

## 4.8: Expected Value and Covariance Matrices

The main purpose of this section is a discussion of expected value and covariance for random matrices and vectors. These topics are somewhat specialized, but are particularly important in multivariate statistical models and for the multivariate normal distribution. This section requires some prerequisite knowledge of linear algebra.

We assume that the various indices  $m, n, p, k$  that occur in this section are positive integers. Also we assume that expected values of real-valued random variables that we reference exist as real numbers, although extensions to cases where expected values are  $\infty$  or  $-\infty$  are straightforward, as long as we avoid the dreaded indeterminate form  $\infty - \infty$ .

### Basic Theory

#### Linear Algebra

We will follow our usual convention of denoting random variables by upper case letters and nonrandom variables and constants by lower case letters. In this section, that convention leads to notation that is a bit nonstandard, since the objects that we will be dealing with are vectors and matrices. On the other hand, the notation we will use works well for illustrating the similarities between results for random matrices and the corresponding results in the one-dimensional case. Also, we will try to be careful to explicitly point out the underlying spaces where various objects live.

Let  $\mathbb{R}^{m \times n}$  denote the space of all  $m \times n$  matrices of real numbers. The  $(i, j)$  entry of  $\mathbf{a} \in \mathbb{R}^{m \times n}$  is denoted  $a_{ij}$  for  $i \in \{1, 2, \dots, m\}$  and  $j \in \{1, 2, \dots, n\}$ . We will identify  $\mathbb{R}^n$  with  $\mathbb{R}^{n \times 1}$ , so that an ordered  $n$ -tuple can also be thought of as an  $n \times 1$  column vector. The transpose of a matrix  $\mathbf{a} \in \mathbb{R}^{m \times n}$  is denoted  $\mathbf{a}^T$ —the  $n \times m$  matrix whose  $(i, j)$  entry is the  $(j, i)$  entry of  $\mathbf{a}$ . Recall the definitions of matrix addition, scalar multiplication, and matrix multiplication. Recall also the standard inner product (or dot product) of  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ :

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x} \cdot \mathbf{y} = \mathbf{x}^T \mathbf{y} = \sum_{i=1}^n x_i y_i \quad (4.8.1)$$

The outer product of  $\mathbf{x}$  and  $\mathbf{y}$  is  $\mathbf{x}\mathbf{y}^T$ , the  $n \times n$  matrix whose  $(i, j)$  entry is  $x_i y_j$ . Note that the inner product is the trace (sum of the diagonal entries) of the outer product. Finally recall the standard norm on  $\mathbb{R}^n$ , given by

$$\|\mathbf{x}\| = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle} = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2} \quad (4.8.2)$$

Recall that inner product is *bilinear*, that is, linear (preserving addition and scalar multiplication) in each argument separately. As a consequence, for  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ ,

$$\|\mathbf{x} + \mathbf{y}\|^2 = \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 + 2\langle \mathbf{x}, \mathbf{y} \rangle \quad (4.8.3)$$

#### Expected Value of a Random Matrix

As usual, our starting point is a random experiment modeled by a probability space  $(\Omega, \mathcal{F}, \mathbb{P})$ . So to review,  $\Omega$  is the set of outcomes,  $\mathcal{F}$  the collection of events, and  $\mathbb{P}$  the probability measure on the sample space  $(\Omega, \mathcal{F})$ . It's natural to define the expected value of a random matrix in a component-wise manner.

Suppose that  $\mathbf{X}$  is an  $m \times n$  matrix of real-valued random variables, whose  $(i, j)$  entry is denoted  $X_{ij}$ . Equivalently,  $\mathbf{X}$  is as a random  $m \times n$  matrix, that is, a random variable with values in  $\mathbb{R}^{m \times n}$ . The *expected value*  $\mathbb{E}(\mathbf{X})$  is defined to be the  $m \times n$  matrix whose  $(i, j)$  entry is  $\mathbb{E}(X_{ij})$ , the expected value of  $X_{ij}$ .

Many of the basic properties of expected value of random variables have analogous results for expected value of random matrices, with matrix operation replacing the ordinary ones. Our first two properties are the critically important *linearity properties*. The first part is the *additive property*—the expected value of a sum is the sum of the expected values.

$\mathbb{E}(\mathbf{X} + \mathbf{Y}) = \mathbb{E}(\mathbf{X}) + \mathbb{E}(\mathbf{Y})$  if  $\mathbf{X}$  and  $\mathbf{Y}$  are random  $m \times n$  matrices.

**Proof**

This is true by definition of the matrix expected value and the ordinary additive property. Note that  $\mathbb{E}(X_{ij} + Y_{ij}) = \mathbb{E}(X_{ij}) + \mathbb{E}(Y_{ij})$ . The left side is the  $(i, j)$  entry of  $\mathbb{E}(\mathbf{X} + \mathbf{Y})$  and the right side is the  $(i, j)$  entry of  $\mathbb{E}(\mathbf{X}) + \mathbb{E}(\mathbf{Y})$ .

The next part of the linearity properties is the *scaling property*—a nonrandom matrix factor can be pulled out of the expected value.

