

CHAPTER 2. THE CLASSICAL LINEAR REGRESSION MODEL (CLRM)

In Chapter 1, we showed how we estimate an LRM by the method of least squares. As noted in Chapter 1, estimation and hypothesis testing are the twin branches of statistical inference. Based on the OLS, we obtained the sample regression, such as the one shown in Equation (1.40). This is of course a sample regression function (SRF) because it is based on a specific sample drawn randomly from the purported population. What can we say about the true population regression function (PRF) from the SRF? In practice, we do not observe the PRF and have to “guess” it from the SRF. To obtain the best possible guess, we need a framework, which is provided by the **classical linear regression model (CLRM)**. The CLRM is based on several assumptions, which are discussed below.

2.1 Assumptions of the CLRM

We now discuss these assumptions. In Chapters 5 and 6, we will examine these assumptions more critically. However, keep in mind that in any scientific inquiry we start with a set of simplified assumptions and gradually proceed to more complex situations.

Assumption 1: The regression model is *linear in the parameters* as in Equation (1.1); it may or may not be linear in the variables, the Y s and X s.

Assumption 2: The regressors are assumed fixed, or nonstochastic, in the sense that their values are fixed in repeated sampling. However, if the regressors are stochastic, we assume that each regressor is independent of the error term or at least uncorrelated with it. We will discuss this assumption in more detail in Chapter 6.

Assumption 3: Given the values of the X variables, the expected, or mean, value of the error term u_i is 0.

$$E(u_i | \mathbf{X}) = 0 \quad (2.1)$$

In matrix notation, we have

$$E(\mathbf{u} | \mathbf{X}) = \mathbf{0} \quad (2.1a)$$

where $\mathbf{0}$ is the **null vector**.

More explicitly,

$$E \begin{bmatrix} u_1 \\ u_2 \\ u_3 \\ \vdots \\ u_n \end{bmatrix} = \begin{bmatrix} E(u_1) \\ E(u_2) \\ E(u_3) \\ \vdots \\ E(u_n) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

Because of this *critical* assumption, and given the values of the regressors, we can write Equation (1.5) as

$$\begin{aligned} E(\mathbf{y} | \mathbf{X}) &= \mathbf{B}\mathbf{X} + E(\mathbf{u} | \mathbf{X}) \\ &= \mathbf{B}\mathbf{X} \end{aligned} \quad (2.2)$$

This is the PRF. In regression analysis, our primary objective is to estimate this function. The PRF thus gives the mean value of the regressand corresponding to the given values of the regressors, noting that conditional on these values the mean value of the error term is 0.

Assumption 4: The variance of each u_i , given the values of \mathbf{X} , is constant or **homoscedastic** (i.e., of equal variance). That is,

$$\text{var}(u_i | \mathbf{X}) = \sigma^2 \quad (2.3)$$

In matrix notation,

$$\text{var}(\mathbf{u} | \mathbf{X}) = E(\mathbf{u}\mathbf{u}') = \sigma^2 \mathbf{I} \quad (2.3a)$$

where \mathbf{I} is an $n \times n$ identity matrix (see also Assumption 5).

If $\text{var}(u_i | \mathbf{X}) = \sigma_i^2$, the error variance is said to be heteroscedastic, or of unequal variance. We will discuss this case in Chapter 5.

Figure 2.1 is a picture of both homoscedasticity and heteroscedasticity.

Assumption 5: There is no correlation between error terms belonging to two different observations. That is,

$$\text{cov}(u_i, u_j | \mathbf{X}) = 0, \quad i \neq j \quad (2.4)$$

where cov stands for covariance, and i and j are two different error terms. Of course, if $i=j$, we get the variance of u_i given in Equation (2.3).

Figure 2.2 shows a likely pattern of autocorrelation.

Figure 2.1 Homoscedasticity and Heteroscedasticity

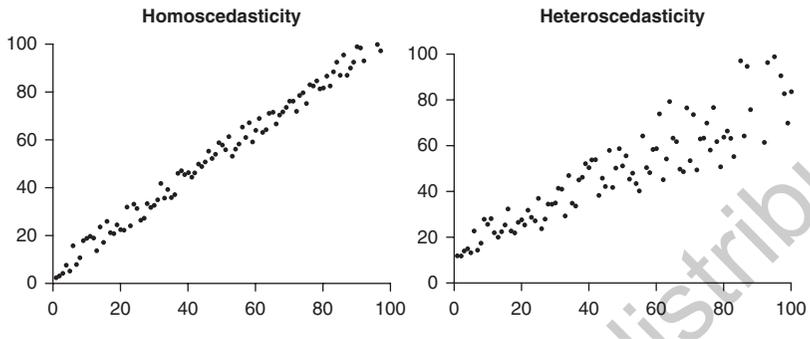
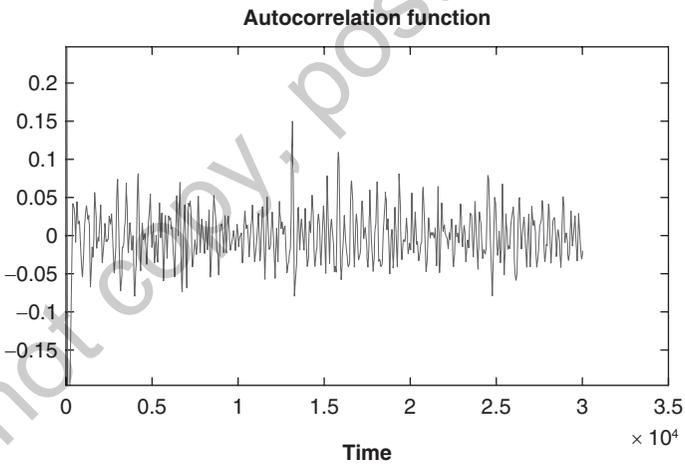


