are the same as for matrix A. However, for  $\lambda_1 = \lambda_2 = 3$ , we have

$$(\lambda_i I_3 - B) x_i = \begin{bmatrix} 0 & -1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 2 \end{bmatrix} \begin{bmatrix} x_{i1} \\ x_{i2} \\ x_{i3} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \ i = 1, 2.$$

which says that  $x_{i2} = 0$  and  $-2x_{i3} = 0$ , or  $x_{i3} = 0$ . Only  $x_{i1}$  may take a free parameter. The eigenvectors corresponding to the eigenvalues  $\lambda_1 = \lambda_2 = 3$  all take the form

$$x_i = \begin{bmatrix} 0 \\ q \\ 0 \end{bmatrix}, i = 1, 2.$$

This space of eigenvectors only has dimension 1. The geometric multiplicity is less than the algebraic multiplicity for  $\lambda_1 = \lambda_2 = 3$ , so we cannot find two independent eigenvectors for these eigenvalues. This matrix is not diagonalizable.

# 10.3.3 Eigenvalues and Matrix Properties

10.3.3.1 Eigenvalues and the Determinant The relationship between the determinant of a matrix and its eigenvalues is easy to see. Denoting  $\det(\lambda I_n - A)$  by  $\rho(\lambda)$ , we have  $\rho(0) = \det(-A) = (-1)^n \det(A)$ . From

$$\rho(\lambda) = \det(\lambda I_n - A) = (\lambda - \lambda_1)(\lambda - \lambda_2) \times \dots \times (\lambda - \lambda_n)$$
(10.24)

we get  $\rho(0) = (-1)^n \lambda_1 \lambda_2 \times \cdots \times \lambda_n$ . It follows that

$$\det(A) = \lambda_1 \lambda_2 \times \dots \times \lambda_n \,. \tag{10.25}$$

A corollary of result (10.25) is that A is has an inverse if and only if none of its eigenvalues is zero.

10.3.3.2 Eigenvalues and Trace To show that the eigenvalues sum to the trace, we focus on the coefficient  $\zeta_{n-1}$  in the polynomial

$$\rho(\lambda) = \det(\lambda I_n - A) = \lambda^n + \zeta_{n-1}\lambda^{n-1} + \dots + \zeta_1\lambda + \zeta_0$$
(10.26)

409

#### Mathematics and Programming for the Quantitative Economist

The Laplace expansion of  $\rho(\lambda) = \det(\lambda I_n - A)$  along the first column is

$$\rho(\lambda) = \begin{vmatrix} \lambda - a_{11} & a_{12} & a_{13} & \dots & a_{1n} \\ a_{21} & \lambda - a_{22} & a_{13} & \dots & a_{1n} \\ a_{31} & a_{32} & \lambda - a_{33} & \dots & a_{1n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & a_{n3} & \dots & \lambda - a_{nn} \end{vmatrix}$$

$$= (\lambda - a_{11})C_{11} + a_{21}C_{21} + \dots + a_{n1}C_{n1}$$
(10.27)

where  $C_{kj}$  is the (k, j)th cofactor of  $\lambda I_n - A$ . You can fully work this out to get all the coefficients  $\zeta_0, \zeta_1, \ldots, \zeta_{n-1}$  in (10.26), but we shall note only that the coefficient on  $\lambda^n$  is 1, and the coefficient  $\zeta_{n-1}$  on  $\lambda^{n-1}$  is -trace(A). To see the latter, note that in the column 1 expansion in (10.27), only the first term  $(\lambda - a_{11})C_{11}$  contains  $\lambda^{n-1}$ . This is because in computing the (k, 1)th cofactor, we remove the kth row and 1st column of  $\lambda I_n - A$ . For k > 1, this removes both the column containing  $\lambda - a_{11}$  and the row containing  $\lambda - a_{kk}$ . The highest power of  $\lambda$  in  $C_{k1}, k > 1$ , will be at most n - 2. We can therefore write  $\rho(\lambda)$  as

$$\rho(\lambda) = (\lambda - a_{11})C_{11} + \text{ order } (n-2) \text{ polynomial in } \lambda.$$
 (10.28)

Repeating this argument as we further expand  $C_{11}$  along the first columns, we get

$$\rho(\lambda) = \prod_{k=1}^{n} (\lambda - a_{kk}) + \text{ order } (n-2) \text{ polynomial in } \lambda.$$
 (10.29)

The  $\lambda^{n-1}$  terms in (10.29) arise only in  $\prod_{i=1}^n (\lambda - a_{kk})$ , and only when the  $a_{kk}$  in each of the  $(\lambda - a_{kk})$  terms is multiplied with the  $\lambda$  in the other n-1  $(\lambda - a_{kk})$  terms. Consolidating these  $\lambda^{n-1}$  terms over  $k = 1, 2, \ldots, n$ , we find that the coefficient on  $\lambda^{n-1}$  to be

$$\zeta_{n-1} = -(a_{11} + a_{22} + \dots + a_{nn}) = -\text{trace}(A).$$

Applying the same argument to (10.24), we find the coefficient of  $\lambda^{n-1}$  there to be  $-(\lambda_1 + \lambda_2 + \dots + \lambda_n)$ . It follows that

$$\operatorname{trace}(A) = \lambda_1 + \lambda_2 + \dots + \lambda_n$$
.

10.3.3.3 Eigenvalues and Rank We know that a full rank  $n \times n$  matrix will have no zero-valued eigenvalues, and the presence of zero-valued eigenvalues implies a rank-deficient matrix, but this result alone does not tell us the exact rank of a rank-deficient matrix. However, if a square matrix is diagonalizable, then  $A = S\Lambda S^{-1}$  where S is full rank. This means that

 $\operatorname{rank}(A) = \operatorname{rank}(\Lambda)$ . Since  $\Lambda$  is diagonal, its rank is just the number of nonzero elements along the diagonal. In other words, if A is diagonalizable, then rank(A) is equal to the number of non-zero eigenvalues it possesses.

We emphasize that this result only holds for diagonalizable matrices. The non-diagonalizable matrix

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

has rank 1, but has no non-zero eigenvalues.

10.3.3.4 Eigendecomposition of Symmetric Matrices Symmetric matrices have the convenient property that (i) they are *always* diagonalizable, and (ii) they always have real-valued eigenvalues. Furthermore, (iii) we can always construct a matrix of (real-valued) eigenvectors that are not just invertible, but also orthonormal. That is, we can write any given symmetric matrix A as

$$A = Q \Lambda Q^{\mathrm{T}}$$
 where  $Q^{\mathrm{T}} Q = I_n$ 

and where the diagonal elements of  $\Lambda$  are real. This eigendecomposition is also called the Spectral Decomposition of A.