Suppose that  $\mathbf{X}$  is a random  $n \times p$  matrix.

1.  $\mathbb{E}(\mathbf{a}\mathbf{X}) = \mathbf{a}\mathbb{E}(\mathbf{X})$  if  $\mathbf{a} \in \mathbb{R}^{m \times n}$ .
2.  $\mathbb{E}(\mathbf{X}\mathbf{a}) = \mathbb{E}(\mathbf{X})\mathbf{a}$  if  $\mathbf{a} \in \mathbb{R}^{p \times n}$ .

**Proof**

1. By the ordinary linearity and scaling properties,  $\mathbb{E}\left(\sum_{j=1}^n a_{ij} X_{jk}\right) = \sum_{j=1}^n a_{ij} \mathbb{E}(X_{jk})$ . The left side is the  $(i, k)$  entry of  $\mathbb{E}(\mathbf{a}\mathbf{X})$  and the right side is the  $(i, k)$  entry of  $\mathbf{a}\mathbb{E}(\mathbf{X})$ .
2. The proof is similar to (a).

Recall that for independent, real-valued variables, the expected value of the product is the product of the expected values. Here is the analogous result for random matrices.

$\mathbb{E}(\mathbf{X}\mathbf{Y}) = \mathbb{E}(\mathbf{X})\mathbb{E}(\mathbf{Y})$  if  $\mathbf{X}$  is a random  $m \times n$  matrix,  $\mathbf{Y}$  is a random  $n \times p$  matrix, and  $\mathbf{X}$  and  $\mathbf{Y}$  are independent.

**Proof**

By the ordinary linearity properties and by the independence assumption,

$$\mathbb{E} \left( \sum_{j=1}^n X_{ij} Y_{jk} \right) = \sum_{j=1}^n \mathbb{E} (X_{ij} Y_{jk}) = \sum_{j=1}^n \mathbb{E} (X_{ij}) \mathbb{E} (Y_{jk}) \quad (4.8.4)$$

The left side is the  $(i, k)$  entry of  $\mathbb{E}(\mathbf{XY})$  and the right side is the  $(i, k)$  entry of  $\mathbb{E}(\mathbf{X})\mathbb{E}(\mathbf{Y})$ .

Actually the previous result holds if  $\mathbf{X}$  and  $\mathbf{Y}$  are simply uncorrelated in the sense that  $X_{ij}$  and  $Y_{jk}$  are uncorrelated for each  $i \in \{1, \dots, m\}$ ,  $j \in \{1, 2, \dots, n\}$  and  $k \in \{1, 2, \dots, p\}$ . We will study covariance of random vectors in the next subsection.

### Covariance Matrices

Our next goal is to define and study the covariance of two random vectors.

Suppose that  $\mathbf{X}$  is a random vector in  $\mathbb{R}^m$  and  $\mathbf{Y}$  is a random vector in  $\mathbb{R}^n$ .

1. The *covariance matrix* of  $\mathbf{X}$  and  $\mathbf{Y}$  is the  $m \times n$  matrix  $\text{cov}(\mathbf{X}, \mathbf{Y})$  whose  $(i, j)$  entry is  $\text{cov}(X_i, Y_j)$  the ordinary covariance of  $X_i$  and  $Y_j$ .
2. Assuming that the coordinates of  $\mathbf{X}$  and  $\mathbf{Y}$  have positive variance, the *correlation matrix* of  $\mathbf{X}$  and  $\mathbf{Y}$  is the  $m \times n$  matrix  $\text{cor}(\mathbf{X}, \mathbf{Y})$  whose  $(i, j)$  entry is  $\text{cor}(X_i, Y_j)$ , the ordinary correlation of  $X_i$  and  $Y_j$

Many of the standard properties of covariance and correlation for real-valued random variables have extensions to random vectors. For the following three results,  $\mathbf{X}$  is a random vector in  $\mathbb{R}^m$  and  $\mathbf{Y}$  is a random vector in  $\mathbb{R}^n$ .

$$\text{cov}(\mathbf{X}, \mathbf{Y}) = \mathbb{E} \left( [\mathbf{X} - \mathbb{E}(\mathbf{X})] [\mathbf{Y} - \mathbb{E}(\mathbf{Y})]^T \right)$$

Proof

By the definition of the expected value of a random vector and by the definition of matrix multiplication, the  $(i, j)$  entry of  $[\mathbf{X} - \mathbb{E}(\mathbf{X})] [\mathbf{Y} - \mathbb{E}(\mathbf{Y})]^T$  is simply  $[X_i - \mathbb{E}(X_i)] [Y_j - \mathbb{E}(Y_j)]$ . The expected value of this entry is  $\text{cov}(X_i, Y_j)$ , which in turn, is the  $(i, j)$  entry of  $\text{cov}(\mathbf{X}, \mathbf{Y})$ .