Figure 2.2 Autocorrelation



Assumptions 4 and 5 can be expressed as

$$\begin{aligned}
 E(\mathbf{u}\mathbf{u}') &= \sigma^2 \mathbf{I} & \text{if } i = j \\
 &= \mathbf{0} & \text{if } i \neq j
 \end{aligned}$$

where $\mathbf{0}$ is the **null matrix** and \mathbf{I} is the identity matrix. We discuss this assumption further in Chapter 5. More compactly, we can express Assumptions 4 and 5 as

$$E(\mathbf{u}\mathbf{u}') = \sigma^2 \mathbf{I} = \begin{pmatrix} \sigma^2 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma^2 \end{pmatrix}$$

Assumption 6: There is *no perfect linear relationship* among the X variables. This is the assumption of no **multicollinearity**. Strictly speaking, *multicollinearity* refers to the existence of more than one exact linear relationship, and *collinearity* refers to the existence of a single exact linear relationship. But this distinction is rarely maintained in practice, and multicollinearity refers to both cases. Imagine what would happen in the wage regression given in Equation (1.5), if we were to include work experience both in years and in months!

In matrix notation, this assumption means that the X matrix is of *full column rank*. In other words, *the columns of the X matrix are linearly independent*. This requires that the number of observations, n , is greater than the number of parameters estimated (i.e., the k regression coefficients). We discuss this assumption further in Chapter 7.

Assumption 7: The regression model used in the analysis is **correctly specified**, that is, there is no (model) **specific error** or **bias**. In practice, this is a tall assumption, but in Chapter 7, we discuss fully the import of this assumption.

Assumption 8: Although not a part of the original CLRM, for statistical inference (hypothesis testing), we assume that the error term u_i follows the normal distribution with 0 mean and (constant) variance σ^2 . Symbolically,

$$u_i \sim N(0, \sigma^2) \quad (2.5)$$

Or in matrix notation,

$$\mathbf{u} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}) \quad (2.5a)$$

The assumption of the normality of the error term is crucial if the sample size is rather small; it is not essential if we have a very large sample. However, we will revisit this assumption in Chapter 7. With this assumption, CLRM is known as the **classical normal linear regression model (CNLRM)**.

Since we are assuming that the X matrix is nonstochastic but \mathbf{u} is stochastic, the regressand Y is also stochastic. In addition, since \mathbf{u} is normally distributed with 0 mean and constant variance, Y inherits the properties of \mathbf{u} . More specifically,

$$\mathbf{y} \sim N(\mathbf{B}\mathbf{X}, \sigma^2 \mathbf{I}) \quad (2.6)$$

That is, the regressand is distributed normally with mean $\mathbf{B}\mathbf{X}$ and the (constant) variance σ^2 .

Under Assumption 8, we can use the method of **maximum likelihood (ML)** as an alternative to OLS. We will discuss ML more thoroughly in Chapter 3 because of its general applicability in many areas of statistics.

With one or more of the preceding assumptions, in this chapter, we discuss the following topics:

1. The sampling distribution of the OLS estimators, \mathbf{b}
2. An estimator of the unknown variance, σ^2
3. The relationship between the residual \mathbf{e} and the error \mathbf{u}
4. Small-sample properties of the OLS estimators
5. Large-sample properties of the OLS estimators

2.2 The Sampling or Probability Distributions of the OLS Estimators

Remember that the population parameters in \mathbf{B} , although unknown, are constants. However, this is not true of the estimated \mathbf{b} coefficients, for their values depend on the sample data at hand. In other words, the \mathbf{b} coefficients are random. As such, we would like to find their sampling or probability distributions to establish properties of the (OLS) estimators.

Recall that

$$\mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}, \text{ using Equation (1.16)}$$

Therefore,

$$\begin{aligned} \mathbf{b} &= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' [\mathbf{X}\mathbf{B} + \mathbf{u}], \text{ using Equation (1.5)} \\ &= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{X}\mathbf{B} + (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{u} \\ &= \mathbf{B} + (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{u} \end{aligned} \quad (2.7)$$

By the definition of covariance, we obtain

$$\begin{aligned}
 \text{cov}(\mathbf{b}) &= E(\mathbf{b} - \mathbf{B})(\mathbf{b} - \mathbf{B})' = E[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u}][(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u}]', \\
 &\quad \text{using Equation (2.7)} \\
 &= E[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u}\mathbf{u}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}] \\
 &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'E(\mathbf{u}\mathbf{u}')\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \\
 &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\sigma^2\mathbf{I}\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \\
 &= \sigma^2(\mathbf{X}'\mathbf{X})^{-1} \tag{2.8}
 \end{aligned}$$

In deriving this expression, we have used properties of the transpose of an inverse matrix, and the assumption that \mathbf{X} is fixed and that the variance of u_i is constant and the u s are uncorrelated. Notice that we can move the expectations operator through \mathbf{X} because it is assumed fixed. The variances of the individual elements of \mathbf{b} are on the main diagonal (running from upper left to lower right), and the off-diagonal elements give the covariances between pairs of coefficients in \mathbf{b} .