We omit proofs of these results, and make do with the  $2 \times 2$  case as an illustration. Suppose

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

You can easily show, by finding the roots of  $det(A - \lambda I_2) = 0$ , that the eigenvalues of A are

$$\lambda_{1,2} = \frac{a+b\pm\sqrt{(a+b)^2-4(ab-c^2)}}{2} = \frac{a+b\pm\sqrt{(a-b)^2+4c^2}}{2}$$

Since  $(a-b)^2+4c^2 \ge 0$ ,  $\lambda_1$  and  $\lambda_2$  are always real-valued. If  $(a-b)^2+4c^2 > 0$ , then  $\lambda_1 \neq \lambda_2$ . Suppose this is the case, and suppose  $x_1$  and  $x_2$  are the corresponding eigenvectors. We have  $Ax_1 = \lambda_1 x_1$  and  $Ax_2 = \lambda_2 x_2$ . Then

$$(\lambda_1 x_1)^{\mathrm{T}} x_2 = (A x_1)^{\mathrm{T}} x_2 = x_1^{\mathrm{T}} A^{\mathrm{T}} x_2 = x_1^{\mathrm{T}} A x_2 = x_1^{\mathrm{T}} (\lambda_2 x_2).$$

Rearranging, we get  $(\lambda_1 - \lambda_2)x_1^{\mathrm{T}}x_2 = 0$ . Since  $\lambda_1 \neq \lambda_2$ , it must be that  $x_1^{\mathrm{T}}x_2 = 0$ , which says that  $x_1$  and  $x_2$  are orthogonal. Then the 2×2 matrix

$$Q = \begin{bmatrix} x_1 & x_2 \\ \|x_1\| & \|x_2\| \end{bmatrix}$$

is the orthogonal matrix such that  $A=Q\Lambda Q^{\rm T}.$  If  $\lambda_1=\lambda_2,$  then it must be  $(a-b)^2+4c^2=0,$  which is possible only if c = 0 and a = b. Then the matrix is

$$A = \begin{bmatrix} a & 0\\ 0 & a \end{bmatrix}$$

with eigenvalues  $\lambda_1 = \lambda_2 = a$ . The equation  $(A - \lambda_i I_n)x_i = 0$  reduces to  $0x_i = 0$ , which says that any vector in  $\mathbb{R}^2$  is an eigenvector, and we are free to select two orthogonal eigenvectors from this space (e.g., select  $x_1 = \begin{bmatrix} 1 & 0 \end{bmatrix}^T$  and  $x_2 = \begin{bmatrix} 0 & 1 \end{bmatrix}^T$ ), and set  $Q = \begin{bmatrix} x_1 & x_2 \end{bmatrix}$  which is certainly orthogonal.

**Example 10.13** The eigendecomposition of symmetric matrices gives us some insight into the transformation made to a vector x when pre-multiplied by a symmetric matrix A. This insight may be helpful to you when thinking about principal components, which we will discuss shortly.

Suppose A is symmetric  $n \times n$  with eigendecomposition

$$A = Q\Lambda Q^{\mathrm{T}}$$

where  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$  is the diagonal matrix of eigenvalues, and Q is an orthogonal matrix whose columns are the corresponding eigenvectors

$$Q = \begin{bmatrix} q_1 & q_2 & \dots & q_n \end{bmatrix}$$
 where  $Aq_i = \lambda_i q_i$ ,  $i = 1, \dots, n$ .

The fact that Q is non-singular (since  $Q^{\mathrm{T}}Q = I_n$ ) means that it has full rank, and its columns span the entire space  $\mathbb{R}^n$ . In other words, the eigenvectors in Q form an orthonormal basis for  $\mathbb{R}^n$ . Any vector  $x \in \mathbb{R}^n$  can therefore be written as a linear combination of the eigenvectors, i.e., for every  $x \in \mathbb{R}^n$ , there is a  $c \in \mathbb{R}^n$  such that

$$x = Qc = \begin{bmatrix} q_1 & q_2 & \dots & q_n \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_n \end{bmatrix} = c_1 q_1 + c_2 q_2 + \dots + c_n q_n \,. \tag{10.30}$$

You can think of the eigenvectors as forming an alternative orthonormal coordinate axes for  $\mathbb{R}^n$ . The orthogonal basis associated with the usual Cartesian coordinate system is the orthonormal set of vectors  $\{e_1, e_2, \ldots, e_n\}$  where  $e_j$  is the vector whose *j*th term is 1, and all other terms 0. The vector  $x \in \mathbb{R}^n$  can be written as

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = x_1 e_1 + x_2 e_2 + \dots + x_n e_n \,.$$

The vectors in Q — the eigenvectors of A — form another orthonormal coordinate axes, under which the "address" of the point x is c. We can obtain x as the linear combination in (10.30).

Consider now the product Ax. We have

$$\begin{split} Ax &= Q\Lambda Q^{\mathrm{T}}Qc = Q\Lambda c \\ &= \begin{bmatrix} q_1 & q_2 & \dots & q_n \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_n \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_n \end{bmatrix} \end{split}$$

$$=\lambda_1c_1q_1+\lambda_2c_2q_2+\cdots+\lambda_nc_nq_n\,.$$

The weights on the eigenvectors in the linear combination (10.30) that forms x are multiplied by the corresponding eigenvalues. In other words, the product Ax, where A is symmetric, transforms the vector x by stretching/compressing the vector along the axes formed by the eigenvectors, by factors equal to the corresponding eigenvalues.

Similar arguments can be made to describe the effect of premultiplying x by any diagonalizable matrix, except that the eigenvectors need not be orthogonal, and you may have to work with complex numbers!

10.3.3.5 Eigenvalues and Positive Definiteness A symmetric  $n \times n$  matrix A is positive definite if  $c^{T}Ac > 0$  for all n-vectors  $c \neq 0_{n}$ . It is positive semidefinite if  $c^{T}Ac \geq 0$  for all n-vectors  $c \neq 0_{n}$ . From the eigendecomposition  $A = QAQ^{T}$ , define

$$d = \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_n \end{bmatrix} = Q^{\mathrm{T}} c$$

Since  $Q^{\mathrm{T}}$  is full rank,  $c \neq 0_n \Leftrightarrow d \neq 0_n$ . Then

$$c^{\mathrm{T}}Ac = c^{\mathrm{T}}Q\Lambda Q^{T}c = d^{\mathrm{T}}\Lambda d = \lambda_{1}d_{1}^{2} + \dots + \lambda_{n}d_{n}^{2}.$$

It follows that A is positive definite if and only if  $\lambda_k > 0$  for all k = 1, 2, ..., n. Suppose  $\lambda_k > 0$  for all k = 1, 2, ..., n. Then

$$c \neq 0_n \ \Rightarrow \ d \neq 0_n \ \Rightarrow \ d_k^2 > 0 \ \text{for at least one} \ k = 1, 2, \dots, n \ \Rightarrow \ c^{\mathrm{T}} A c > 0 \,.$$

If  $\lambda_k \leq 0$  for some k = 1, 2, ..., n, then we can find  $d \neq 0_n$ , and therefore  $c \neq 0_n$ , such that  $c^{\mathrm{T}}Ac \leq 0$ . A similar argument shows that A is positive semidefinite if and only if  $\lambda_k \geq 0$  for all k = 1, 2, ..., n.