Thus, the covariance of  $\mathbf{X}$  and  $\mathbf{Y}$  is the expected value of the outer product of  $\mathbf{X} - \mathbb{E}(\mathbf{X})$  and  $\mathbf{Y} - \mathbb{E}(\mathbf{Y})$ . Our next result is the computational formula for covariance: the expected value of the outer product of  $\mathbf{X}$  and  $\mathbf{Y}$  minus the outer product of the expected values.

$$\text{cov}(\mathbf{X}, \mathbf{Y}) = \mathbb{E}(\mathbf{XY}^T) - \mathbb{E}(\mathbf{X})[\mathbb{E}(\mathbf{Y})]^T.$$

Proof

The  $(i, j)$  entry of  $\mathbb{E}(\mathbf{XY}^T) - \mathbb{E}(\mathbf{X})[\mathbb{E}(\mathbf{Y})]^T$  is  $\mathbb{E}(X_i Y_j) - \mathbb{E}(X_i) \mathbb{E}(Y_j)$ , which by the standard computational formula, is  $\text{cov}(X_i, Y_j)$ , which in turn is the  $(i, j)$  entry of  $\text{cov}(\mathbf{X}, \mathbf{Y})$ .

The next result is the matrix version of the symmetry property.

$$\text{cov}(\mathbf{Y}, \mathbf{X}) = [\text{cov}(\mathbf{X}, \mathbf{Y})]^T.$$

Proof

The  $(i, j)$  entry of  $\text{cov}(\mathbf{X}, \mathbf{Y})$  is  $\text{cov}(X_i, Y_j)$ , which is the  $(j, i)$  entry of  $\text{cov}(\mathbf{Y}, \mathbf{X})$ .

In the following result,  $\mathbf{0}$  denotes the  $m \times n$  zero matrix.

$\text{cov}(\mathbf{X}, \mathbf{Y}) = \mathbf{0}$  if and only if  $\text{cov}(X_i, Y_j) = 0$  for each  $i$  and  $j$ , so that each coordinate of  $\mathbf{X}$  is uncorrelated with each coordinate of  $\mathbf{Y}$ .

Proof

This follows immediately from the definition of  $\text{cov}(\mathbf{X}, \mathbf{Y})$ .

Naturally, when  $\text{cov}(\mathbf{X}, \mathbf{Y}) = \mathbf{0}$ , we say that the random vectors  $\mathbf{X}$  and  $\mathbf{Y}$  are *uncorrelated*. In particular, if the random vectors are independent, then they are uncorrelated. The following results establish the *bi-linear* properties of covariance.

The additive properties.

1.  $\text{cov}(\mathbf{X} + \mathbf{Y}, \mathbf{Z}) = \text{cov}(\mathbf{X}, \mathbf{Z}) + \text{cov}(\mathbf{Y}, \mathbf{Z})$  if  $\mathbf{X}$  and  $\mathbf{Y}$  are random vectors in  $\mathbb{R}^m$  and  $\mathbf{Z}$  is a random vector in  $\mathbb{R}^n$ .
2.  $\text{cov}(\mathbf{X}, \mathbf{Y} + \mathbf{Z}) = \text{cov}(\mathbf{X}, \mathbf{Y}) + \text{cov}(\mathbf{X}, \mathbf{Z})$  if  $\mathbf{X}$  is a random vector in  $\mathbb{R}^m$ , and  $\mathbf{Y}$  and  $\mathbf{Z}$  are random vectors in  $\mathbb{R}^n$ .

Proof

1. From the ordinary additive property of covariance,  $\text{cov}(X_i + Y_i, Z_j) = \text{cov}(X_i, Z_j) + \text{cov}(Y_i, Z_j)$ . The left side is the  $(i, j)$  entry of  $\text{cov}(\mathbf{X} + \mathbf{Y}, \mathbf{Z})$  and the right side is the  $(i, j)$  entry of  $\text{cov}(\mathbf{X}, \mathbf{Z}) + \text{cov}(\mathbf{Y}, \mathbf{Z})$ .
2. The proof is similar to (a), using the additivity of covariance in the second argument.

The scaling properties

1.  $\text{cov}(\mathbf{aX}, \mathbf{Y}) = \mathbf{a} \text{cov}(\mathbf{X}, \mathbf{Y})$  if  $\mathbf{X}$  is a random vector in  $\mathbb{R}^n$ ,  $\mathbf{Y}$  is a random vector in  $\mathbb{R}^p$ , and  $\mathbf{a} \in \mathbb{R}^{m \times n}$ .
2.  $\text{cov}(\mathbf{X}, \mathbf{aY}) = \text{cov}(\mathbf{X}, \mathbf{Y}) \mathbf{a}^T$  if  $\mathbf{X}$  is a random vector in  $\mathbb{R}^m$ ,  $\mathbf{Y}$  is a random vector in  $\mathbb{R}^n$ , and  $\mathbf{a} \in \mathbb{R}^{k \times n}$ .