Since

$$\mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} \tag{1.16}$$

and the \mathbf{X} matrix is fixed, \mathbf{b} is a linear function of \mathbf{y} . Using Assumption 8, we know that \mathbf{y} is normally distributed. It is a property of the normal distribution that any linear function of a normally distributed variable is also normally distributed. Therefore, \mathbf{b} is ipso facto normally distributed as follows:

$$\mathbf{b} \sim N(\mathbf{B}, \sigma^2(\mathbf{X}'\mathbf{X})^{-1}) \tag{2.9}$$

That is, \mathbf{b} is normally distributed with \mathbf{B} as its mean (see Equation 2.20) and the variance established in Equation (2.8). In other words, under the normality assumption, the **sampling distribution** of the OLS estimator is normal, as shown in Equation (2.9). This finding will aid us in testing hypotheses about any element of \mathbf{B} or any linear combination thereof. It may be noted that a sampling distribution is a probability distribution of an estimator or of any test statistic. In other words, it describes the variation in the values of a statistic over all possible samples, here the variation in \mathbf{b} over all possible samples.¹

¹Suppose we draw several independent samples and for each sample we compute a (test) statistic, such as the mean, and draw a frequency distribution of all these statistics. Roughly speaking, this frequency distribution is the sampling distribution of that statistic. In our case, under the assumed conditions, the probability or sampling distribution of any component of \mathbf{b} is normal as shown in Equation (2.9).

For any single element of \mathbf{b} , b_k , we can express Equation (2.9) as

$$b_k \sim N(B_k, \sigma^2 x^{kk}) \quad (2.9a)$$

where x^{kk} is the k th diagonal element of $(\mathbf{X}'\mathbf{X})^{-1}$. The square root of $\sigma^2 x^{kk}$ will give the standard error of b_k (see Figure 2.3).

However, before we can engage in hypothesis testing, we need to estimate the unknown σ^2 . Remember that σ^2 refers to the variance of the error term \mathbf{u} . Since we do not observe \mathbf{u} directly, we have to rely on the estimated residuals, \mathbf{e} , to learn about the true variance. Toward that end, we need to establish the relationship between \mathbf{u} and \mathbf{e} . Recall that

$$\mathbf{e} = \mathbf{y} - \hat{\mathbf{y}} \quad (2.10)$$

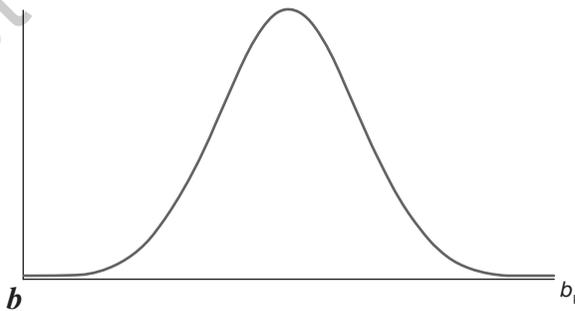
Substituting for $\hat{\mathbf{y}}$ from (1.23), we obtain

$$\begin{aligned} \mathbf{e} &= \mathbf{y} - \mathbf{X}\mathbf{b} \\ &= \mathbf{y} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}, \text{ substituting for } \mathbf{b} \text{ from Eq. (1.16)} \\ &= \mathbf{M}\mathbf{y} \\ &= \mathbf{M}(\mathbf{X}\mathbf{B} + \mathbf{u}) \end{aligned}$$

where

$$\begin{aligned} \mathbf{M} &= [\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'] \\ &= \mathbf{M}\mathbf{u}, \text{ because } \mathbf{M}\mathbf{X}\mathbf{B} = \mathbf{0}.^2 \end{aligned} \quad (2.11)$$

Figure 2.3 The Distribution of b_k , a Component of the Vector \mathbf{b}



² $\mathbf{M}\mathbf{X}\mathbf{B} = [\mathbf{X}\mathbf{B} - \mathbf{X}(\mathbf{X}\mathbf{X}')^{-1}\mathbf{X}'\mathbf{X}\mathbf{B}] = \mathbf{X}\mathbf{B} - \mathbf{I}\mathbf{X}\mathbf{B} = \mathbf{0}$, where \mathbf{I} is the identity matrix.

As noted in Chapter 1, M is a very important matrix in the analysis of LRMs. It is an **idempotent matrix**, a square matrix with the property that $M=M^2$. For further properties of the idempotent matrices, see Appendix A on linear algebra.

Since M is constant because it is a function of (fixed) X , we can write

$$\begin{aligned} E(\mathbf{e}) &= E(\mathbf{M}\mathbf{u}) \\ &= \mathbf{M}E(\mathbf{u}) \\ &= \mathbf{0} \end{aligned} \quad (2.12)$$

because $E(\mathbf{u})=\mathbf{0}$, by Assumption 1. We have thus shown that the expectation of each element of \mathbf{e} is 0.

Now,

$$\begin{aligned} \text{cov}(\mathbf{e}) &= E(\mathbf{e}\mathbf{e}') = E(\mathbf{M}\mathbf{u}\mathbf{u}'\mathbf{M}') \\ &= \mathbf{M}E(\mathbf{u}\mathbf{u}')\mathbf{M}' \\ &= \sigma^2 \mathbf{I}\mathbf{M} = \sigma^2 \mathbf{M} \end{aligned} \quad (2.13)$$

recalling the properties of M . This equation gives the covariance matrix of \mathbf{e} .

Since \mathbf{e} is a linear function of \mathbf{u} and since \mathbf{u} is normally distributed by Assumption 8, we have

$$\mathbf{e} \sim N(\mathbf{0}, \sigma^2 \mathbf{M}) \quad (2.14)$$

Therefore, like the mean of \mathbf{u} , the mean of \mathbf{e} is 0, but unlike \mathbf{u} , the residuals are heteroscedastic as well as autocorrelated.³

What Equations (2.13) and (2.14) show is that the residuals e_1, e_2, \dots, e_n have zero mean values, generally have different variances, and have nonzero covariances. Remember that in the (population) CLRM errors, u_1, u_2, \dots, u_n have zero expectations, have constant variance, and are not autocorrelated (by assumption). In other words, the properties that hold for \mathbf{u} generally do not hold for \mathbf{e} , except for zero expectations.