10.3.3.6 *Eigenvalues and Rank, Again* Since all symmetric matrices are diagonalizable, and the rank of a diagonalizable matrix is equal to its number of non-zero eigenvalues, we have the convenient property that the

rank of a symmetric matrix is equal to the number of non-zero eigenvalues that it possesess. If matrix A is non-diagonalizable, we can make use of the fact that  $A^{T}A$  is symmetric (and therefore diagonalizable), and rank $(A) = \operatorname{rank}(A^{T}A)$ , to get the result that the rank of A is equal to the number of non-zero eigenvalues in  $A^{T}A$ .

10.3.3.7 *Further Facts about Eigenvalues* We state a few additional facts about eigenvalues (and eigenvectors) that you are asked to prove as exercises. These apply to all matrices regardless of symmetry.

(a) The transpose of a matrix has the same eigenvalues as the matrix itself.

(b) If A is  $n \times n$  and non-singular with eigenvalues  $\lambda_k$ , k = 1, 2, ..., n, then the eigenvalues of  $A^{-1}$  are  $1/\lambda_k$ , k = 1, 2, ..., n.

(c) The eigenvalues of idempotent matrices (matrices such that AA = A) take values 1 and 0 only. An implication of this is that the trace of an idempotent matrix is equal to the number of non-zero eigenvalues it possesses. Projection matrices are examples of idempotent matrices.

(d) Although idempotent matrices need not be symmetric, e.g.,

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

is idempotent but not symmetric, idempotent matrices are always diagonalizable.

(e) The rank of an idempotent matrix is equal to its trace.

Example 10.14 The eigenvalues of

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

are  $\lambda_1 = 1 + 2\sqrt{2}$ ,  $\lambda_2 = 1$  and  $\lambda_3 = 1 - 2\sqrt{2}$ , with corresponding eigenvectors

$$x_1 = \begin{bmatrix} s \\ \sqrt{2}s \\ -s \end{bmatrix}, x_2 = \begin{bmatrix} s \\ 0 \\ s \end{bmatrix}, x_3 = \begin{bmatrix} s \\ -\sqrt{2}s \\ -s \end{bmatrix}.$$

You can easily verify that the product of the eigenvalues is equal to the determinant of A, and the sum of the eigenvalues is equal to the trace of A. Since A is symmetric, the eigenvectors should be mutually orthogonal, and you can verify this to be the case. Selecting s = 1 (so that the eigenvectors are real) and then normalizing, we get the specific orthonormal eigenvectors

$$x_{1} = \begin{bmatrix} \frac{1}{2} \\ \frac{1}{\sqrt{2}} \\ -\frac{1}{2} \end{bmatrix}, x_{2} = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ 0 \\ \frac{1}{\sqrt{2}} \end{bmatrix}, x_{3} = \begin{bmatrix} \frac{1}{2} \\ -\frac{1}{\sqrt{2}} \\ -\frac{1}{2} \end{bmatrix}.$$

As an exercise, you should verify (by hand or otherwise) that  $A=Q\Lambda Q^{\rm T}$  where

$$Q = \begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} = \begin{bmatrix} \frac{1}{2} & \frac{1}{\sqrt{2}} & \frac{1}{2} \\ \frac{1}{\sqrt{2}} & 0 & -\frac{1}{\sqrt{2}} \\ -\frac{1}{2} & \frac{1}{\sqrt{2}} & -\frac{1}{2} \end{bmatrix} \text{ and } \Lambda = \begin{bmatrix} 1 + 2\sqrt{2} & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 - 2\sqrt{2} \end{bmatrix}.$$

Note that your statistical package may multiply one or more of these eigenvectors by -1. Since there are no zero eigenvalues, the matrix A has full rank. Since there are eigenvalues of both signs, the matrix A is neither positive semidefinite or negative semidefinite.

### 10.3.4 Principal Component Analysis

One reason why the eigendecompositions of symmetric matrices are of particular interest is that variance matrices are symmetric, and eigenvalues and eigenvectors of variance matrices have useful interpretations.

Suppose x is a k-vector of zero-mean random variables with the  $k \times k$  variance matrix var(x). Suppose  $var(x) = Q\Lambda Q^{\mathrm{T}}$  where

$$\Lambda = \operatorname{diag}(\lambda_1, \lambda_2, \dots, \lambda_k) \text{ and } Q = \begin{bmatrix} q_1 & q_2 & \dots & q_k \end{bmatrix},$$

and where

$$q_j = \begin{bmatrix} q_{1j} \\ q_{2j} \\ \vdots \\ q_{kj} \end{bmatrix}, \ j = 1, 2, \dots, k$$

are the k orthonormal eigenvectors corresponding to the eigenvalues. We will label the eigenvalues in descending order  $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_k \geq 0$  (there are no negative eigenvalues since the variance matrix is positive semidefinite). Let  $y = Q^T x$ , i.e.,

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_k \end{bmatrix} = \begin{bmatrix} q_1^{\mathrm{T}} x \\ q_2^{\mathrm{T}} x \\ \vdots \\ q_k^{\mathrm{T}} x \end{bmatrix} = Q^{\mathrm{T}} x \,,$$

is an k-vector of random variables, each a linear combination of the random variables in x. This random vector y has variance matrix

$$var(y) = var(Q^{\mathrm{T}}x) = Q^{\mathrm{T}}var(x)Q = Q^{\mathrm{T}}Q\Lambda Q^{\mathrm{T}}Q = \Lambda$$
.

That is, the random variables in y are mutually uncorrelated, with variances equal to the corresponding eigenvalues down the diagonal of  $\Lambda$ . Moreover, it is fairly straightforward to show (see next chapter) that

Mathematics and Programming for the Quantitative Economist

- $y_1 = q_1^{\mathrm{T}} x$  is the linear combination of the variables in x with the largest variance subject to  $||q_1|| = 1$ ,
- $y_2 = q_2^{\mathrm{T}} x$  is the linear combination of the variables in x with the next largest variance subject to  $||q_2|| = 1$  and  $q_2^{\mathrm{T}} q_1 = 0$ ,
- $y_3 = q_3^{\mathrm{T}} x$  is the linear combination of the variables in x with the third largest variance subject to  $||q_3|| = 1$ ,  $q_3^{\mathrm{T}}q_1$  and  $q_3^{\mathrm{T}}q_2$ ,

and so on.