Proof

1. Using the ordinary linearity properties of covariance in the first argument, we have

$$\text{cov}\left(\sum_{j=1}^n a_{ij}X_j, Y_k\right) = \sum_{j=1}^n a_{ij}\text{cov}(X_j, Y_k) \quad (4.8.5)$$

The left side is the  $(i, k)$  entry of  $\text{cov}(\mathbf{aX}, \mathbf{Y})$  and the right side is the  $(i, k)$  entry of  $\mathbf{a}\text{cov}(\mathbf{X}, \mathbf{Y})$ .  
2. The proof is similar to (a), using the linearity of covariance in the second argument.

### Variance-Covariance Matrices

Suppose that  $\mathbf{X}$  is a random vector in  $\mathbb{R}^n$ . The covariance matrix of  $\mathbf{X}$  with itself is called the *variance-covariance matrix* of  $\mathbf{X}$ :

$$\text{vc}(\mathbf{X}) = \text{cov}(\mathbf{X}, \mathbf{X}) = \mathbb{E}\left([\mathbf{X} - \mathbb{E}(\mathbf{X})][\mathbf{X} - \mathbb{E}(\mathbf{X})]^T\right) \quad (4.8.6)$$

Recall that for an ordinary real-valued random variable  $X$ ,  $\text{var}(X) = \text{cov}(X, X)$ . Thus the variance-covariance matrix of a random vector in some sense plays the same role that variance does for a random variable.

$\text{vc}(\mathbf{X})$  is a symmetric  $n \times n$  matrix with  $(\text{var}(X_1), \text{var}(X_2), \dots, \text{var}(X_n))$  on the diagonal.

Proof

Recall that  $\text{cov}(X_i, X_j) = \text{cov}(X_j, X_i)$ . Also, the  $(i, i)$  entry of  $\text{vc}(\mathbf{X})$  is  $\text{cov}(X_i, X_i) = \text{var}(X_i)$ .

The following result is the formula for the variance-covariance matrix of a sum, analogous to the formula for the variance of a sum of real-valued variables.

$\text{vc}(\mathbf{X} + \mathbf{Y}) = \text{vc}(\mathbf{X}) + \text{cov}(\mathbf{X}, \mathbf{Y}) + \text{cov}(\mathbf{Y}, \mathbf{X}) + \text{vc}(\mathbf{Y})$  if  $\mathbf{X}$  and  $\mathbf{Y}$  are random vectors in  $\mathbb{R}^n$ .

Proof

This follows from the additive property of covariance:

$$\text{vc}(\mathbf{X} + \mathbf{Y}) = \text{cov}(\mathbf{X} + \mathbf{Y}, \mathbf{X} + \mathbf{Y}) = \text{cov}(\mathbf{X}, \mathbf{X}) + \text{cov}(\mathbf{X}, \mathbf{Y}) + \text{cov}(\mathbf{Y}, \mathbf{X}) + \text{cov}(\mathbf{Y}, \mathbf{Y}) \quad (4.8.7)$$

Recall that  $\text{var}(aX) = a^2\text{var}(X)$  if  $X$  is a real-valued random variable and  $a \in \mathbb{R}$ . Here is the analogous result for the variance-covariance matrix of a random vector.

$\text{vc}(\mathbf{aX}) = \mathbf{a}\text{vc}(\mathbf{X})\mathbf{a}^T$  if  $\mathbf{X}$  is a random vector in  $\mathbb{R}^n$  and  $\mathbf{a} \in \mathbb{R}^{m \times n}$ .

Proof

This follows from the scaling property of covariance:

$$\text{vc}(\mathbf{aX}) = \text{cov}(\mathbf{aX}, \mathbf{aX}) = \mathbf{a}\text{cov}(\mathbf{X}, \mathbf{X})\mathbf{a}^T \quad (4.8.8)$$

Recall that if  $X$  is a random variable, then  $\text{var}(X) \geq 0$ , and  $\text{var}(X) = 0$  if and only if  $X$  is a constant (with probability 1). Here is the analogous result for a random vector:

Suppose that  $\mathbf{X}$  is a random vector in  $\mathbb{R}^n$ .

1.  $\text{vc}(\mathbf{X})$  is either positive semi-definite or positive definite.
2.  $\text{vc}(\mathbf{X})$  is positive semi-definite but not positive definite if and only if there exists  $\mathbf{a} \in \mathbb{R}^n$  and  $c \in \mathbb{R}$  such that, with probability 1,  $\mathbf{a}^T \mathbf{X} = \sum_{i=1}^n a_i X_i = c$

Proof

1. From the previous result,  $0 \leq \text{var}(\mathbf{a}^T \mathbf{X}) = \text{vc}(\mathbf{a}^T \mathbf{X}) = \mathbf{a}^T \text{vc}(\mathbf{X}) \mathbf{a}$  for every  $\mathbf{a} \in \mathbb{R}^n$ . Thus, by definition,  $\text{vc}(\mathbf{X})$  is either positive semi-definite or positive definite.
2. In light of (a),  $\text{vc}(\mathbf{X})$  is positive semi-definite but not positive definite if and only if there exists  $\mathbf{a} \in \mathbb{R}^n$  such that  $\mathbf{a}^T \text{vc}(\mathbf{X}) \mathbf{a} = \text{var}(\mathbf{a}^T \mathbf{X}) = 0$ . But in turn, this is true if and only if  $\mathbf{a}^T \mathbf{X}$  is constant with probability 1.