³Actually, the distribution of \mathbf{e} is degenerate as its variance–covariance matrix is singular. On this, see Vogelvang, B. (2005). *Econometrics: Theory and applications with Eviews* (chapter 4). Harlow, England: Pearson-Addison Wesley.

Although we have assumed that the variance of \mathbf{u} (not of \mathbf{e}) is constant, equal to σ^2 , we are yet to estimate it from the sample data. Toward that end, we proceed as follows.

Even though we do not observe \mathbf{u} , we observe \mathbf{e} (after the regression is estimated). Naturally, we will have to estimate the unknown variance from the estimated \mathbf{e} . From Equation (2.11), we know that

$$\mathbf{e} = \mathbf{M}\mathbf{u} \quad (2.11)$$

Therefore,

$$\begin{aligned} E(\mathbf{e}'\mathbf{e}) &= E(\mathbf{u}'\mathbf{M}'\mathbf{M}\mathbf{u}) \\ &= E(\mathbf{u}'\mathbf{M}\mathbf{u}) \end{aligned} \quad (2.15)$$

because of the properties of \mathbf{M} . Now,

$$\begin{aligned} E(\mathbf{e}'\mathbf{e}) &= E(\mathbf{u}'\mathbf{M}\mathbf{u}) \\ &= E[\text{tr}(\mathbf{u}'\mathbf{M}\mathbf{u})], \text{ since } \mathbf{u}'\mathbf{M}\mathbf{u} \text{ is a scalar} \\ &= E[\text{tr}(\mathbf{M}\mathbf{u}\mathbf{u}')] \text{ changing the order of multiplication inside} \\ &\quad \text{the trace} \\ &= \text{tr}[E(\mathbf{M}\mathbf{u}\mathbf{u}')], \text{ since the trace and expectations operators are} \\ &\quad \text{both linear} \\ &= \text{tr}[\mathbf{M}(\sigma^2\mathbf{I})] \\ &= \sigma^2 \text{tr}(\mathbf{M}) \\ &= \sigma^2 [\text{tr}(\mathbf{I}) - \text{tr}(\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')], \text{ using the definition of } \mathbf{M} \\ &= \sigma^2 [\text{tr}(\mathbf{I}_n) - \text{tr}((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X})] \\ &= \sigma^2 [\text{tr}(\mathbf{I}_n) - \text{tr}(\mathbf{I}_k)], \text{ since } (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X} = \mathbf{I}_k \\ &= \sigma^2 (n - k), \text{ since } \text{tr}(\mathbf{I}) = n \text{ and } \text{tr}(\mathbf{I}_k) = k \end{aligned} \quad (2.16)$$

The notation $\text{tr}(\mathbf{M})$ means the trace of the matrix \mathbf{M} , which is simply the sum of the entries of the main diagonal of \mathbf{M} . In deriving the steps in Equation (2.16), we have made use of several properties of the *trace* of a matrix, such as the fact that trace is a linear operator and if \mathbf{AB} and \mathbf{BA} are both square matrices, then $\text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA})$.

As a result, we can now write

$$E\left(\frac{\mathbf{e}'\mathbf{e}}{n-k}\right)=\sigma^2 \quad (2.17)$$

If we now define

$$S^2 = \frac{\mathbf{e}'\mathbf{e}}{n-k} = \frac{\sum e_i^2}{n-k} \quad (2.18)$$

then

$$E(S^2)=\sigma^2 \quad (2.19)$$

In words, S^2 is an unbiased estimator of the true error variance σ^2 . S , the square root of S^2 , is called the **standard error (se) of the estimate** or the **standard error of the regression**. In practice, therefore, we use S^2 in place of σ^2 .

2.3 Properties of OLS Estimators: The Gauss–Markov Theorem⁴

The OLS estimators possess some ideal or optimum properties, which are contained in the well-known **Gauss–Markov theorem**:⁵ Given the assumptions of the classical regression model, in the class of unbiased linear estimators, the least-squares estimators have minimum variance; that is, they are **best linear unbiased estimators, BLUE** for short. In other words, no other linear, unbiased estimator of \mathbf{B} can have a smaller variance than the OLS estimator given in Equation (2.8).

To establish this theorem, first note that \mathbf{b} , the OLS estimator of \mathbf{B} , is a linear function of the regressand \mathbf{y} , as we have established in Chapter 1 (see Equation 1.16).⁶ To prove that \mathbf{b} is unbiased, we proceed as follows:

$$\begin{aligned} \mathbf{b} &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} \quad (1.16) \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'[\mathbf{X}\mathbf{B} + \mathbf{u}], \text{ substituting for } \mathbf{y} \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\mathbf{B} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u} \\ &= \mathbf{B} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u} \end{aligned}$$

⁴In Appendix C, we discuss both small-sample and large-sample properties of OLS and ML estimators.

⁵Although known as the *Gauss–Markov theorem*, the least-squares approach of Gauss antedates (1821) the minimum-variance approach of Markov (1900).

⁶See the discussion in Darnell, A. C. (1994). *A dictionary of econometrics* (p. 155). Cheltenham, England: Edward Elgar.