This argument carries over to sample variance matrices. Suppose X is an  $n \times k$  matrix with columns containing the n centered observations of each of k variables. "Centered observations" means that the sample means of each of the variable have been removed, i.e.,

$$X = \begin{bmatrix} \tilde{x}_{*1} & \tilde{x}_{*2} & \dots & \tilde{x}_{*k} \end{bmatrix} \text{ where } \tilde{x}_{*j} = \begin{bmatrix} x_{1j} - \overline{x}_j \\ x_{2j} - \overline{x}_j \\ \vdots \\ x_{nj} - \overline{x}_j \end{bmatrix} = \begin{bmatrix} \tilde{x}_{1j} \\ \tilde{x}_{2j} \\ \vdots \\ \tilde{x}_{nj} \end{bmatrix}$$

and where  $\overline{x}_j = (1/n) \sum_i^n x_{ij}$  is the sample mean of the variable j observations. Then

$$\frac{1}{n-1}X^{\mathrm{T}}X = \begin{bmatrix} \frac{1}{n-1}\sum_{i=1}^{n} \tilde{x}_{i1}^{2} & \frac{1}{n-1}\sum_{i=1}^{n} \tilde{x}_{i1}\tilde{x}_{i2} & \dots & \frac{1}{n-1}\sum_{i=1}^{n} \tilde{x}_{i1}\tilde{x}_{ik} \\ \frac{1}{n-1}\sum_{i=1}^{n} \tilde{x}_{i2}\tilde{x}_{i1} & \frac{1}{n-1}\sum_{i=1}^{n} \tilde{x}_{i2}^{2} & \dots & \frac{1}{n-1}\sum_{i=1}^{n} \tilde{x}_{i1}\tilde{x}_{ik} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{1}{n-1}\sum_{i=1}^{n} \tilde{x}_{in}\tilde{x}_{i1} & \frac{1}{n-1}\sum_{i=1}^{n} \tilde{x}_{in}\tilde{x}_{i2} & \dots & \frac{1}{n-1}\sum_{i=1}^{n} \tilde{x}_{in}^{2} \end{bmatrix}$$

contains all of the sample variances and sample covariances of the variables. If  $\frac{1}{n-1}X^{\mathrm{T}}X = Q\Lambda Q^{\mathrm{T}}$  where the diagonal elements of the diagonal matrix  $\Lambda$  are the eigenvalues of  $\frac{1}{n-1}X^{\mathrm{T}}X$  and the columns of Q are the corresponding orthonormal eigenvectors, then Y where

$$Y = XQ$$

i.e., 
$$\begin{bmatrix} y_1 & y_2 & \dots & y_k \end{bmatrix} = X \begin{bmatrix} q_1 & q_2 & \dots & q_k \end{bmatrix}$$
  
=  $\begin{bmatrix} Xq_1 & Xq_2 & \dots & Xq_k \end{bmatrix}$ 

is a  $n \times k$  matrix containing observations on k variables, each of which a linear combination of the observations of the k variables in X, and whose sample variance matrix is

$$\frac{1}{n-1}Y^{\mathrm{T}}Y = \frac{1}{n-1}Q^{\mathrm{T}}X^{\mathrm{T}}XQ = Q^{\mathrm{T}}\left(\left(\frac{1}{n-1}\right)X^{\mathrm{T}}X\right)Q = Q^{\mathrm{T}}Q\Lambda Q^{\mathrm{T}}Q = \Lambda$$

That is, the columns of Y are mutually uncorrelated indices each computed as a linear combination of the columns of X. Furthermore, the first column  $y_1$  is the index that contains the most variation across observations i = 1, 2, ..., n, the second column contains the second most variation, and so on. We call these indices the **principal components** of X.<sup>6</sup>

**Example 10.15** Principal Component Analysis (PCA) can help us understand the way observational units differ in a data set. Take for example the  $51 \times 14$  data matrix in causes-of-death-by-state.csv, which contains the age-adjusted number of deaths per 100,000, all races, both sexes, all ages, over the period 2016-2020, across the 51 US states plus District of Columbia (rows) and 14 causes of death (columns): accidents & adverse effects (accident), Alzheimer's disease (Alzheimers), cancer, cerebrovascular diseases (cerebrovascular), chronic lower respiratory disease (respiratory), chronic liver disease & cirrhosis (liver), diabetes mellitus (diabetes), heart disease (heart), momicide & legal intervention (homicide), influenza, kidney disease - nephritis & nephrosis (kidney), pneumonia, septicemia, suicide & self-inflicted injury (suicide).<sup>7</sup>

Suppose we want to create a few "causes of death" indices that capture meaningful differences in causes of death across states. One way is to group the causes into interpretable groups, perhaps vascular diseases (heart and cerebrovascular), trauma-related causes (accident, homicide, suicide), and disease-related (all others) and create an index for each group. Another way is to use PCA to create uncorrelated cause-of-death indices that measure the ways the states differ most in terms of causes of death.

The PCA method proceeds as follows:

First we center the data matrix so that each column has zero sample mean. We sometimes also scale the data so that the columns have unit sample variance. This is because PCA quite naturally places more weight on variables with larger sample variances. In our example, we do not scale the columns because the data is age-adjusted number of deaths per 100,000 people, so variation due to different age profiles and populations across states have already been accounted for. The remaining differences in variation across states per disease are meaningful for our purpose, so we do not further standardize the data. Let X denote the centered data set.

Second, we find the eigendecomposition of the sample variance matrix

$$(1/n)X^{\mathrm{T}}X = Q\Lambda Q^{\mathrm{T}}$$

$$\frac{1}{n-1}Y^{\mathrm{T}}Y = \frac{1}{n-1}Q^{\mathrm{T}}X^{\mathrm{T}}XQ = \frac{1}{n-1}Q^{\mathrm{T}}Q\Lambda Q^{\mathrm{T}}Q = \frac{1}{n-1}\Lambda\,.$$

 $<sup>^6\</sup>mathrm{If}$  instead of the eigendecomposition of  $\frac{1}{n-1}X^\mathrm{T}X$  we used the eigendecomposition  $X^\mathrm{T}X=Q\Lambda Q^\mathrm{T},$  then

The principal components, i.e., the columns of XQ do not change, but the diagonal elements of  $\Lambda$  now contain n-1 times the sample variance of the principal components. <sup>7</sup>Compiled from https://hdpulse.nimhd.nih.gov.

and calculate the k principal components

$$Y = XQ$$
.

The first column of Y is a linear combination of all of the columns of X with weights ("loadings") found down the first column of Q. It is the 1st principal component of the data, a vector containing the "cause of death" index score for each state that accounts for the most variation in the data across states. The second column of Y contains the 2nd principal component, and so on. The eigenvalues are the sample variances of each of the principal components.