Recall that since  $\text{vc}(\mathbf{X})$  is either positive semi-definite or positive definite, the eigenvalues and the determinant of  $\text{vc}(\mathbf{X})$  are nonnegative. Moreover, if  $\text{vc}(\mathbf{X})$  is positive semi-definite but not positive definite, then one of the coordinates of  $\mathbf{X}$  can be written as a linear transformation of the other coordinates (and hence can usually be eliminated in the underlying model). By contrast, if  $\text{vc}(\mathbf{X})$  is positive definite, then this cannot happen;  $\text{vc}(\mathbf{X})$  has positive eigenvalues and determinant and is invertible.

### Best Linear Predictor

Suppose that  $\mathbf{X}$  is a random vector in  $\mathbb{R}^m$  and that  $\mathbf{Y}$  is a random vector in  $\mathbb{R}^n$ . We are interested in finding the function of  $\mathbf{X}$  of the form  $\mathbf{a} + \mathbf{bX}$ , where  $\mathbf{a} \in \mathbb{R}^n$  and  $\mathbf{b} \in \mathbb{R}^{n \times m}$ , that is closest to  $\mathbf{Y}$  in the mean square sense. Functions of this form are analogous to linear functions in the single variable case. However, unless  $\mathbf{a} = \mathbf{0}$ , such functions are not *linear transformations* in the sense of linear algebra, so the correct term is *affine function* of  $\mathbf{X}$ . This problem is of fundamental importance in statistics when random vector  $\mathbf{X}$ , the *predictor vector* is observable, but not random vector  $\mathbf{Y}$ , the *response vector*. Our discussion here generalizes the one-dimensional case, when  $X$  and  $Y$  are random variables. That problem was solved in the section on Covariance and Correlation. We will assume that  $\text{vc}(\mathbf{X})$  is positive definite, so that  $\text{vc}(\mathbf{X})$  is invertible, and none of the coordinates of  $\mathbf{X}$  can be written as an affine function of the other coordinates. We write  $\text{vc}^{-1}(\mathbf{X})$  for the inverse instead of the clunkier  $[\text{vc}(\mathbf{X})]^{-1}$ .

As with the single variable case, the solution turns out to be the affine function that has the same expected value as  $\mathbf{Y}$ , and whose covariance with  $\mathbf{X}$  is the same as that of  $\mathbf{Y}$ .

Define  $L(\mathbf{Y} | \mathbf{X}) = \mathbb{E}(\mathbf{Y}) + \text{cov}(\mathbf{Y}, \mathbf{X})\text{vc}^{-1}(\mathbf{X})[\mathbf{X} - \mathbb{E}(\mathbf{X})]$ . Then  $L(\mathbf{Y} | \mathbf{X})$  is the only affine function of  $\mathbf{X}$  in  $\mathbb{R}^n$  satisfying

1.  $\mathbb{E}[L(\mathbf{Y} | \mathbf{X})] = \mathbb{E}(\mathbf{Y})$
2.  $\text{cov}[L(\mathbf{Y} | \mathbf{X}), \mathbf{X}] = \text{cov}(\mathbf{Y}, \mathbf{X})$

Proof

From linearity,

$$\mathbb{E}[L(\mathbf{Y} | \mathbf{X})] = E(\mathbf{Y}) + \text{cov}(\mathbf{Y}, \mathbf{X})\text{vc}^{-1}(\mathbf{X})[\mathbb{E}(\mathbf{X}) - \mathbb{E}(\mathbf{X})] = 0 \quad (4.8.9)$$

From linearity and the fact that a constant vector is independent (and hence uncorrelated) with any random vector,

$$\text{cov}[L(\mathbf{Y} | \mathbf{X}), \mathbf{X}] = \text{cov}(\mathbf{Y}, \mathbf{X})\text{vc}^{-1}(\mathbf{X})\text{cov}(\mathbf{X}, \mathbf{X}) = \text{cov}(\mathbf{Y}, \mathbf{X})\text{vc}^{-1}(\mathbf{X})\text{vc}(\mathbf{X}) = \text{cov}(\mathbf{Y}, \mathbf{X}) \quad (4.8.10)$$

Conversely, suppose that  $\mathbf{U} = \mathbf{a} + \mathbf{b}\mathbf{X}$  for some  $\mathbf{a} \in \mathbb{R}^n$  and  $\mathbf{b} \in \mathbb{R}^{m \times n}$ , and that  $\mathbb{E}(\mathbf{U}) = \mathbb{E}(\mathbf{Y})$  and  $\text{cov}(\mathbf{U}, \mathbf{X}) = \text{cov}(\mathbf{Y}, \mathbf{X})$ . From the second equation, again using linearity and the uncorrelated property of constant vectors, we get  $\mathbf{b}\text{cov}(\mathbf{X}, \mathbf{X}) = \text{cov}(\mathbf{Y}, \mathbf{X})$  and therefore  $\mathbf{b} = \text{cov}(\mathbf{Y}, \mathbf{X})\text{vc}^{-1}(\mathbf{X})$ . Then from the first equation,  $\mathbf{a} + \mathbf{b}\mathbb{E}(\mathbf{X}) = \mathbf{Y}$  so  $\mathbf{a} = \mathbb{E}(\mathbf{Y}) - \mathbf{b}\mathbb{E}(\mathbf{X})$ .