Now

$$\begin{aligned} E(\mathbf{b}) &= \mathbf{B} + (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'E(\mathbf{u}) \\ &= \mathbf{B} \end{aligned} \quad (2.20)$$

In words, the expected value of \mathbf{b} is equal to \mathbf{B} , thus proving that \mathbf{b} is unbiased. (Recall the definition of unbiased estimator.) Note that $E(\mathbf{u}|\mathbf{X})=\mathbf{0}$ by assumption.

To prove that in the class of unbiased linear estimators the least-squares estimators have the least variance (i.e., they are efficient), we proceed as follows:

Let \mathbf{b}^* be another linear estimator of \mathbf{B} such that

$$\mathbf{b}^* = [\mathbf{A} + (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}']\mathbf{y} \quad (2.21)$$

where \mathbf{A} is some nonstochastic $k \times n$ matrix, similar to \mathbf{X} . Simplifying, we obtain

$$\begin{aligned} \mathbf{b}^* &= \mathbf{A}\mathbf{y} + (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y} \\ &= \mathbf{A}\mathbf{y} + \mathbf{b} \end{aligned} \quad (2.22)$$

where \mathbf{b} is the least-squares estimator given in Equation (1.16).

Now

$$\begin{aligned} E(\mathbf{b}^*) &= [\mathbf{A} + (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}']E(\mathbf{y}) \\ &= [\mathbf{A} + (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'](\mathbf{X}\mathbf{B}) \\ &= (\mathbf{A}\mathbf{X} + \mathbf{I})\mathbf{B} \end{aligned} \quad (2.23)$$

Now $E(\mathbf{b}^*)=\mathbf{B}$ if and only if $\mathbf{A}\mathbf{X}=\mathbf{0}$. In other words, for the linear estimator \mathbf{b}^* to be unbiased, $\mathbf{A}\mathbf{X}$ must be $\mathbf{0}$.

Thus,

$$\begin{aligned} \mathbf{b}^* &= [\mathbf{A} + (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'][\mathbf{X}\mathbf{B} + \mathbf{u}], \text{ substituting for } (\mathbf{y}) \\ &= \mathbf{B} + [\mathbf{A} + (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}']\mathbf{u}, \text{ because } \mathbf{A}\mathbf{X} = \mathbf{0} \end{aligned}$$

Given that \mathbf{u} has zero mean and constant variance ($=\sigma^2\mathbf{I}$), we can now find the variance of \mathbf{b}^* as follows:

$$\begin{aligned} \text{cov}(\mathbf{b}^*) &= E[\mathbf{A} + (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}']\mathbf{u}\mathbf{u}'[\mathbf{A} + (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}']' \\ &= [\mathbf{A} + (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}']E(\mathbf{u}\mathbf{u}')[\mathbf{A} + (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}']' \\ &= \sigma^2[\mathbf{A}\mathbf{A}' + (\mathbf{X}'\mathbf{X})^{-1}] \\ &= \sigma^2(\mathbf{X}'\mathbf{X})^{-1} + \mathbf{A}\mathbf{A}'\sigma^2 \\ &= \text{var}(\mathbf{b}) + \mathbf{A}\mathbf{A}'\sigma^2 \end{aligned} \quad (2.24)$$

Since AA' is a positive semidefinite matrix, Equation (2.24) shows that the covariance matrix of \mathbf{b}^* is equal to the covariance matrix of \mathbf{b} plus a positive semidefinite matrix. That is, $\text{cov}(\mathbf{b}^*) > \text{cov}(\mathbf{b})$, unless $\mathbf{A} = \mathbf{0}$. This shows that in the class of unbiased linear estimators, the least-square estimator \mathbf{b} has the least variance, that is, it is efficient compared with any other linear unbiased estimator of \mathbf{B} .

It is important to note that in establishing the Gauss–Markov theorem we do not have to assume that the error term \mathbf{u} follows a particular probability distribution, such as the normal. To establish the theorem, we only need Assumptions 1 to 5.

It is also important to note that if one or more assumptions underlying the Gauss–Markov theorem are not satisfied, the OLS estimators will not be BLUE. Also, bear in mind that the Gauss–Markov theorem holds only for linear estimators, that is, linear functions of the observation vector \mathbf{y} . There are situations where nonlinear (in-the-parameter) estimators are more efficient than the linear estimators. In this book, we do not deal with nonlinear estimators, for that requires a separate book.⁷

To sum up, we have shown that under the Gauss–Markov assumptions, \mathbf{b} , the least-square estimator of \mathbf{B} , is BLUE, that is, in the class of unbiased linear estimators, \mathbf{b} has the least variance. We also showed how to estimate \mathbf{B} and the variance of the estimated \mathbf{B} .

2.4 Estimating Linear Functions of the OLS Parameters

We have shown how to estimate \mathbf{B} , that is, each of its elements. Suppose we want to estimate some linear function of the elements of \mathbf{B} , that is, that of $B_1, B_2, B_3, \dots, B_k$. More specifically, suppose we want to estimate $\mathbf{t}'\mathbf{B}$, where \mathbf{t}' is a $1 \times k$ vector of real numbers and \mathbf{B} is a $k \times 1$ vector of the parameters in \mathbf{B} . It can be shown that the BLUE of $\mathbf{t}'\mathbf{B}$ is $\mathbf{t}'\mathbf{b}$, where \mathbf{b} is the least-square estimator of \mathbf{B} (see also Appendix C).

What this means is that whether we estimate all the elements of \mathbf{B} , or one of its elements, or estimate a linear combination ($\mathbf{t}'\mathbf{B}$), we can use the OLS regression.