We find that the first two principal components account for just over 83 percent of the total variation in the data (defined as the sum of variances of all of the principal components). The first three principal components accounts for just under 90 percent of the total variation. Figure xx(a) below plots the loadings per disease on the first two principal components. We see that the first principal component emphasizes death due to heart disease whereas the second principal component emphasizes accidents. Both place some weight on extent cancer and respiratory diseases.<sup>8</sup> All other diseases are given much lower weights. Figure xx(b) plots PC1 and PC2 scores for each state. A group of states which includes many of the "southern states" (Mississippi, Oklahoma, Alabama, Arkansas, Louisiana, Tennessee, Kentucky, West Virginia) stand out for high PC1 scores. West Virginia stands out for also having a high PC2 score.

PCA sometimes produces indices that are hard to interpret. It nonetheless adds a useful data analytic perspective, and can also serve as a dimension reduction technique. Instead on focusing of all 14 causes of death, we can get insights on most of the variation in the data by examining just two or three indices.



<sup>&</sup>lt;sup>8</sup>The second principal component places substantial weight on heart disease, but opposite in sign to accident, cancer and respiratory. We can interpret this as an "adjustment" for heart disease, which was accounted for in the first principal component.

# 10.3.5 The Singular Value Decomposition

We come finally to the **singular value decomposition** (SVD) which is an eigendecomposition that applies to all matrices, regardless of "shape" or rank. Let A be a  $n \times k$  matrix of rank  $r \leq \min\{n, k\}$ . The  $k \times k$  matrix  $A^{\mathrm{T}}A$  is symmetric, positive semidefinite and rank r and so it has r positive eigenvalues  $\sigma_1^2 \geq \sigma_2^2 \geq \cdots \geq \sigma_r^2 > 0$  and k - r zero-valued eigenvalues  $\sigma_{r+1}^2 = \ldots \sigma_k^2 = 0$ . Let  $v_1, v_2, \ldots v_r$  be the eigenvectors associated with the positive eigenvalues

$$A^{\mathrm{T}}Av_{j} = \sigma_{j}^{2}v_{j}, \ j = 1, 2, \dots, r.$$
(10.31)

Let  $v_{r+1},\ldots,v_k$  be the eigenvectors corresponding to the zero-valued eigenvalues. These satisfy

$$A^{\mathrm{T}}Av_{j} = 0_{k}, \ j = r+1, \dots, k.$$
 (10.32)

The set of eigenvectors  $v_j, \; j=1,2,\ldots,k$  are mutually orthonormal. We note that (10.31) implies

$$\|Av_{j}\|^{2} = (Av_{j})^{\mathrm{T}}(Av_{j}) = v_{j}^{\mathrm{T}}A^{\mathrm{T}}Av_{j} = \sigma_{j}^{2}v_{j}^{\mathrm{T}}v_{j} = \sigma_{j}^{2}$$
(10.33)

for  $j = 1, 2, \dots, r$  whereas (10.32) implies

$$||Av_j||^2 = (Av_j)^{\mathrm{T}}(Av_j) = v_j^{\mathrm{T}}A^{\mathrm{T}}Av_j = 0 \text{ or } Av_j = 0$$
(10.34)

for j = r + 1, ..., k. The set of all vectors v such that Av = 0 forms the dimension k - r null space of A, and  $\{v_{r+1}, ..., v_k\}$  are a set of k - r orthonormal vectors satisfying (10.34). Therefore,  $\{v_{r+1}, ..., v_k\}$  forms an orthonormal basis for N(A, k - r). It follows from the Fundamental Theorem of Linear Algebra that  $\{v_1, ..., v_r\}$  forms an orthonormal bases for  $C(A^{\mathrm{T}})$ .

The  $n \times n$  matrix  $AA^{T}$  is also symmetric, positive semidefinite and rank r so it too has r positive eigenvalues, plus n - r zero-valued eigenvalues. For j = 1, 2, ..., r, (10.31) implies

$$AA^{\mathrm{T}}Av_{j} = \sigma_{j}^{2}Av_{j}. \qquad (10.35)$$

Dividing throughout by  $\sigma_i$ , we can write

$$AA^{\mathrm{T}}\left(\frac{Av_{j}}{\sigma_{j}}\right) = \sigma_{j}^{2}\left(\frac{Av_{j}}{\sigma_{j}}\right), \ j = 1, 2, \dots, r, \qquad (10.36)$$

which we can write as

$$AA^{\mathrm{T}}u_{j} = \sigma_{j}^{2}u_{j}$$
 where  $u_{j} = \frac{Av_{j}}{\sigma_{j}}$ ,  $j = 1, 2, ..., r$ . (10.37)

This says that  $\sigma_j$  are also the non-zero eigenvalues of  $AA^{\mathrm{T}}$ , and  $u_j$  as defined in (10.37) are orthonormal eigenvectors associated with these eigenvalues. The  $u_j$  vectors, j = 1, 2, ..., r, are unit vectors by definition, and they are orthogonal because

$$u_j^{\mathrm{T}} u_h = \frac{1}{\sigma_j} \frac{1}{\sigma_h} v_j^{\mathrm{T}} A^{\mathrm{T}} A v_h = \frac{1}{\sigma_j} \frac{1}{\sigma_h} v_j^{\mathrm{T}} (A^{\mathrm{T}} A v_h) = \frac{1}{\sigma_j} \frac{1}{\sigma_h} v_j^{\mathrm{T}} (\sigma_h^2 v_h) = 0$$

for all  $j \neq h, j, h = 1, 2, ..., r$ . The remaining eigenvalues of  $AA^{T}$  are zero, and the associated set of eigenvectors comprises vectors u such that

$$AA^{\mathrm{T}}u = 0$$

This implies  $u^{\mathrm{T}}AA^{\mathrm{T}}u = (A^{\mathrm{T}}u)^{\mathrm{T}}(A^{\mathrm{T}}u) = 0$ , or  $A^{\mathrm{T}}u = 0$ . These eigenvectors form the null space  $N(A^{\mathrm{T}}, n-r)$ , and we can select n-k orthonormal eigenvectors  $\{u_{r+1}, \ldots, u_n\}$  from this space. These form an orthonormal basis for  $N(A^{\mathrm{T}}, n-r)$ , and are all orthogonal to  $\{u_1, \ldots, u_r\}$ , which form an orthonormal basis for C(A, r).

Finally, from the definition of  $u_j$  in (10.37), we have

$$Av_{j} = \sigma_{j}u_{j}, \, j = 1, 2, \dots, r \tag{10.38}$$

which we can write in matrix form as

$$AV_{1} = A \begin{bmatrix} v_{1} & \dots & v_{r} \end{bmatrix} = \begin{bmatrix} u_{1} & \dots & u_{r} \end{bmatrix} \begin{bmatrix} \sigma_{1} & 0 & \dots & 0 \\ 0 & \sigma_{2} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_{r} \end{bmatrix}$$
(10.39)  
=  $U_{1}\Sigma_{1}$ .