A simple corollary is the  $\mathbf{Y} - L(\mathbf{Y} | \mathbf{X})$  is uncorrelated with any affine function of  $\mathbf{X}$ :

If  $\mathbf{U}$  is an affine function of  $\mathbf{X}$  then

1.  $\text{cov}[\mathbf{Y} - L(\mathbf{Y} | \mathbf{X}), \mathbf{U}] = \mathbf{0}$
2.  $\mathbb{E}(\langle \mathbf{Y} - L(\mathbf{Y} | \mathbf{X}), \mathbf{U} \rangle) = 0$

Proof

Suppose that  $\mathbf{U} = \mathbf{a} + \mathbf{b}\mathbf{X}$  where  $\mathbf{a} \in \mathbb{R}^n$  and  $\mathbf{b} \in \mathbb{R}^{m \times n}$ . For simplicity, let  $\mathbf{L} = L(\mathbf{Y} | \mathbf{X})$

1. From the previous result,  $\text{cov}(\mathbf{Y}, \mathbf{X}) = \text{cov}(\mathbf{L}, \mathbf{X})$ . Hence using linearity,

$$\text{cov}(\mathbf{Y} - \mathbf{L}, \mathbf{U}) = \text{cov}(\mathbf{Y} - \mathbf{L}, \mathbf{a}) + \text{cov}(\mathbf{Y} - \mathbf{L}, \mathbf{X})\mathbf{b}^T = \mathbf{0} + [\text{cov}(\mathbf{Y}, \mathbf{X}) - \text{cov}(\mathbf{L}, \mathbf{X})] = \mathbf{0} \quad (4.8.11)$$

2. Recall that  $\langle \mathbf{Y} - \mathbf{L}, \mathbf{U} \rangle$  is the trace of  $\text{cov}(\mathbf{Y} - \mathbf{L}, \mathbf{U})$  and hence has expected value 0 by part (a).

The variance-covariance matrix of  $L(\mathbf{Y} | \mathbf{X})$ , and its covariance matrix with  $\mathbf{Y}$  turn out to be the same, again analogous to the single variable case.

Additional properties of  $L(\mathbf{Y} | \mathbf{X})$ :

1.  $\text{cov}[\mathbf{Y}, L(\mathbf{Y} | \mathbf{X})] = \text{cov}(\mathbf{Y}, \mathbf{X})\text{vc}^{-1}(\mathbf{X})\text{cov}(\mathbf{X}, \mathbf{Y})$
2.  $\text{vc}[L(\mathbf{Y} | \mathbf{X})] = \text{cov}(\mathbf{Y}, \mathbf{X})\text{vc}^{-1}(\mathbf{X})\text{cov}(\mathbf{X}, \mathbf{Y})$

Proof

Recall that  $L(\mathbf{Y} | \mathbf{X}) = \mathbb{E}(\mathbf{Y}) + \text{cov}(\mathbf{Y}, \mathbf{X})\text{vc}^{-1}(\mathbf{X})[\mathbf{X} - \mathbb{E}(\mathbf{X})]$

1. Using basic properties of covariance,

$$\text{cov}[\mathbf{Y}, L(\mathbf{Y} | \mathbf{X})] = \text{cov}[\mathbf{Y}, \mathbf{X} - \mathbb{E}(\mathbf{X})] [\text{cov}(\mathbf{Y}, \mathbf{X})\text{vc}^{-1}(\mathbf{X})]^T = \text{cov}(\mathbf{Y}, \mathbf{X})\text{vc}^{-1}(\mathbf{X})\text{cov}(\mathbf{X}, \mathbf{Y}) \quad (4.8.12)$$

2. Using basic properties of variance-covariance,

$$\text{vc}[L(\mathbf{Y} | \mathbf{X})] = \text{vc}[\text{cov}(\mathbf{Y}, \mathbf{X})\text{vc}^{-1}(\mathbf{X})\mathbf{X}] = \text{cov}(\mathbf{Y}, \mathbf{X})\text{vc}^{-1}(\mathbf{X})\text{vc}(\mathbf{X})[\text{cov}(\mathbf{Y}, \mathbf{X})\text{vc}^{-1}(\mathbf{X})]^T = \text{cov}(\mathbf{Y}, \mathbf{X})\text{vc}^{-1}(\mathbf{X})\text{cov}(\mathbf{X}, \mathbf{Y}) \quad (4.8.13)$$

Next is the fundamental result that  $L(\mathbf{Y} | \mathbf{X})$  is the affine function of  $\mathbf{X}$  that is closest to  $\mathbf{Y}$  in the mean square sense.