Let $\lambda = \mathbf{t}'\mathbf{B}$. By choosing \mathbf{t} appropriately, we can make λ equal to any element of \mathbf{B} , or to the sum of the elements of \mathbf{B} that might be of interest to researchers. Suppose in Equation (1.2), we want the coefficient of \mathbf{B}_1 equal to 4, and the coefficient of \mathbf{B}_5 equal to -1 , and the rest of the coefficients to be all zeros. In other words, we want $\lambda = 4\mathbf{B}_1 - \mathbf{B}_5$. Here, $\mathbf{t}' = (4, 0, 0, 0, -1, 0, 0, 0, \dots)'$.

Using the definition of variance, we can now find the variance of the estimated $\lambda (= \hat{\lambda})$, which is

⁷For examples of nonlinear estimators and their applications, see Gujarati, D. (2015). *Econometrics by example* (2nd ed.). London, England: Palgrave Macmillan.

$$\text{var}(\hat{\lambda}) = \mathbf{t}'(\text{var}(\mathbf{b}))\mathbf{t} = \sigma^2 \mathbf{t}'(\mathbf{X}'\mathbf{X})^{-1} \mathbf{t} \quad (2.25)$$

But keep in mind that in general the variance of $\hat{\lambda}$ depends on every element of the covariance matrix of \mathbf{b} , the estimator of \mathbf{B} . However, if some elements of the vector \mathbf{t} are equal to zero, $\text{var}(\hat{\lambda})$ does not depend on the corresponding rows and columns of the covariance matrix $\sigma^2(\mathbf{X}'\mathbf{X})^{-1}$.

As an example, consider $\lambda = 4B_1 - B_5$. In this case,

$$\begin{aligned} \text{var}(\hat{\lambda}) &= t_1^2 \text{var}(b_1) + t_5^2 \text{var}(b_5) + 2t_1 t_5 \text{cov}(b_1, b_5) \\ &= 16 \text{var}(b_1) + \text{var}(b_5) - 8 \text{cov}(b_1, b_5) \end{aligned} \quad (2.26)$$

Notice in this example only the variances of b_1 and b_5 and their covariances are involved, as the values of the other parameters in the k -variable regression (1.2) are assumed to be zero. But if there are more nonzero coefficients, the variances and their pairwise covariances will also be involved in computing the variance of the linear combination $\mathbf{t}'\mathbf{B}$.

2.5 Large-Sample Properties of OLS Estimators

2.5.1 Consistency of OLS Estimators

We have shown that the OLS estimators of the CLRM are unbiased, which is a small, or finite sample, property. We can also show that the OLS estimators are consistent, that is, they converge to their true values as the sample size increases indefinitely. Convergence is a large-sample property.

Proof: A sufficient condition for an unbiased estimator to be consistent is for its variance to converge to zero as the sample size n increases indefinitely. For the OLS estimator \mathbf{b} , we have already shown that its variance is

$$\text{cov}(\mathbf{b}) = \sigma^2 (\mathbf{X}'\mathbf{X})^{-1} \quad (2.8)$$

which we can write as

$$\text{cov}(\mathbf{b}) = \frac{\sigma^2}{n} (n^{-1} \mathbf{X}'\mathbf{X})^{-1} \quad (2.27)$$

To see the behavior of this expression as $n \rightarrow \infty$, we have

$$\begin{aligned} \text{plim}_{n \rightarrow \infty} \text{cov}(\mathbf{b}) &= \text{plim}_{n \rightarrow \infty} \left[\frac{\sigma^2}{n} (n^{-1} \mathbf{X}'\mathbf{X})^{-1} \right] \\ &= \text{plim}_{n \rightarrow \infty} \frac{\sigma^2}{n} \lim_{n \rightarrow \infty} (n^{-1} \mathbf{X}'\mathbf{X})^{-1} \end{aligned} \quad (2.28)$$

where *plim* is probability limit (see Appendix C for details).

We have assumed that the elements of the matrix X are bounded, which means the second term in the preceding equation is bounded for all n . Therefore, the second term above can be replaced by a matrix of finite constants. Now, the limit of the first term in Equation (2.28) tends to zero as n increases indefinitely. As a result,

$$\underset{n \rightarrow \infty}{\text{plim}} \text{cov}(\mathbf{b}) = 0 \quad (2.29)$$

which establishes that \mathbf{b} is a consistent estimator of \mathbf{B} . In establishing this result, we have used some of the properties of the *plim* operator.

2.5.2 Consistency of the OLS Estimator of the Error Variance

We have proved that S^2 is an unbiased estimator of σ^2 . Assuming values of u_i are independent and identically distributed (iid), we can prove that S^2 is also a consistent estimator of σ^2 . The proof is as follows:

$$\begin{aligned} S^2 &= \frac{(\mathbf{y} - \mathbf{Xb})'(\mathbf{y} - \mathbf{Xb})}{n - k} \\ &= \frac{\mathbf{u}'(\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')\mathbf{u}}{n - k} \\ &= \left(\frac{n}{n - k}\right) \left(\frac{\mathbf{u}'\mathbf{u}}{n} - \frac{\mathbf{u}'\mathbf{X}}{n} \cdot \left(\frac{\mathbf{X}'\mathbf{X}}{n}\right)^{-1} \cdot \frac{\mathbf{X}'\mathbf{u}}{n}\right) \end{aligned} \quad (2.30)$$

Note: $\mathbf{e} = \mathbf{My} = \mathbf{Mu}$, where $\mathbf{M} = [\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']$. Also note how the entries are manipulated by multiplying or dividing them by the sample size or the adjusted sample size without affecting the basic relationships.