Since  $Av_i = 0$  for j = r + 1, ..., k, we can extend (10.39) to

$$\begin{split} & \underset{n \times k}{\overset{A}{\underset{k \times k}{}}} = A \begin{bmatrix} V_1 & V_2 \end{bmatrix} \\ & = \begin{bmatrix} U_1 & U_2 \end{bmatrix} \begin{bmatrix} \Sigma_r & 0_{r \times (k-r)} \\ 0_{(n-r) \times r} & 0_{(n-r) \times (k-r)} \end{bmatrix} = \underset{n \times n}{\overset{\Sigma}{\underset{n \times k}{}}} \end{split}$$

where  $V_2 = \begin{bmatrix} v_{r+1} & \dots & v_k \end{bmatrix}$  and  $U_2 = \begin{bmatrix} u_{r+1} & \dots & u_n \end{bmatrix}$ .

All this gives the singular value decomposition for a  $n \times k$  matrix A:

**Theorem 10.1** (Singular Value Decomposition) All  $n \times k$  matrices of rank  $r \leq \min\{n, k\}$  can be written as

$$A = U\Sigma V^{\mathrm{T}} \tag{10.40}$$

where

(i) the  $n \times k$  matrix  $\Sigma$  is "diagonal" with non-zero elements  $(\Sigma)_{jj} = \sigma_j$ , j = 1, 2, ..., r, and 0 everywhere else. The  $\sigma_j$  are the square roots of the r (common) positive eigenvalues of the  $A^{\mathrm{T}}A$  and  $AA^{\mathrm{T}}$  and are called the **singular values** of A. We usually label the singular values such that  $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_r$ .

(ii) The columns of the  $n \times n$  matrix U are the orthonormal eigenvectors of  $AA^{\mathrm{T}}$ . The first r columns are the eigenvectors corresponding to the singular values, and form an orthonormal basis for C(A, r). The remaining n - r columns are the eigenvectors corresponding to the zero eigenvectors of  $AA^{\mathrm{T}}$  and form an orthonormal basis for  $N(A^{\mathrm{T}}, n - r)$ .

(iii) The columns of the  $k \times k$  matrix V are the orthonormal eigenvectors of  $A^{\mathrm{T}}A$ . The first r columns are the eigenvectors corresponding to the singular values, and form an orthonormal basis for  $C(A^{\mathrm{T}}, r)$ . The remaining k - r columns are the eigenvectors corresponding to the zero eigenvectors of  $A^{\mathrm{T}}A$  and form an orthonormal basis for N(A, k - r).

If A is a  $n \times k$  data matrix containing n centered observations of k variables, then  $\frac{1}{n-1}A^{\mathrm{T}}A$  is the sample variance matrix containing the sample variances of, and sample covariances between the k variable across observational units, whereas  $\frac{1}{k-1}AA^{\mathrm{T}}$  contains the sample variances of, and sample covariance between each of the observational units across variables. The singular value decomposition provides all the relevant information for principal component analysis of both A and  $A^{\mathrm{T}}$ . We have

$$A^{\mathrm{T}}A = V\Sigma^{\mathrm{T}}U^{\mathrm{T}}U\Sigma V^{\mathrm{T}} = V\Sigma^{\mathrm{T}}\Sigma V^{\mathrm{T}}.$$
(10.41)

The matrix  $\Sigma^{\mathrm{T}}\Sigma$  is  $k \times k$  diagonal with  $\sigma_j^2$  down the first r diagonal terms. (10.41) is the eigendecomposition of  $A^{\mathrm{T}}A$ . The columns of V can be used to construct principal components of A. Also

$$AA^{\mathrm{T}} = U\Sigma V^{\mathrm{T}} V\Sigma^{\mathrm{T}} U^{\mathrm{T}} = U\Sigma\Sigma^{\mathrm{T}} U^{\mathrm{T}} .$$
(10.42)

The matrix  $\Sigma^{\mathrm{T}}\Sigma$  is  $n \times n$  diagonal with  $\sigma_j^2$  down the first r diagonal terms. This is the eigendecomposition of  $AA^{\mathrm{T}}$ . The columns of U can be used to construct principal components of  $A^{\mathrm{T}}$ .

Example 10.17 Completion of death-by-causes data matrix

**Example 10.18** The singular value decomposition can also be written as

$$A = \begin{bmatrix} u_1 & u_2 & \dots & u_n \end{bmatrix} \begin{bmatrix} \sigma_1 & 0 & \dots & 0 & 0 & \dots & 0 \\ 0 & \sigma_2 & \dots & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_r & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 & 0 & \dots & 0 \end{bmatrix} \begin{bmatrix} v_1^T \\ v_2^T \\ \vdots \\ v_k^T \end{bmatrix}$$

we get

422

$$A = \sum_{j=1}^{r} \sigma_j u_j v_j^{\mathrm{T}} \,. \tag{10.43}$$

This can be used to obtain "approximations" of large data matrices. For instance, consider a full column rank  $n \times k$  matrix A may contain pixel data for an image. This can be written as  $A = \sum_{j=1}^{k} \sigma_j u_j v_j^{\mathrm{T}}$ . If we do not need a very high resolution image, we can "compress" the image with the image matrix  $A = \sum_{j=1}^{r} \sigma_j u_j v_j^{\mathrm{T}}$  where r might be considerably smaller than k. This reduces the number of data points from nk to r(n + k). In Fig xx we show an image at various levels of compression.

# 10.3.6 Exercises

Ex. 10.8 The matrix

$$A = \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}$$

is a "rotation matrix" which rotates every  $x \in \mathbb{R}^2$  anticlockwise by an angle of  $\pi/4$  or 45°. Find its eigenvalues and the corresponding eigenvectors. Find its eigendecomposition.

**Ex. 10.9** Show that the matrix  $A = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}$  is non-diagonalizable.

Ex. 10.10 Prove the statements in Section 10.3.3.7.



# 10.5 Solutions to Exercises

Ex. 10.1: We have

$$\begin{split} A^{\mathrm{T}}A &= \begin{bmatrix} \cos\theta & \sin\theta \\ -\sin\theta & \cos\theta \end{bmatrix} \begin{bmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{bmatrix} \\ &= \begin{bmatrix} \cos^{2}\theta + \sin^{2}\theta & -\cos\theta\sin\theta + \sin\theta\cos\theta \\ -\sin\theta\cos\theta + \cos\theta\sin\theta & \sin^{2}\theta + \cos^{2}\theta \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}. \end{split}$$

If

$$y = Ax = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} x_1 \cos \theta - x_2 \sin \theta \\ x_1 \sin \theta + x_2 \cos \theta \end{bmatrix}$$

then

$$\begin{split} y\|^2 &= (x_1 \cos \theta - x_2 \sin \theta)^2 + (x_1 \sin \theta + x_2 \cos \theta)^2 \\ &= (x_1^2 + x_2^2)(\sin^2 \theta + \cos^2 \theta) = (x_1^2 + x_2^2) = \|x\|^2 \end{split}$$

Furthermore,

$$\begin{aligned} x \cdot y &= x_1^2 \cos \theta - x_1 x_2 \sin \theta + x_1 x_2 \sin \theta + x_2^2 \cos \theta \\ &= (x_1^2 + x_2^2) \cos \theta = \|x\|^2 \cos \theta = \|x\| \|y\| \cos \theta \,. \end{aligned}$$

This says that the angle between x and y is  $\theta$ .