Suppose that  $\mathbf{U} \in \mathbb{R}^n$  is an affine function of  $\mathbf{X}$ . Then

1.  $\mathbb{E}(\|\mathbf{Y} - L(\mathbf{Y} | \mathbf{X})\|^2) \leq \mathbb{E}(\|\mathbf{Y} - \mathbf{U}\|^2)$
2. Equality holds in (a) if and only if  $\mathbf{U} = L(\mathbf{Y} | \mathbf{X})$  with probability 1.

Proof

Again, let  $\mathbf{L} = L(\mathbf{Y} | \mathbf{X})$  for simplicity and let  $\mathbf{U} \in \mathbb{R}^n$  be an affine function of  $\mathbf{X}$ .

1. Using the linearity of expected value, note that

$$\mathbb{E}(\|\mathbf{Y} - \mathbf{U}\|^2) = \mathbb{E}(\|(\mathbf{Y} - \mathbf{L}) + (\mathbf{L} - \mathbf{U})\|^2) = \mathbb{E}(\|\mathbf{Y} - \mathbf{L}\|^2) + 2\mathbb{E}(\langle \mathbf{Y} - \mathbf{L}, \mathbf{L} - \mathbf{U} \rangle) + \mathbb{E}(\|\mathbf{L} - \mathbf{U}\|^2) \quad (4.8.14)$$

But  $\mathbf{L} - \mathbf{U}$  is an affine function of  $\mathbf{X}$  and hence the middle term is 0 by our previous corollary. Hence

$$\mathbb{E}(\|\mathbf{Y} - \mathbf{U}\|^2) = \mathbb{E}(\|\mathbf{L} - \mathbf{U}\|^2) + \mathbb{E}(\|\mathbf{Y} - \mathbf{L}\|^2) \geq \mathbb{E}(\|\mathbf{L} - \mathbf{U}\|^2)$$

2. From (a), equality holds in the inequality if and only if  $\mathbb{E}(\|\mathbf{L} - \mathbf{U}\|^2) = 0$  if and only if  $\mathbb{P}(\mathbf{L} = \mathbf{U}) = 1$ .

The variance-covariance matrix of the difference between  $\mathbf{Y}$  and the best affine approximation is given in the next theorem.

$$\text{vc}[\mathbf{Y} - L(\mathbf{Y} | \mathbf{X})] = \text{vc}(\mathbf{Y}) - \text{cov}(\mathbf{Y}, \mathbf{X})\text{vc}^{-1}(\mathbf{X})\text{cov}(\mathbf{X}, \mathbf{Y})$$

Proof

Again, we abbreviate  $L(\mathbf{Y} | \mathbf{X})$  by  $\mathbf{L}$ . Using basic properties of variance-covariance matrices,

$$\text{vc}(\mathbf{Y} - \mathbf{L}) = \text{vc}(\mathbf{Y}) - \text{cov}(\mathbf{Y}, \mathbf{L}) - \text{cov}(\mathbf{L}, \mathbf{Y}) + \text{vc}(\mathbf{L}) \quad (4.8.15)$$

But  $\text{cov}(\mathbf{Y}, \mathbf{L}) = \text{cov}(\mathbf{L}, \mathbf{Y}) = \text{vc}(\mathbf{L}) = \text{cov}(\mathbf{Y}, \mathbf{X})\text{vc}^{-1}(\mathbf{X})\text{cov}(\mathbf{Y}, \mathbf{X})$ . Substituting gives the result.

The actual mean square error when we use  $L(\mathbf{Y} | \mathbf{X})$  to approximate  $\mathbf{Y}$ , namely  $\mathbb{E}(\|\mathbf{Y} - L(\mathbf{Y} | \mathbf{X})\|^2)$ , is the trace (sum of the diagonal entries) of the variance-covariance matrix above. The function of  $\mathbf{x}$  given by

$$L(\mathbf{Y} | \mathbf{X} = \mathbf{x}) = \mathbb{E}(\mathbf{Y}) + \text{cov}(\mathbf{Y}, \mathbf{X}) \text{vc}^{-1}(\mathbf{X}) [\mathbf{x} - \mathbb{E}(\mathbf{X})] \quad (4.8.16)$$

is known as the (distribution) linear regression function. If we observe  $\mathbf{x}$  then  $L(\mathbf{Y} | \mathbf{X} = \mathbf{x})$  is our best affine prediction of  $\mathbf{Y}$ .

Multiple linear regression is more powerful than it may at first appear, because it can be applied to non-linear transformations of the random vectors. That is, if  $g: \mathbb{R}^m \rightarrow \mathbb{R}^j$  and  $h: \mathbb{R}^n \rightarrow \mathbb{R}^k$  then  $L[h(\mathbf{Y}) | g(\mathbf{X})]$  is the affine function of  $g(\mathbf{X})$  that is closest to  $h(\mathbf{Y})$  in the mean square sense. Of course, we must be able to compute the appropriate means, variances, and covariances.