Taking the *plim* of both sides of Equation (2.30), we obtain

$$\begin{aligned} \underset{n \rightarrow \infty}{\text{plim}} S^2 &= \underset{n \rightarrow \infty}{\text{plim}} \left(\frac{n}{n - k}\right) \left(\underset{n \rightarrow \infty}{\text{plim}} \frac{\mathbf{u}'\mathbf{u}}{n} - \underset{n \rightarrow \infty}{\text{plim}} \frac{\mathbf{u}'\mathbf{X}}{n} \cdot \underset{n \rightarrow \infty}{\text{plim}} \left(\frac{\mathbf{X}'\mathbf{X}}{n}\right)^{-1} \cdot \underset{n \rightarrow \infty}{\text{plim}} \frac{\mathbf{X}'\mathbf{u}}{n}\right) \\ &= 1(\sigma^2 - 0 \cdot \mathcal{Q}^{-1} \cdot 0), \text{ where } \underset{n \rightarrow \infty}{\text{plim}} \left(\frac{\mathbf{X}'\mathbf{X}}{n}\right)^{-1} = \mathcal{Q}^{-1} \\ &= \sigma^2 \end{aligned} \quad (2.31)$$

which establishes the result. Note that for large n , $(n - k) \approx n$.

In deriving the preceding result, we have used **Khinchine's theorem** (see Appendix B) as well as the properties of the *plim*.

2.5.3 Independence of the OLS

Estimators and the Residual Term, e

What this says is that each element of \mathbf{b} is uncorrelated with each element of the least-squares residual vector \mathbf{e} . The proof is as follows:

Recall that

$$\mathbf{b} = \mathbf{B} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u} \quad (2.7)$$

$$\mathbf{e} = \mathbf{M}\mathbf{u} \quad (2.11)$$

are both linear functions of the error term \mathbf{u} .

Now the covariance of \mathbf{b} and \mathbf{e} is

$$\begin{aligned} \text{cov}(\mathbf{b}, \mathbf{e}) &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\text{var}(\mathbf{u})\mathbf{M}' \\ &= \sigma^2(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{M}' \\ &= 0, \quad \text{since } \mathbf{M}\mathbf{X} = \mathbf{0} \leftrightarrow \mathbf{X}'\mathbf{M}' \end{aligned} \quad (2.32)$$

This shows that \mathbf{b} and \mathbf{e} are uncorrelated.

It may be noted that if we assume that \mathbf{u} is normally distributed, \mathbf{b} and \mathbf{e} are not only uncorrelated but also independent. In Chapter 3, we will consider the **normal linear regression model**, which explicitly assumes that the error term \mathbf{u} is normally distributed and will see the consequence of the normality assumption.

2.5.4 Large-Sample Distribution of \mathbf{b} :

Asymptotic Normality of the OLS Estimators

It can be shown that⁸

$$\mathbf{b} \text{ asy} \sim N(\mathbf{B}, \sigma^2(\mathbf{X}'\mathbf{X})^{-1}) \quad (2.33)$$

where asy means asymptotically (i.e., $n \rightarrow \infty$).

In other words, as the sample size n increases indefinitely, \mathbf{b} is approximately normally distributed with mean equal to \mathbf{B} and variance equal to

⁸The proof is rather complicated and can be found in Theil, H. (1971). *Principles of econometrics* (pp. 380–381). New York, NY: Wiley; see also Mittlehammer, R. C. (1996). *Mathematical statistics for economics and business* (pp. 443–447). New York, NY: Springer.

$\sigma^2(\mathbf{X}'\mathbf{X})^{-1}$. Each element of \mathbf{b} is individually normally distributed with variance equal to the appropriate element of the variance matrix $\sigma^2(\mathbf{X}'\mathbf{X})^{-1}$. *This result holds whether \mathbf{u} is normally distributed or not.*

It may be noted that if the errors u_i are not iid, even then \mathbf{b} is asymptotically normally distributed as in (2.33) under certain conditions.⁹

2.5.5 Asymptotic Normality of \mathbf{S}^2

If in addition to the classical assumptions, it is assumed that values of u_i are iid and have bounded fourth-order moments about the origin, \mathbf{S}^2 is asymptotically normally distributed. These results also hold true even if the u_i values are not iid.¹⁰

2.6 Summary

The CLRM, $\mathbf{y}=\mathbf{X}\mathbf{B}+\mathbf{u}$, is the foundation of regression analysis. It is based on several assumptions. The basic assumptions are that (1) the data matrix \mathbf{X} is nonstochastic, (2) it is of full column rank, (3) the expected value of the error term is zero, and (4) the covariance matrix of the error term $E(\mathbf{u}\mathbf{u}')=\sigma^2\mathbf{I}$. This means the error variance is constant and equal to σ^2 and that the error terms are mutually uncorrelated.

We used the method of OLS to estimate the parameters of an LRM. One reason for using OLS is that it does not require us to make assumptions about the probability distribution of the error term, and it is comparatively easy to estimate. Parameters of the CLRM estimated by OLS are called OLS estimators. OLS estimators have several desirable statistical properties such as (1) they are unbiased and (2) among all linear unbiased estimators of \mathbf{B} , they have minimum variances. This is called the Gauss–Markov theorem. These are small-sample properties.

OLS estimators have these asymptotic, or large-sample, properties: (1) The OLS estimators of \mathbf{B} as well as the estimator of the error variance are consistent estimators and (2) the OLS estimators asymptotically follow the normal distribution.

Exercises

2.1 Consider the bivariate regression: $Y_i = B_1 + B_2X_i + u_i$. Under the classical linear regression assumptions, show that

⁹See Mittlehammer, R. C. (1996). *Mathematical statistics for economics and business* (p. 445). New York, NY: Springer.