Ex. 10.2: We use the notation  $\sum_{i}$  for  $\sum_{i=1}^{n}$ . From

$$\hat{\beta} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix} = (X^{\mathrm{T}}X)^{-1}X^{\mathrm{T}}y = \begin{bmatrix} n & \sum_i x_i \\ \sum_i x_i & \sum_i x_i^2 \end{bmatrix}^{-1} \begin{bmatrix} \sum_i y_i \\ \sum_i x_i y_i \end{bmatrix}$$

we have  $\hat{\beta}_0 = \frac{\sum_i x_i^2 \sum_i y_i - \sum x_i \sum x_i y_i}{n \sum_i x_i^2 - (\sum_i x_i)^2}$  and  $\hat{\beta}_1 = \frac{n \sum_i x_i y_i - \sum x_i \sum y_i}{n \sum_i x_i^2 - (\sum_i x_i)^2}$ . Using the fact the  $\sum_i (x_i - \overline{x})(y_i - \overline{y}) = \sum_i x_i y_i - \overline{x}\overline{y}$ , we get with some algebra that  $\hat{\beta}_1 = \frac{\sum_i (x_i - \overline{x})(y_i - \overline{y})}{\sum_i (x_i - \overline{x})^2}$ . Then using this expression for  $\hat{\beta}_1$ , show that  $\overline{y} - \hat{\beta}_1 \overline{x}$ works out to the expression for  $\hat{\beta}_0$  shown above.

Ex. 10.3: Projecting  $\hat{y} = X (X^{\mathrm{T}} X)^{-1} X^{\mathrm{T}} y$  onto the column space of X gives

$$\begin{split} X(X^{\mathrm{T}}X)^{-1}X^{\mathrm{T}}\hat{y} &= X(X^{\mathrm{T}}X)^{-1}X^{\mathrm{T}}\left[X(X^{\mathrm{T}}X)^{-1}X^{\mathrm{T}}y\right] \\ &= X(X^{\mathrm{T}}X)^{-1}X^{\mathrm{T}}y = \hat{y} \end{split}$$

Since  $\hat{y}$  is already in the column space of X, projecting it again onto the column space of X just returns  $\hat{y}$ .

Ex. 10.4: (a) The orthogonal projection of y onto x is

$$\hat{y} = x(x^{\mathrm{T}}x)^{-1}x^{\mathrm{T}}\hat{y} \,.$$

Since  $x^{\mathrm{T}}x$  and  $x^{\mathrm{T}}\hat{y}$  her are, we can write the above as  $\hat{y} = \frac{x^{\mathrm{T}}y}{x^{\mathrm{T}}x}x$ . The orthogonal  $\hat{e}$  can be found by  $\hat{e} = y - \hat{y} = y - \frac{x^{\mathrm{T}}y}{x^{\mathrm{T}}x}x$ .

423

(b) If  $x = i_n$ , then the orthogonal projection becomes

$$\hat{y} = \frac{i_n^{\mathrm{T}} y}{i_n^{\mathrm{T}} i_n} i_n = \frac{\sum_i y_i}{n} i_n = \overline{y} \, i_n = \begin{bmatrix} \overline{y} \\ \overline{y} \\ \vdots \\ \overline{y} \end{bmatrix} \, .$$

Ex. 10.5: If the columns of X are orthogonal, then  $X^T X$  is the diagonal matrix diag  $(x_1 \cdot x_1, x_2 \cdot x_2, \ldots, x_k \cdot x_k)$  where  $x_i$  is the *i*th column of X,  $i = 1, 2, \ldots, k$ . The formula for  $\hat{y}$  becomes

$$\begin{split} \hat{y} &= X \hat{\beta} = X (X^{\mathrm{T}} X)^{-1} X^{\mathrm{T}} y \\ &= \begin{bmatrix} x_1 & x_2 & \dots & x_k \end{bmatrix} \begin{bmatrix} x_1 \cdot x_1 & 0 & \dots & 0 \\ 0 & x_2 \cdot x_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & x_k \cdot x_k \end{bmatrix}^{-1} \begin{bmatrix} x_1^{\mathrm{T}} \\ x_2^{\mathrm{T}} \\ \vdots \\ x_k^{\mathrm{T}} \end{bmatrix} y \\ &= x_1 \frac{x_1 \cdot y}{x_1 \cdot x_1} + x_2 \frac{x_2 \cdot y}{x_2 \cdot x_2} + \dots + x_k \frac{x_k \cdot y}{x_k \cdot x_k} \,. \end{split}$$

If X is orthogonal, then  $X^{\mathrm{T}}X = I$ . Then  $\hat{\beta} = X^{\mathrm{T}}y$ , and the orthogonal projection is simply  $\hat{y} = XX^{\mathrm{T}}y$ .

Ex. 10.6: (a) We have (symmetry)

$$H^{\rm T} = I_n^{\rm T} - \frac{2}{\|v\|^2} (v^{\rm T})^{\rm T} v^{\rm T} = I_n - \frac{2}{\|v\|^2} v v^{\rm T} = H \,.$$

and (orthogonality)

$$\begin{split} H^{\mathrm{T}}H &= HH = \left(I_n - \frac{2}{\|v\|^2}vv^{\mathrm{T}}\right) \left(I_n - \frac{2}{\|v\|^2}vv^{\mathrm{T}}\right) \\ &= I_nI_n - \frac{4}{\|v\|^2}vv^{\mathrm{T}} + \frac{4}{\|v\|^4}v(v^{\mathrm{T}}v)v^{\mathrm{T}} \\ &= I_n - \frac{4}{\|v\|^2}vv^{\mathrm{T}} + \frac{4\|v\|^2}{\|v\|^4}vv^{\mathrm{T}} = I_n. \end{split}$$

(b) Similarly,

$$\begin{split} \widetilde{H}^{\mathrm{T}} &= \begin{bmatrix} I_p^{\mathrm{T}} & 0_{p \times (n-p)} \\ 0_{(n-p) \times p} & H_{n-p}^{\mathrm{T}} \end{bmatrix} = \begin{bmatrix} I_p & 0_{p \times (n-p)} \\ 0_{(n-p) \times p} & H_{n-p} \end{bmatrix} = \widetilde{H} \ . \\ \widetilde{H}^{\mathrm{T}} \widetilde{H} &= \widetilde{H} \widetilde{H} = \begin{bmatrix} I_p & 0_{p \times (n-p)} \\ 0_{(n-p) \times p} & H_{n-p} \end{bmatrix} \begin{bmatrix} I_p & 0_{p \times (n-p)} \\ 0_{(n-p) \times p} & H_{n-p} \end{bmatrix} \\ &= \begin{bmatrix} I_p & 0_{p \times (n-p)} \\ 0_{(n-p) \times p} & H_{n-p} H_{n-p} \end{bmatrix} = \begin{bmatrix} I_p & 0_{p \times (n-p)} \\ 0_{(n-p) \times p} & I_{n-p} \end{bmatrix} . \end{split}$$

Ex. 10.7: (a) Show that  $a_i \cdot a_j \neq 0$  where  $a_i$  and  $a_j$  are different columns of A, by direct multiplication. (b) Show  $Q^{\mathrm{T}}Q = I$  by direct multiplication.