Moreover, *Non-linear regression* with a single, real-valued predictor variable can be thought of as a special case of multiple linear regression. Thus, suppose that  $X$  is the predictor variable,  $Y$  is the response variable, and that  $(g_1, g_2, \dots, g_n)$  is a sequence of real-valued functions. We can apply the results of this section to find the linear function of  $(g_1(X), g_2(X), \dots, g_n(X))$  that is closest to  $Y$  in the mean square sense. We just replace  $X_i$  with  $g_i(X)$  for each  $i$ . Again, we must be able to compute the appropriate means, variances, and covariances to do this.

## Examples and Applications

Suppose that  $(X, Y)$  has probability density function  $f$  defined by  $f(x, y) = x + y$  for  $0 \leq x \leq 1, 0 \leq y \leq 1$ . Find each of the following:

1.  $\mathbb{E}(X, Y)$
2.  $\text{vc}(X, Y)$

Answer

1.  $\left(\frac{7}{12}, \frac{7}{12}\right)$
2.  $\begin{bmatrix} \frac{11}{144} & -\frac{1}{144} \\ -\frac{1}{144} & \frac{11}{144} \end{bmatrix}$

Suppose that  $(X, Y)$  has probability density function  $f$  defined by  $f(x, y) = 2(x + y)$  for  $0 \leq x \leq y \leq 1$ . Find each of the following:

1.  $\mathbb{E}(X, Y)$
2.  $\text{vc}(X, Y)$

Answer

1.  $\left(\frac{5}{12}, \frac{3}{4}\right)$
2.  $\begin{bmatrix} \frac{43}{720} & \frac{1}{48} \\ \frac{1}{48} & \frac{3}{80} \end{bmatrix}$

Suppose that  $(X, Y)$  has probability density function  $f$  defined by  $f(x, y) = 6x^2y$  for  $0 \leq x \leq 1, 0 \leq y \leq 1$ . Find each of the following:

1.  $\mathbb{E}(X, Y)$
2.  $\text{vc}(X, Y)$

Answer

Note that  $X$  and  $Y$  are independent.

1.  $\left(\frac{3}{4}, \frac{2}{3}\right)$
2.  $\begin{bmatrix} \frac{3}{80} & 0 \\ 0 & \frac{1}{18} \end{bmatrix}$

Suppose that  $(X, Y)$  has probability density function  $f$  defined by  $f(x, y) = 15x^2y$  for  $0 \leq x \leq y \leq 1$ . Find each of the following:

1.  $\mathbb{E}(X, Y)$
2.  $\text{vc}(X, Y)$
3.  $L(Y | X)$
4.  $L[Y | (X, X^2)]$
5. Sketch the regression curves on the same set of axes.

Answer

1.  $\left(\frac{5}{8}, \frac{5}{6}\right)$
2.  $\begin{bmatrix} \frac{17}{448} & \frac{5}{336} \\ \frac{5}{336} & \frac{5}{252} \end{bmatrix}$
3.  $\frac{10}{17} + \frac{20}{51}X$
4.  $\frac{49}{76} + \frac{10}{57}X + \frac{7}{38}X^2$

Suppose that  $(X, Y, Z)$  is uniformly distributed on the region  $\{(x, y, z) \in \mathbb{R}^3 : 0 \leq x \leq y \leq z \leq 1\}$ . Find each of the following:

1.  $\mathbb{E}(X, Y, Z)$
2.  $\text{vc}(X, Y, Z)$

3.  $L[Z | (X, Y)]$
4.  $L[Y | (X, Z)]$
5.  $L[X | (Y, Z)]$
6.  $L[(Y, Z) | X]$

Answer

1.  $(\frac{1}{4}, \frac{1}{2}, \frac{3}{4})$
2.  $\begin{bmatrix} \frac{3}{80} & \frac{1}{40} & \frac{1}{80} \\ \frac{1}{40} & \frac{1}{20} & \frac{1}{40} \\ \frac{1}{80} & \frac{1}{40} & \frac{3}{80} \end{bmatrix}$
3.  $\frac{1}{2} + \frac{1}{2}Y$ . Note that there is no  $X$  term.
4.  $\frac{1}{2}X + \frac{1}{2}Z$ . Note that this is the midpoint of the interval  $[X, Z]$ .
5.  $\frac{1}{2}Y$ . Note that there is no  $Z$  term.
6.  $\begin{bmatrix} \frac{1}{3} + \frac{2}{3}X \\ \frac{2}{3} + \frac{1}{3}X \end{bmatrix}$

Suppose that  $X$  is uniformly distributed on  $(0, 1)$ , and that given  $X$ , random variable  $Y$  is uniformly distributed on  $(0, X)$ . Find each of the following:

1.  $\mathbb{E}(X, Y)$
2.  $\text{vc}(X, Y)$

Answer

1.  $(\frac{1}{2}, \frac{1}{4})$
2.  $\begin{bmatrix} \frac{1}{12} & \frac{1}{24} \\ \frac{1}{24} & \frac{7}{144} \end{bmatrix}$

This page titled [4.8: Expected Value and Covariance Matrices](#) is shared under a [CC BY 2.0](#) license and was authored, remixed, and/or curated by [Kyle Siegrist \(Random Services\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.