¹⁰See Mittlehammer, R. C. (1996). *Mathematical statistics for economics and business* (pp. 448–449). New York, NY: Springer.

$$\text{a. } \text{cov}(b_1, b_2) = -\bar{X} \frac{\sigma^2}{\Sigma(X_i - \bar{X})^2}$$

$$\text{b. } \text{cov}(\bar{Y}, b_2) = 0$$

2.2 Show that for the model in Exercise 2.1,

$$\text{RSS} = \frac{\Sigma x_i^2 \Sigma y_i^2 - (\Sigma x_i y_i)^2}{\Sigma x_i^2}$$

where RSS is the residual sum of squares and

$$x_i = (X_i - \bar{X}); \quad y_i = (Y_i - \bar{Y}); \quad x_i y_i = (X_i - \bar{X})(Y_i - \bar{Y})$$

2.3 Verify the following properties of OLS estimators:

- The OLS regression line (plane) passes through the sample means of the regressand and the regressors.
- The mean values of the actual Y and the estimated $Y (= \hat{Y})$ are the same.
- In the CLRM with intercept, the mean value of the residuals (\bar{e}) is zero.
- As a result of the preceding property, the k -variable sample CLRM can be expressed as

$$y_i = b_2 x_{2i} + b_3 x_{3i} + \dots + b_k x_{ki} + e_i$$

$$\text{where } y_i = (Y_i - \bar{Y}); \quad x_{ki} = (X_{ki} - \bar{X}_k)$$

2.4 Consider the following bivariate regression model:

$$Y_i^* = B_1^* + B_2^* X_i^* + u_i$$

where

$$Y_i^* = \frac{Y_i - \bar{Y}}{s_Y}; \quad X_i^* = \frac{X_i - \bar{X}}{s_X}$$

where s_Y and s_X are the sample standard deviations of Y and X . Y_i^* and X_i^* are known as **standardized variables**, often known as **Z scores**. Since the units of measurement of the Z scores in the numerator and the denominator are the same, they are called “pure” or “unitless” numbers.

- Show that a standardized variable has a zero mean and unit variance.
 - What are the formulas to estimate B_1^* and B_2^* ?
 - What is the relationship between B_1^* and B_1 and between B_2^* and B_2 ?
- 2.5 The sample correlation coefficient between variables Y and X , r_{XY} , is defined as

$$r_{XY} = \frac{\sum x_i y_i}{\sqrt{\sum x_i^2 \sum y_i^2}}$$

where

$$x_i = (X_i - \bar{X}); \quad y_i = (Y_i - \bar{Y})$$

If we standardize variables as in Exercise 2.4, does it affect the correlation coefficient between X and Y ? Show the necessary calculations.

- 2.6 Consider variables X_1 , X_2 , and X_3 . Now consider the following correlation coefficients:

$$\begin{aligned} r_{12} &= \text{correlation coefficient between } X_1 \text{ and } X_2 \\ r_{13} &= \text{correlation coefficient between } X_1 \text{ and } X_3 \\ r_{23} &= \text{correlation coefficient between } X_2 \text{ and } X_3 \end{aligned}$$

$$r_{12.3} = \frac{r_{12} - r_{13}r_{23}}{\sqrt{1 - r_{13}^2} \sqrt{1 - r_{23}^2}}$$

$r_{12.3}$ is called the **partial correlation coefficient** between X_1 and X_2 holding the influence of the variable X_3 . The concept of partial correlation is akin to the concept of a partial regression coefficient.

- What other partial correlation coefficients can you compute?
 - If we standardize the three variables as in Exercise 2.4, would the correlation coefficients among the standardized variables be different from the unstandardized variables?
 - Would partial correlation coefficients be affected by standardizing the variables? Explain.
- 2.7 Consider the following LRM:

$$Y_i = B_1 + B_2 X_{2i} + B_3 X_{3i} + B_4 X_{4i} + B_5 X_{5i} + u_i$$

How would you test the following hypotheses?

- $B_2 = B_3 = B_4 = B_5 = B$, that is, all partial regression coefficients are the same.
- $B_2 = B_3$ and $B_4 = B_5$
- $B_2 + B_3 = 2B_4$

2.8 Remember that the hat matrix, H , is expressed as

$$H = X(X'X)^{-1}X'$$

Show that the residual vector e can also be expressed as

$$e = (I - H)y$$

2.9 Prove that the matrices H and $(I - H)$ are idempotent.

2.10* For the following matrix, compute its eigenvalues:

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

(*Optional)

2.11 Consider the following regression model (see Chapter 7, Equation (7.30)):

$$Y_i = B_1 + B_2X_i + B_3X_i^2 + u_i$$

Models like this are called polynomial regression models, here a second-degree polynomial.

- Is this an LRM?
- Can OLS be used to estimate the parameters of this model?
- Since X_i^2 is the square of X_i , does this model suffer from perfect collinearity?

2.12 Consider the following model:

$$Y_i = B_1 + B_2X_{2i} + B_3X_{3i} + B_4X_{4i} + u_i$$

You are told that $B_2 = 1$.

- a. In this case, is it legitimate to estimate the following regression?

$$(Y_i - X_{2i}) = B_1 + B_3X_{3i} + B_4X_{4i} + u_i$$

This model is called a **restricted** linear regression, whereas the preceding model is called an **unrestricted** linear regression (see Chapter 4, Appendix 4A for further details).

- b. How would you estimate the restricted regression, taking into account the restriction that $B_2=1$?

Do not copy, post, or distribute