Ex. 10.8: The eigenvalue equation is

$$\det(A - \lambda I_2) = \det \begin{bmatrix} \frac{1}{\sqrt{2}} - \lambda & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} - \lambda \end{bmatrix} = \lambda^2 - \sqrt{2}\lambda + 1 = 0$$

which you can solve for the eigenvalues  $\lambda_1 = \frac{1}{\sqrt{2}} + \frac{i}{\sqrt{2}}$  and  $\lambda_2 = \frac{1}{\sqrt{2}} - \frac{i}{\sqrt{2}}$ . Remark: Note that the eigenvalues are complex conjugates. Furthermore, it

Remark: Note that the eigenvalues are complex conjugates. Furthermore, it is not surprising that there are no real eigenvalues, since every real vector gets rotated by  $\pi/4$ , so no real vectors stay in the same or opposite direction when premultiplied by A.

The eigenvalues associated with  $\lambda_1$  are the complex-valued vectors  $x_1$  such that

$$(A - \lambda_1)x_1 = \begin{bmatrix} -\frac{i}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{i}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} x_{11} \\ x_{12} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

These vectors satisfy  $-ix_{11} = x_{12}$  and  $x_{11} = ix_{12}$ . Both of these equations say the same thing since 1/i = -i. Letting  $x_{11} = s$  we get

$$x_1 = \begin{bmatrix} s \\ -i \, s \end{bmatrix}, \, s \in \mathbb{C} \, .$$

For  $\lambda_2 = \frac{1}{\sqrt{2}} - \frac{i}{\sqrt{2}}$ , we have

$$(A - \lambda_2)x_2 = \begin{bmatrix} \frac{i}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{i}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} x_{21} \\ x_{22} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

These vectors satisfy  $ix_{21} = x_{22}$  and  $x_{21} = -ix_{22}$ . Again, both of these equations say the same thing. Letting  $x_{21} = s$  we get

$$x_2 = \begin{bmatrix} s \\ i \, s \end{bmatrix}, \, s \in \mathbb{C} \, .$$

(Note: Statistical libraries will report a normalized version of  $x_1$  and  $x_2$ . For instance, they might set  $s = \frac{1}{\sqrt{2}}$ .) The eigendecomposition is

$$\underbrace{\begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}}_{A} = \underbrace{\begin{bmatrix} s & s \\ -is & is \end{bmatrix}}_{S} \underbrace{\begin{bmatrix} \frac{1}{\sqrt{2}} + \frac{i}{\sqrt{2}} & 0 \\ 0 & \frac{1}{\sqrt{2}} - \frac{i}{\sqrt{2}} \end{bmatrix}}_{\Lambda} \underbrace{\begin{bmatrix} is & -s \\ is & s \end{bmatrix} \frac{1}{2is^2}}_{S^{-1}} .$$

Ex. 10.9: The eigenvalues of  $A = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}$  solve the equation det  $\lambda I_n - A = \lambda^2 = 0$  so both eigenvalues are zero (alternatively, just note that A is upper triangular with zeros along its diagonal). The eigenvectors are therefore simply the null space of A, i.e., vectors  $x = \begin{bmatrix} x_{11} & x_{21} \end{bmatrix}^T$  such that Ax = 0, which are all vectors of the form  $x = \begin{bmatrix} s & 0 \end{bmatrix}^T$ . This is a 1-dimensional subspace. The geometric multiplicity of the eigenvalue 0 is 1 which is less than its algebraic multiplicity. Therefore A is non-diagonalizable.

#### Mathematics and Programming for the Quantitative Economist

Ex. 10.10: (a) A and  $A^{\mathrm{T}}$  share the same characteristic polynomial, since  $\det(\lambda I_n - A) = \det((\lambda I_n - A)^{\mathrm{T}}) = \det(\lambda I_n - A^{\mathrm{T}})$ . Therefore A and  $A^{\mathrm{T}}$  share the same characteris. They do *not* necessarily share the same eigenvectors!

(b) If  $Ax = \lambda x$  and A is non-singular, then  $A^{-1}Ax = A^{-1}(\lambda x)$ , so  $x = \lambda A^{-1}x$ . This implies  $A^{-1}x = (1/\lambda)x$ , which shows that the eigenvalues of  $A^{-1}$  are the reciprocals of the eigenvalues of A.

(c) If A is idempotent (AA = A), then  $Ax = \lambda x$  implies  $AAx = A(\lambda x) = \lambda Ax = \lambda^2 x$ , therefore we have  $\lambda^2 = \lambda$  which implies  $\lambda$  is zero or one.

(d) If A is idempotent, it only has unit and zero eigenvalues. Suppose there are s unit eigenvectors, and n-s zero eigenvectors. Suppose A has rank r. Consider the eigenvalue  $\lambda = 1$ . The corresponding eigenspace is the space of all vectors x such that  $(I_n - A)x = 0_n$ . This is exactly the column space of A, C(A, r): If  $x \in C(A, r)$ , i.e., x = Ay for some  $y \in \mathbb{R}^n$ , then  $(I_n - A)Ay = Ay - A^2y = Ay - Ay = 0_n$ . If  $(I_n - A)x = 0_n$ , then Ax = x which says that  $x \in C(A, r)$ . Since geometric multiplicity of eigenvalues cannot exceed algebraic multiplicity, we have  $r \leq s$ . Consider the eigenvalue  $\lambda = 0$ . The corresponding eigenspace is the space of all vectors x such that Ax = 0. This is the null space of A, N(A, n - r), which has dimension n - r. Again, we have  $n - r \leq n - s$ , which implies  $r \geq s$ .

Therefore r = s, and the geometric multiplicities of eigenvalues 1 and 0 are equal to their respective algebraic multiplicities, so the matrix is diagonalizable. (e) The sum of eigenvalues of a matrix is equal to its trace, which for idempotent

matrices is equal to the number of unit eigenvalues, which from part (d) is equal to its rank.